

Fixed-Point Designer™

User's Guide



MATLAB®

R2023a



How to Contact MathWorks



Latest news: www.mathworks.com
Sales and services: www.mathworks.com/sales_and_services
User community: www.mathworks.com/matlabcentral
Technical support: www.mathworks.com/support/contact_us



Phone: 508-647-7000



The MathWorks, Inc.
1 Apple Hill Drive
Natick, MA 01760-2098

Fixed-Point Designer™ User's Guide

© COPYRIGHT 2013–2023 by The MathWorks, Inc.

The software described in this document is furnished under a license agreement. The software may be used or copied only under the terms of the license agreement. No part of this manual may be photocopied or reproduced in any form without prior written consent from The MathWorks, Inc.

FEDERAL ACQUISITION: This provision applies to all acquisitions of the Program and Documentation by, for, or through the federal government of the United States. By accepting delivery of the Program or Documentation, the government hereby agrees that this software or documentation qualifies as commercial computer software or commercial computer software documentation as such terms are used or defined in FAR 12.212, DFARS Part 227.72, and DFARS 252.227-7014. Accordingly, the terms and conditions of this Agreement and only those rights specified in this Agreement, shall pertain to and govern the use, modification, reproduction, release, performance, display, and disclosure of the Program and Documentation by the federal government (or other entity acquiring for or through the federal government) and shall supersede any conflicting contractual terms or conditions. If this License fails to meet the government's needs or is inconsistent in any respect with federal procurement law, the government agrees to return the Program and Documentation, unused, to The MathWorks, Inc.

Trademarks

MATLAB and Simulink are registered trademarks of The MathWorks, Inc. See www.mathworks.com/trademarks for a list of additional trademarks. Other product or brand names may be trademarks or registered trademarks of their respective holders.

Patents

MathWorks products are protected by one or more U.S. patents. Please see www.mathworks.com/patents for more information.

Revision History

| | | |
|----------------|-------------|---|
| March 2013 | Online only | New for Version 4.0 (R2013a) |
| September 2013 | Online only | Revised for Version 4.1 (R2013b) |
| March 2014 | Online only | Revised for Version 4.2 (R2014a) |
| October 2014 | Online Only | Revised for Version 4.3 (R2014b) |
| March 2015 | Online Only | Revised for Version 5.0 (R2015a) |
| September 2015 | Online Only | Revised for Version 5.1 (R2015b) |
| October 2015 | Online only | Rereleased for Version 5.0.1 (Release 2015aSP1) |
| March 2016 | Online Only | Revised for Version 5.2 (R2016a) |
| September 2016 | Online only | Revised for Version 5.3 (R2016b) |
| March 2017 | Online only | Revised for Version 5.4 (R2017a) |
| September 2017 | Online only | Revised for Version 6.0 (R2017b) |
| March 2018 | Online only | Revised for Version 6.1 (R2018a) |
| September 2018 | Online only | Revised for Version 6.2 (R2018b) |
| March 2019 | Online only | Revised for Version 6.3 (R2019a) |
| September 2019 | Online only | Revised for Version 6.4 (R2019b) |
| March 2020 | Online only | Revised for Version 7.0 (R2020a) |
| September 2020 | Online only | Revised for Version 7.1 (R2020b) |
| March 2021 | Online only | Revised for Version 7.2 (R2021a) |
| September 2021 | Online only | Revised for Version 7.3 (R2021b) |
| March 2022 | Online only | Revised for Version 7.4 (R2022a) |
| September 2022 | Online only | Revised for Version 7.5 (R2022b) |
| March 2023 | Online only | Revised for Version 7.6 (R2023a) |

Fixed-Point Designer for MATLAB Code

| | |
|----------|--|
| | Fixed-Point Concepts |
| 1 | |
| | Fixed-Point Data Types 1-2 |
| | Scaling 1-3 |
| | Compute Slope and Bias 1-4 |
| | What Is Slope Bias Scaling? 1-4 |
| | Compute Slope and Bias 1-4 |
| | Precision and Range 1-6 |
| | Range 1-6 |
| | Precision 1-7 |
| | Arithmetic Operations 1-9 |
| | Modulo Arithmetic 1-9 |
| | Two's Complement 1-9 |
| | Addition and Subtraction 1-10 |
| | Multiplication 1-11 |
| | Casts 1-15 |
| | fi Objects and C Integer Data Types 1-18 |
| | Integer Data Types 1-18 |
| | Unary Conversions 1-19 |
| | Binary Conversions 1-20 |
| | Overflow Handling 1-21 |
| | Working with fi Objects |
| 2 | |
| | Ways to Construct fi Objects 2-2 |
| | Use fi Constructor to Create fi Objects 2-2 |
| | Use the Insert fi Constructor Dialog to Build fi Object Constructors |
| | 2-6 |
| | Determine Property Precedence 2-7 |
| | Create fi Objects For Use in a Types Table 2-8 |

| | |
|--|-------------|
| Cast fi Objects | 2-10 |
| Overwriting by Assignment | 2-10 |
| Ways to Cast with MATLAB Software | 2-10 |
| Set fi Object Properties | 2-15 |
| Set Fixed-Point Properties at Object Creation | 2-15 |
| Use Subscripted Assignment to Set Real-World Value of fi Object .. | 2-15 |
| Direct Property Referencing to Read fi Object Properties | 2-16 |
| Best Practices for Code Generation | 2-16 |
| Remove Local fimath Properties from fi Object | 2-18 |

Working with fimath Objects

3

| | |
|---|-------------|
| fimath Object Construction | 3-2 |
| fimath Object Syntaxes | 3-2 |
| Building fimath Object Constructors in a GUI | 3-3 |
| fimath Object Properties | 3-4 |
| Math, Rounding, and Overflow Properties | 3-4 |
| How Properties are Related | 3-7 |
| Setting fimath Object Properties | 3-8 |
| fimath Properties Usage for Fixed-Point Arithmetic | 3-10 |
| fimath Rules for Fixed-Point Arithmetic | 3-10 |
| Binary-Point Arithmetic | 3-11 |
| [Slope Bias] Arithmetic | 3-14 |
| fimath for Rounding and Overflow Modes | 3-16 |
| fimath for Sharing Arithmetic Rules | 3-17 |
| Default fimath Usage to Share Arithmetic Rules | 3-17 |
| Local fimath Usage to Share Arithmetic Rules | 3-17 |
| fimath ProductMode and SumMode | 3-19 |
| Example Setup | 3-19 |
| FullPrecision | 3-19 |
| KeepLSB | 3-20 |
| KeepMSB | 3-21 |
| SpecifyPrecision | 3-22 |
| How Functions Use fimath | 3-24 |
| Functions that use then discard attached fimath | 3-24 |
| Functions that ignore and discard attached fimath | 3-24 |
| Functions that do not perform math | 3-24 |

4

- Set fi Object Display Preferences Using fipref** 4-2
- Underflow and Overflow Logging Using fipref** 4-3
 - Logging Overflows and Underflows as Warnings 4-3
 - Accessing Logged Information with Functions 4-5
- Data Type Override Preferences Using fipref** 4-7
 - Overriding the Data Type of fi Objects 4-7
 - Data Type Override for Fixed-Point Scaling 4-8

Working with numerictype Objects

5

- numerictype Object Construction** 5-2
 - numerictype Object Syntaxes 5-2
 - Example: Construct a numerictype Object with Property Name and Property Value Pairs 5-2
 - Example: Copy a numerictype Object 5-3
 - Example: Build numerictype Object Constructors in a GUI 5-4
- numerictype Object Properties** 5-5
 - Data Type and Scaling Properties 5-5
 - How Properties are Related 5-7
 - Set numerictype Object Properties 5-8
- numerictype of Fixed-Point Objects** 5-9
 - Valid Values for numerictype Object Properties 5-9
 - Properties That Affect the Slope 5-10
 - Stored Integer Value and Real World Value 5-10
- numerictype Objects Usage to Share Data Type and Scaling Settings of fi objects** 5-12
 - Example 1 5-12
 - Example 2 5-12

Working with quantizer Objects

6

- Transformations for Quantized Data** 6-2

| | |
|--|------|
| Fixed-Point Conversion Workflows | 7-2 |
| Choosing a Conversion Workflow | 7-2 |
| Automated Workflow | 7-2 |
| Manual Workflow | 7-2 |
| | |
| Automated Fixed-Point Conversion | 7-4 |
| Automated Fixed-Point Conversion Capabilities | 7-4 |
| Code Coverage | 7-5 |
| Proposing Data Types | 7-7 |
| Locking Proposed Data Types | 7-10 |
| Viewing Functions | 7-10 |
| Viewing Variables | 7-17 |
| Log Data for Histogram | 7-19 |
| Function Replacements | 7-21 |
| Validating Types | 7-21 |
| Testing Numerics | 7-22 |
| Detecting Overflows | 7-22 |
| | |
| Debug Numerical Issues in Fixed-Point Conversion Using Variable Logging | 7-23 |
| Prerequisites | 7-23 |
| Convert to Fixed Point Using Default Configuration | 7-26 |
| Determine Where Numerical Issues Originated | 7-29 |
| Adjust fimath Settings | 7-30 |
| Adjust Word Length Settings | 7-31 |
| Replace Constant Functions | 7-32 |
| | |
| MATLAB Language Features Supported for Automated Fixed-Point Conversion | 7-35 |
| MATLAB Language Features Supported for Automated Fixed-Point Conversion | 7-35 |
| MATLAB Language Features Not Supported for Automated Fixed-Point Conversion | 7-36 |
| | |
| Generated Fixed-Point Code | 7-37 |
| Location of Generated Fixed-Point Files | 7-37 |
| Minimizing fi-casts to Improve Code Readability | 7-37 |
| Avoiding Overflows in the Generated Fixed-Point Code | 7-38 |
| Controlling Bit Growth | 7-38 |
| Avoiding Loss of Range or Precision | 7-39 |
| Handling Non-Constant mpower Exponents | 7-40 |
| | |
| Fixed-Point Code for MATLAB Classes | 7-42 |
| Automated Conversion Support for MATLAB Classes | 7-42 |
| Unsupported Constructs | 7-42 |
| Coding Style Best Practices | 7-42 |
| | |
| Automated Fixed-Point Conversion Best Practices | 7-44 |
| Create a Test File | 7-44 |
| Prepare Your Algorithm for Code Acceleration or Code Generation | 7-45 |

| | |
|--|-------------|
| Check for Fixed-Point Support for Functions Used in Your Algorithm | 7-45 |
| Manage Data Types and Control Bit Growth | 7-46 |
| Convert to Fixed Point | 7-46 |
| Use the Histogram to Fine-Tune Data Type Settings | 7-47 |
| Optimize Your Algorithm | 7-47 |
| Avoid Explicit Double and Single Casts | 7-49 |
| Replacing Functions Using Lookup Table Approximations | 7-50 |
| Custom Plot Functions | 7-51 |
| Generate Fixed-Point MATLAB Code for Multiple Entry-Point Functions | 7-52 |
| Convert Code Containing Global Data to Fixed Point | 7-56 |
| Workflow | 7-56 |
| Declare Global Variables | 7-56 |
| Define Global Data | 7-56 |
| Define Constant Global Data | 7-57 |
| Convert Code Containing Global Variables to Fixed-Point | 7-60 |
| Convert Code Containing Structures to Fixed Point | 7-64 |
| Convert Identical Functions Called with Different Data Types | 7-67 |
| Data Type Issues in Generated Code | 7-71 |
| Enable the Highlight Option in the Fixed-Point Converter App | 7-71 |
| Enable the Highlight Option at the Command Line | 7-71 |
| Stowaway Doubles | 7-71 |
| Stowaway Singles | 7-71 |
| Expensive Fixed-Point Operations | 7-71 |
| System Objects Supported by Fixed-Point Converter App | 7-73 |
| Convert dsp.FIRFilter Object to Fixed-Point Using the Fixed-Point Converter App | 7-74 |
| Create DSP Filter Function and Test Bench | 7-74 |
| Convert the Function to Fixed-Point | 7-75 |

Automated Conversion Using Fixed-Point Converter App

8

| | |
|--|-------------|
| Specify Type Proposal Options | 8-2 |
| Detect Overflows | 8-5 |
| Propose Data Types Based on Simulation Ranges | 8-13 |
| Propose Data Types Based on Derived Ranges | 8-24 |

| | |
|---|-------------|
| View and Modify Variable Information | 8-35 |
| View Variable Information | 8-35 |
| Modify Variable Information | 8-35 |
| Revert Changes | 8-36 |
| Promote Sim Min and Sim Max Values | 8-36 |
| Replace the exp Function with a Lookup Table | 8-38 |
| Convert Fixed-Point Conversion Project to MATLAB Scripts | 8-45 |
| Replace a Custom Function with a Lookup Table | 8-47 |
| Visualize Differences Between Floating-Point and Fixed-Point Results | 8-52 |
| Copy Relevant Files | 8-52 |
| Prerequisites | 8-52 |
| Enable Plotting Using the Simulation Data Inspector | 8-62 |
| Add Global Variables by Using the App | 8-63 |
| Automatically Define Input Types by Using the App | 8-64 |
| Define Constant Input Parameters Using the App | 8-65 |
| Define or Edit Input Parameter Type by Using the App | 8-66 |
| Define or Edit an Input Parameter Type | 8-66 |
| Specify a String Scalar Input Parameter | 8-67 |
| Specify an Enumerated Type Input Parameter | 8-67 |
| Specify a Fixed-Point Input Parameter | 8-68 |
| Specify a Structure Input Parameter | 8-68 |
| Specify a Cell Array Input Parameter | 8-69 |
| Define Input Parameter by Example by Using the App | 8-73 |
| Define an Input Parameter by Example | 8-73 |
| Specify Input Parameters by Example | 8-74 |
| Specify a String Scalar Input Parameter by Example | 8-75 |
| Specify a Structure Type Input Parameter by Example | 8-75 |
| Specify a Cell Array Type Input Parameter by Example | 8-76 |
| Specify an Enumerated Type Input Parameter by Example | 8-77 |
| Specify a Fixed-Point Input Parameter by Example | 8-78 |
| Specify an Input from an Entry-Point Function Output Type | 8-79 |
| Specify Global Variable Type and Initial Value Using the App | 8-80 |
| Why Specify a Type Definition for Global Variables? | 8-80 |
| Specify a Global Variable Type | 8-80 |
| Define a Global Variable by Example | 8-80 |
| Define or Edit Global Variable Type | 8-81 |
| Define Global Variable Initial Value | 8-81 |
| Define Global Variable Constant Value | 8-82 |
| Remove Global Variables | 8-82 |
| Specify Properties of Entry-Point Function Inputs Using the App .. | 8-83 |
| Why Specify Input Properties? | 8-83 |
| Specify an Input Definition Using the App | 8-83 |

| | |
|---|-------------|
| Detect Unexecuted and Constant-Folded Code | 8-84 |
| What Is Unexecuted Code? | 8-84 |
| Detect Unexecuted Code | 8-84 |
| Fix Unexecuted Code | 8-87 |

Automated Conversion Using Programmatic Workflow

9

| | |
|---|-------------|
| Propose Data Types Based on Simulation Ranges | 9-2 |
| Propose Data Types Based on Derived Ranges | 9-6 |
| Detect Overflows | 9-12 |
| Replace the exp Function with a Lookup Table | 9-16 |
| Replace a Custom Function with a Lookup Table | 9-18 |
| Visualize Differences Between Floating-Point and Fixed-Point Results | 9-20 |
| Copy Relevant Files | 9-20 |
| Prerequisites | 9-20 |
| Enable Plotting Using the Simulation Data Inspector | 9-25 |

Single-Precision Conversion

10

| | |
|---|--------------|
| Generate Single-Precision MATLAB Code | 10-2 |
| Prerequisites | 10-2 |
| Create a Folder and Copy Relevant Files | 10-2 |
| Set Up the Single-Precision Configuration Object | 10-3 |
| Generate Single-Precision MATLAB Code | 10-3 |
| View the Type Proposal Report | 10-4 |
| View Generated Single-Precision MATLAB Code | 10-4 |
| View Potential Data Type Issues | 10-5 |
| Compare the Double-Precision and Single-Precision Variables | 10-5 |
| MATLAB Language Features Supported for Single-Precision Conversion | 10-8 |
| MATLAB Language Features Supported for Single-Precision Conversion | 10-8 |
| MATLAB Language Features Not Supported for Single-Precision Conversion | 10-9 |
| Single-Precision Conversion Best Practices | 10-10 |
| Use Integers for Index Variables | 10-10 |
| Limit Use of assert Statements | 10-10 |
| Initialize MATLAB Class Properties in Constructor | 10-10 |

| | |
|---|-------|
| Provide a Test File That Calls Your MATLAB Function | 10-10 |
| Prepare Your Code for Code Generation | 10-11 |
| Use the -args Option to Specify Input Properties | 10-11 |
| Test Numerics and Log I/O Data | 10-11 |

Fixed-Point Conversion — Manual Conversion

11

| | |
|--|--------------|
| Manual Fixed-Point Conversion Workflow | 11-2 |
| Manual Fixed-Point Conversion Best Practices | 11-3 |
| Create a Test File | 11-3 |
| Prepare Your Algorithm for Code Acceleration or Code Generation | 11-4 |
| Check for Fixed-Point Support for Functions Used in Your Algorithm | 11-5 |
| Manage Data Types and Control Bit Growth | 11-6 |
| Separate Data Type Definitions from Algorithm | 11-6 |
| Convert to Fixed Point | 11-7 |
| Optimize Data Types | 11-9 |
| Optimize Your Algorithm | 11-12 |
| Fixed-Point Design Exploration in Parallel | 11-15 |
| Real-Time Image Acquisition, Image Processing, and Fixed-Point Blob Analysis for Target Practice Analysis | 11-20 |

Code Acceleration and Code Generation from MATLAB for Fixed-Point Algorithms

12

| | |
|---|-------------|
| Code Acceleration and Code Generation from MATLAB | 12-2 |
| Requirements for Generating Compiled C Code Files | 12-3 |
| Functions Supported for Code Acceleration or C Code Generation | 12-4 |
| Workflow for Fixed-Point Code Acceleration and Generation | 12-5 |
| Accelerate Code Using fiaccel | 12-6 |
| Speeding Up Fixed-Point Execution with fiaccel | 12-6 |
| Running fiaccel | 12-6 |
| Generated Files and Locations | 12-6 |
| Data Type Override Using fiaccel | 12-9 |
| Specifying Default fimath Values for MEX Functions | 12-9 |

| | |
|--|--------------|
| File Infrastructure and Paths Setup | 12-11 |
| Compile Path Search Order | 12-11 |
| Naming Conventions | 12-11 |
| Detect and Debug Code Generation Errors | 12-14 |
| Debugging Strategies | 12-14 |
| Error Detection at Design Time | 12-14 |
| Error Detection at Compile Time | 12-15 |
| Set Up C Compiler and Compilation Options | 12-16 |
| Set Up C Compiler | 12-16 |
| C Code Compiler Configuration Object | 12-16 |
| Compilation Options Modification at the Command Line Using Dot Notation | 12-16 |
| How fiaccel Resolves Conflicting Options | 12-17 |
| MEX Configuration Dialog Box Options | 12-18 |
| See Also | 12-20 |
| Specify Configuration Parameters in Command-Line Workflow Interactively | 12-21 |
| Create and Modify Configuration Objects by Using the Dialog Box | 12-21 |
| Additional Functionalities in the Dialog Box | 12-21 |
| Best Practices for Accelerating Fixed-Point Code | 12-24 |
| Recommended Compilation Options for fiaccel | 12-24 |
| Build Scripts | 12-24 |
| Check Code Interactively Using MATLAB Code Analyzer | 12-25 |
| Separating Your Test Bench from Your Function Code | 12-25 |
| Preserving Your Code | 12-25 |
| File Naming Conventions | 12-26 |
| Code Generation Reports | 12-27 |
| Report Generation | 12-27 |
| Report Location | 12-27 |
| Errors and Warnings | 12-27 |
| Files and Functions | 12-27 |
| MATLAB Source | 12-28 |
| MATLAB Variables | 12-29 |
| Code Insights | 12-30 |
| Report Limitations | 12-30 |
| Generate C Code from Code Containing Global Data | 12-31 |
| Workflow Overview | 12-31 |
| Declaring Global Variables | 12-31 |
| Defining Global Data | 12-31 |
| Synchronizing Global Data with MATLAB | 12-32 |
| Limitations of Using Global Data | 12-34 |
| Define Input Properties Programmatically in MATLAB File | 12-35 |
| How to Use assert | 12-35 |
| Rules for Using assert Function | 12-38 |
| Specifying Properties of Primary Fixed-Point Inputs | 12-38 |
| Specifying Properties of Cell Arrays | 12-39 |

| | |
|---|--------------|
| Specifying Class and Size of Scalar Structure | 12-40 |
| Specifying Class and Size of Structure Array | 12-41 |
| Specify Cell Array Inputs at the Command Line | 12-42 |
| Specify Cell Array Inputs by Example | 12-42 |
| Specify the Type of the Cell Array Input | 12-42 |
| Make a Homogeneous Copy of a Type | 12-43 |
| Make a Heterogeneous Copy of a Type | 12-44 |
| Specify Variable-Size Cell Array Inputs | 12-44 |
| Specify Constant Cell Array Inputs | 12-45 |
| Specify Global Cell Arrays at the Command Line | 12-47 |
| Control Run-Time Checks | 12-48 |
| Types of Run-Time Checks | 12-48 |
| When to Disable Run-Time Checks | 12-48 |
| How to Disable Run-Time Checks | 12-49 |
| Fix Run-Time Stack Overflows | 12-50 |
| Code Generation with MATLAB Coder | 12-51 |
| Fixed-Point FIR Code Example Parameter Values | 12-52 |
| Accelerate Code for Variable-Size Data | 12-54 |
| Disable Support for Variable-Size Data | 12-54 |
| Control Dynamic Memory Allocation | 12-54 |
| Accelerate Code for MATLAB Functions with Variable-Size Data .. | 12-55 |
| Accelerate Code for a MATLAB Function That Expands a Vector in a Loop | 12-56 |
| Code Generation Readiness Tool | 12-61 |
| Issues Tab | 12-61 |
| Files Tab | 12-62 |
| Check Code Using the Code Generation Readiness Tool | 12-64 |
| Run Code Generation Readiness Tool at the Command Line | 12-64 |
| Run the Code Generation Readiness Tool From the Current Folder Browser | 12-64 |
| See Also | 12-64 |
| Check Code Using the MATLAB Code Analyzer | 12-65 |
| Fix Errors Detected at Code Generation Time | 12-66 |
| See Also | 12-66 |
| Avoid Multiword Operations in Generated Code | 12-67 |
| Find Potential Data Type Issues in Generated Code | 12-69 |
| Data Type Issues Overview | 12-69 |
| Enable Highlighting of Potential Data Type Issues | 12-69 |
| Find and Address Cumbersome Operations | 12-69 |
| Find and Address Expensive Rounding | 12-70 |
| Find and Address Expensive Comparison Operations | 12-71 |
| Find and Address Multiword Operations | 12-71 |

13

| | |
|---|--------------|
| fi Objects with Simulink | 13-2 |
| View and Edit fi objects in Model Explorer | 13-2 |
| Reading Fixed-Point Data from the Workspace | 13-3 |
| Writing Fixed-Point Data to the Workspace | 13-3 |
| Setting the Value and Data Type of Block Parameters | 13-6 |
| Logging Fixed-Point Signals | 13-6 |
| Accessing Fixed-Point Block Data During Simulation | 13-6 |
| fi Objects with DSP System Toolbox | 13-7 |
| Reading Fixed-Point Signals from the Workspace | 13-7 |
| Writing Fixed-Point Signals to the Workspace | 13-7 |
| Ways to Generate Code | 13-10 |

Calling Functions for Code Generation

14

| | |
|--|--------------|
| Resolution of Function Calls for Code Generation | 14-2 |
| Key Points About Resolving Function Calls | 14-2 |
| Compile Path Search Order | 14-2 |
| When to Use the Code Generation Path | 14-2 |
| Resolution of File Types on Code Generation Path | 14-4 |
| Compilation Directive %#codegen | 14-5 |
| Use MATLAB Engine to Execute a Function Call in Generated Code | 14-6 |
| When To Declare a Function as Extrinsic | 14-6 |
| Use the coder.extrinsic Construct | 14-7 |
| Call MATLAB Functions Using feval | 14-9 |
| Working with mxArray | 14-9 |
| Restrictions on Using Extrinsic Functions | 14-11 |
| Code Generation for Recursive Functions | 14-12 |
| Compile-Time Recursion | 14-12 |
| Run-Time Recursion | 14-12 |
| Disallow Recursion | 14-13 |
| Disable Run-Time Recursion | 14-13 |
| Recursive Function Limitations for Code Generation | 14-13 |
| Force Code Generator to Use Run-Time Recursion | 14-14 |
| Treat the Input to the Recursive Function as a Nonconstant | 14-14 |
| Make the Input to the Recursive Function Variable-Size | 14-15 |
| Assign Output Variable Before the Recursive Call | 14-16 |
| Avoid Duplicate Functions in Generated Code | 14-17 |
| Issue | 14-17 |

| | |
|----------------|-------|
| Cause | 14-17 |
| Solution | 14-17 |

Code Generation for MATLAB Classes

15

| | |
|--|-------|
| MATLAB Classes Definition for Code Generation | 15-2 |
| Language Limitations | 15-2 |
| Code Generation Features Not Compatible with Classes | 15-2 |
| Defining Class Properties for Code Generation | 15-3 |
| Inheritance from Built-In MATLAB Classes Not Supported | 15-5 |
| Classes That Support Code Generation | 15-7 |
| Generate Code for MATLAB Value Classes | 15-8 |
| Generate Code for MATLAB Handle Classes and System Objects | 15-12 |
| Code Generation for Handle Class Destructors | 15-15 |
| Guidelines and Restrictions | 15-15 |
| Behavioral Differences of Objects in Generated Code and in MATLAB | 15-16 |
| Class Does Not Have Property | 15-18 |
| Solution | 15-18 |
| Handle Object Limitations for Code Generation | 15-19 |
| A Variable Outside a Loop Cannot Refer to a Handle Object Allocated Inside a Loop | 15-19 |
| A Handle Object That a Persistent Variable Refers To Must Be a Singleton Object | 15-20 |
| References to Handle Objects Can Appear Undefined | 15-21 |
| System Objects in MATLAB Code Generation | 15-23 |
| Usage Rules and Limitations for System Objects for Generating Code | 15-23 |
| System Objects in codegen | 15-25 |
| System Objects in the MATLAB Function Block | 15-25 |
| System Objects in the MATLAB System Block | 15-25 |
| System Objects and MATLAB Compiler Software | 15-25 |
| Specify Objects as Inputs | 15-26 |
| Consistency Between coder.ClassType Object and Class Definition File | 15-27 |
| Limitations for Using Objects as Entry-Point Function Inputs | 15-27 |
| Work Around Language Limitation: Code Generation Does Not Support Object Arrays | 15-29 |
| Issue | 15-29 |
| Possible Solutions | 15-29 |

| | |
|---|--------------|
| Data Definition Considerations for Code Generation | 16-2 |
| Code Generation for Complex Data | 16-8 |
| Restrictions When Defining Complex Variables | 16-8 |
| Code Generation for Complex Data with Zero-Valued Imaginary Parts | 16-8 |
| Results of Expressions That Have Complex Operands | 16-11 |
| Results of Complex Multiplication with Nonfinite Values | 16-11 |
| Encoding of Characters in Code Generation | 16-12 |
| Array Size Restrictions for Code Generation | 16-13 |
| Code Generation for Constants in Structures and Arrays | 16-14 |
| Code Generation for Strings | 16-16 |
| Limitations | 16-16 |
| Differences Between Generated Code and MATLAB Code | 16-16 |
| Define String Scalar Inputs | 16-17 |
| Define String Scalar Types at the Command Line | 16-17 |
| Code Generation for Sparse Matrices | 16-19 |
| Input Definition | 16-19 |
| Code Generation Guidelines | 16-19 |
| Code Generation Limitations | 16-19 |
| Specify Array Layout in Functions and Classes | 16-21 |
| Specify Array Layout in a Function | 16-21 |
| Query Array Layout of a Function | 16-22 |
| Specify Array Layout in a Class | 16-22 |
| Code Design for Row-Major Array Layout | 16-25 |
| Linear Indexing Uses Column-Major Array Layout | 16-25 |
| Generate Code for Growing Arrays and Cell Arrays with end + 1 Indexing | 16-27 |
| Grow Array with (end + 1) Indexing | 16-27 |
| Grow Cell Array with {end + 1} Indexing | 16-28 |

| | |
|--|-------------|
| Code Generation for Variable Length Argument Lists | 17-2 |
| Generate Code for arguments Block That Validates Input and Output Arguments | 17-3 |
| Supported Features | 17-3 |

| | |
|---|-------|
| Names Must Be Compile-Time Constants | 17-3 |
| Using the Structure That Holds Name-Value Arguments | 17-4 |
| Differences Between Generated Code and MATLAB Code | 17-5 |
| Input Type Specification and arguments blocks | 17-5 |
| Specify Number of Entry-Point Function Input or Output Arguments to Generate | 17-7 |
| Control Number of Input Arguments | 17-7 |
| Control the Number of Output Arguments | 17-7 |
| Code Generation for Anonymous Functions | 17-9 |
| Anonymous Function Limitations for Code Generation | 17-9 |
| Code Generation for Nested Functions | 17-10 |
| Nested Function Limitations for Code Generation | 17-10 |

18 | Defining MATLAB Variables for C/C++ Code Generation

| | |
|--|-------|
| Variables Definition for Code Generation | 18-2 |
| Best Practices for Defining Variables for C/C++ Code Generation | 18-3 |
| Explicitly Define Variables Before Using Them | 18-3 |
| Use Caution When Reassigning Variable Properties | 18-4 |
| Define Variable Numeric Data Types | 18-5 |
| Define Matrices Before Assigning Indexed Variables | 18-5 |
| Index Arrays by Using Constant Value Vectors | 18-5 |
| Eliminate Redundant Copies of Variables in Generated Code | 18-7 |
| When Redundant Copies Occur | 18-7 |
| How to Eliminate Redundant Copies by Defining Uninitialized Variables | 18-7 |
| Defining Uninitialized Variables | 18-7 |
| Reassignment of Variable Properties | 18-9 |
| Reuse the Same Variable with Different Properties | 18-10 |
| When You Can Reuse the Same Variable with Different Properties | 18-10 |
| When You Cannot Reuse Variables | 18-10 |
| Limitations of Variable Reuse | 18-11 |
| Supported Variable Types | 18-13 |
| Edit and Represent Coder Type Objects and Properties | 18-14 |
| Object Properties | 18-14 |
| Legacy Representation of Coder Type Objects | 18-15 |

| | |
|--|--------------|
| When to Generate Code from MATLAB Algorithms | 19-2 |
| When Not to Generate Code from MATLAB Algorithms | 19-2 |
| Which Code Generation Feature to Use | 19-3 |
| Prerequisites for C/C++ Code Generation from MATLAB | 19-4 |
| MATLAB Code Design Considerations for Code Generation | 19-5 |
| See Also | 19-5 |
| Differences Between Generated Code and MATLAB Code | 19-6 |
| Functions that have Multiple Possible Outputs | 19-7 |
| Passing Input Argument Name at Run Time | 19-7 |
| Empty Repeating Input Argument | 19-9 |
| Output Argument Validation of Conditionally-Assigned Outputs . . . | 19-9 |
| Writing to ans Variable | 19-10 |
| Logical Short-Circuiting | 19-11 |
| Loop Index Overflow | 19-11 |
| Indexing for Loops by Using Single Precision Operands | 19-12 |
| Index of an Unentered for Loop | 19-13 |
| Character Size | 19-14 |
| Order of Evaluation in Expressions | 19-14 |
| Name Resolution While Constructing Function Handles | 19-15 |
| Termination Behavior | 19-16 |
| Size of Variable-Size N-D Arrays | 19-16 |
| Size of Empty Arrays | 19-17 |
| Size of Empty Array That Results from Deleting Elements of an Array | 19-17 |
| Growing Variable-Size Column Cell Array That is Initialized as Scalar at Run Time | 19-17 |
| Binary Element-Wise Operations with Single and Double Operands | 19-18 |
| Floating-Point Numerical Results | 19-19 |
| NaN and Infinity | 19-19 |
| Negative Zero | 19-19 |
| Code Generation Target | 19-20 |
| MATLAB Class Property Initialization | 19-20 |
| MATLAB Classes in Nested Property Assignments That Have Set Methods | 19-20 |
| MATLAB Handle Class Destructors | 19-20 |
| Variable-Size Data | 19-21 |
| Complex Numbers | 19-21 |
| Converting Strings with Consecutive Unary Operators to double . | 19-21 |
| Display Function | 19-21 |
| Potential Differences Reporting | 19-23 |
| Addressing Potential Differences Messages | 19-23 |
| Disabling and Enabling Potential Differences Reporting for MATLAB Coder | 19-23 |
| Disabling and Enabling Potential Differences Reporting for Fixed-Point Designer | 19-24 |

| | |
|--|--------------|
| Potential Differences Messages | 19-25 |
| Automatic Dimension Incompatibility | 19-25 |
| mtimes No Dynamic Scalar Expansion | 19-25 |
| Matrix-Matrix Indexing | 19-26 |
| Vector-Vector Indexing | 19-26 |
| Loop Index Overflow | 19-27 |
| | |
| MATLAB Language Features Supported for C/C++ Code Generation | 19-29 |
| MATLAB Features That Code Generation Supports | 19-29 |
| MATLAB Language Features That Code Generation Does Not Support | 19-30 |

Code Generation for Enumerated Data

20

| | |
|---|--------------|
| Code Generation for Enumerations | 20-2 |
| Define Enumerations for Code Generation | 20-2 |
| Allowed Operations on Enumerations | 20-4 |
| MATLAB Toolbox Functions That Support Enumerations | 20-5 |
| | |
| Customize Enumerated Types in Generated Code | 20-7 |
| Specify a Default Enumeration Value | 20-8 |
| Specify a Header File | 20-8 |
| Include Class Name Prefix in Generated Enumerated Type Value Names | 20-9 |
| Generate C++11 Code Containing Ordinary C Enumeration | 20-10 |

Code Generation for Categorical Arrays

21

| | |
|--|-------------|
| Code Generation for Categorical Arrays | 21-2 |
| Define Categorical Arrays for Code Generation | 21-2 |
| Allowed Operations on Categorical Arrays | 21-2 |
| MATLAB Toolbox Functions That Support Categorical Arrays | 21-3 |
| | |
| Define Categorical Array Inputs | 21-6 |
| Define Categorical Array Inputs at the Command Line | 21-6 |
| Representation of Categorical Arrays | 21-6 |
| | |
| Categorical Array Limitations for Code Generation | 21-8 |

Code Generation for Datetime Arrays

22

| | |
|---|------|
| Code Generation for Datetime Arrays | 22-2 |
| Define Datetime Arrays for Code Generation | 22-2 |
| Allowed Operations on Datetime Arrays | 22-2 |
| MATLAB Toolbox Functions That Support Datetime Arrays | 22-2 |
| Define Datetime Array Inputs | 22-5 |
| Define Datetime Array Inputs at the Command Line | 22-5 |
| Representation of Datetime Arrays | 22-5 |
| Datetime Array Limitations for Code Generation | 22-7 |

Code Generation for Duration Arrays

23

| | |
|---|------|
| Code Generation for Duration Arrays | 23-2 |
| Define Duration Arrays for Code Generation | 23-2 |
| Allowed Operations on Duration Arrays | 23-2 |
| MATLAB Toolbox Functions That Support Duration Arrays | 23-3 |
| Define Duration Array Inputs | 23-6 |
| Define Duration Array Inputs at the Command Line | 23-6 |
| Representation of Duration Arrays | 23-6 |
| Duration Array Limitations for Code Generation | 23-8 |

Code Generation for Function Handles

24

| | |
|--|------|
| Function Handle Limitations for Code Generation | 24-2 |
|--|------|

Code Generation for MATLAB Structures

25

| | |
|---|------|
| Structure Definition for Code Generation | 25-2 |
| Structure Operations Allowed for Code Generation | 25-3 |
| Define Scalar Structures for Code Generation | 25-4 |
| Restrictions When Defining Scalar Structures by Assignment | 25-4 |
| Adding Fields in Consistent Order on Each Control Flow Path | 25-4 |
| Restriction on Adding New Fields After First Use | 25-4 |

| | |
|--|-------|
| Define Arrays of Structures for Code Generation | 25-6 |
| Ensuring Consistency of Fields | 25-6 |
| Using repmat to Define an Array of Structures with Consistent Field Properties | 25-6 |
| Defining an Array of Structures by Using struct | 25-6 |
| Defining an Array of Structures Using Concatenation | 25-7 |
| Index Substructures and Fields | 25-8 |
| Assign Values to Structures and Fields | 25-10 |
| Pass Large Structures as Input Parameters | 25-11 |

Functions, Classes, and System Objects Supported for Code Generation

26

| | |
|--|------|
| Functions and Objects Supported for C/C++ Code Generation | 26-2 |
|--|------|

Code Generation for Tables

27

| | |
|--|-------|
| Code Generation for Tables | 27-2 |
| Define Tables for Code Generation | 27-2 |
| Allowed Operations on Tables | 27-2 |
| MATLAB Toolbox Functions That Support Tables | 27-3 |
| Define Table Inputs | 27-5 |
| Define Table Inputs at the Command Line | 27-5 |
| Representation of Tables | 27-5 |
| Table Limitations for Code Generation | 27-8 |
| Creating Tables Limitations | 27-8 |
| Modifying Tables Limitations | 27-8 |
| Using Table Functions Limitations | 27-10 |

Code Generation for Timetables

28

| | |
|--|------|
| Code Generation for Timetables | 28-2 |
| Define Timetables for Code Generation | 28-2 |
| Allowed Operations on Timetables | 28-2 |
| MATLAB Toolbox Functions That Support Timetables | 28-3 |

| | |
|--|--------------|
| Define Timetable Inputs | 28-6 |
| Define Timetable Inputs at the Command Line | 28-6 |
| Representation of Timetables | 28-6 |
| Timetable Limitations for Code Generation | 28-9 |
| Creating Timetables Limitations | 28-9 |
| Modifying Timetables Limitations | 28-10 |
| Using Timetable Functions Limitations | 28-12 |

Code Generation for Variable-Size Data

29

| | |
|---|--------------|
| Code Generation for Variable-Size Arrays | 29-2 |
| Memory Allocation for Variable-Size Arrays | 29-2 |
| Enabling and Disabling Support for Variable-Size Arrays | 29-3 |
| Variable-Size Arrays in a Code Generation Report | 29-3 |
| Control Memory Allocation for Variable-Size Arrays | 29-4 |
| Provide Upper Bounds for Variable-Size Arrays | 29-4 |
| Disable Dynamic Memory Allocation | 29-4 |
| Configure Code Generator to Use Dynamic Memory Allocation for Arrays Bigger Than a Threshold | 29-4 |
| Control Dynamic Memory Allocation for Fixed-Size Arrays | 29-6 |
| Enable Dynamic Memory Allocation for Fixed-Size Arrays | 29-6 |
| Dynamic Memory Allocation Threshold for Fixed-Size Arrays | 29-6 |
| Generating Code for Fixed-Size Arrays | 29-6 |
| Usage Notes and Limitations | 29-7 |
| Specify Upper Bounds for Variable-Size Arrays | 29-8 |
| Specify Upper Bounds for Variable-Size Inputs | 29-8 |
| Specify Upper Bounds for Local Variables | 29-8 |
| Define Variable-Size Data for Code Generation | 29-10 |
| Use a Matrix Constructor with Nonconstant Dimensions | 29-10 |
| Assign Multiple Sizes to the Same Variable | 29-10 |
| Growing an Array by Using (end + 1) | 29-11 |
| Define Variable-Size Data Explicitly by Using coder.varsize | 29-12 |
| Diagnose and Fix Variable-Size Data Errors | 29-15 |
| Diagnosing and Fixing Size Mismatch Errors | 29-15 |
| Diagnosing and Fixing Errors in Detecting Upper Bounds | 29-17 |
| Incompatibilities with MATLAB in Variable-Size Support for Code Generation | 29-18 |
| Incompatibility with MATLAB for Scalar Expansion | 29-18 |
| Incompatibility with MATLAB in Determining Size of Variable-Size N-D Arrays | 29-19 |
| Incompatibility with MATLAB in Determining Size of Empty Arrays | 29-19 |
| Incompatibility with MATLAB in Determining Class of Empty Arrays | 29-21 |

| | |
|---|--------------|
| Incompatibility with MATLAB in Matrix-Matrix Indexing | 29-21 |
| Incompatibility with MATLAB in Vector-Vector Indexing | 29-22 |
| Incompatibility with MATLAB in Matrix Indexing Operations for Code Generation | 29-22 |
| Incompatibility with MATLAB in Concatenating Variable-Size Matrices | 29-23 |
| Differences When Curly-Brace Indexing of Variable-Size Cell Array Inside Concatenation Returns No Elements | 29-23 |
| Variable-Sizing Restrictions for Code Generation of Toolbox Functions | 29-25 |
| Common Restrictions | 29-25 |
| Toolbox Functions with Restrictions for Variable-Size Data | 29-25 |
| Generate Code With Implicit Expansion Enabled | 29-30 |
| Output Size | 29-30 |
| Additional Code Generation | 29-30 |
| Performance Variation | 29-32 |
| Optimize Implicit Expansion in Generated Code | 29-34 |
| Disable Implicit Expansion in Specified Function by Using coder.noImplicitExpansionInFunction | 29-36 |
| Disable Implicit Expansion for Specific Binary Operation by Using coder.sameSizeBinaryOp | 29-37 |
| Disable Implicit Expansion in your Project | 29-38 |
| Representation of Arrays in Generated Code | 29-39 |
| Customize Interface Generation | 29-41 |
| Control Memory Allocation for Fixed-Size Arrays | 29-43 |
| Enable Dynamic Memory Allocation for All Fixed-Size Arrays | 29-43 |
| Enable Dynamic Memory Allocation for Arrays Bigger Than a Threshold | 29-43 |
| Resolve Error: Size Mismatches | 29-45 |
| Issue | 29-45 |
| Possible Solutions | 29-45 |

Code Generation for Cell Arrays

30

| | |
|---|-------------|
| Code Generation for Cell Arrays | 30-2 |
| Homogeneous vs. Heterogeneous Cell Arrays | 30-2 |
| Controlling Whether a Cell Array Is Homogeneous or Heterogeneous | 30-2 |
| Cell Arrays in Reports | 30-3 |
| Control Whether a Cell Array Is Variable-Size | 30-4 |
| Define Cell Array Inputs | 30-6 |

| | |
|--|--------------|
| Cell Array Limitations for Code Generation | 30-7 |
| Cell Array Element Assignment | 30-7 |
| Variable-Size Cell Arrays | 30-8 |
| Definition of Variable-Size Cell Array by Using cell | 30-8 |
| Cell Array Indexing | 30-11 |
| Growing a Cell Array by Using {end + 1} | 30-12 |
| Cell Array Contents | 30-13 |
| Passing Cell Arrays to External C/C++ Functions | 30-13 |

Primary Functions

31

| | |
|--|--------------|
| Specify Properties of Entry-Point Function Inputs | 31-2 |
| Why You Must Specify Input Properties | 31-2 |
| Properties to Specify | 31-2 |
| Rules for Specifying Properties of Primary Inputs | 31-3 |
| Methods for Defining Properties of Primary Inputs | 31-4 |
| Define Input Properties by Example at the Command Line | 31-4 |
| Specify Constant Inputs at the Command Line | 31-6 |
| Specify Variable-Size Inputs at the Command Line | 31-7 |
| Input Type Specification and arguments blocks | 31-8 |
| | |
| Define Input Properties Programmatically in the MATLAB File .. | 31-10 |
| How to Use assert with MATLAB Coder | 31-10 |
| Rules for Using assert Function | 31-14 |
| Specifying General Properties of Primary Inputs | 31-15 |
| Specifying Properties of Primary Fixed-Point Inputs | 31-16 |
| Specifying Properties of Cell Arrays | 31-16 |
| Specifying Class and Size of Scalar Structure | 31-17 |
| Specifying Class and Size of Structure Array | 31-18 |
| | |
| Create and Edit Input Types by Using the Coder Type Editor | 31-19 |
| Open the Coder Type Editor | 31-19 |
| Common Editor Actions | 31-19 |
| Type Browser Pane | 31-20 |
| Type Properties Pane | 31-21 |
| MATLAB Code Pane | 31-22 |

System Objects Supported for Code Generation

32

| | |
|---|-------------|
| Code Generation for System Objects | 32-2 |
|---|-------------|

| | |
|--|--------------|
| Half Precision Code Generation Support | 33-2 |
| Generate Native Half-Precision C Code Using MATLAB Coder ... | 33-13 |
| Generate Native Half-Precision C Code for ARM® Cortex®-A with GCC Compiler | 33-13 |
| Generate Native Half-Precision C Code for ARM Cortex-A with Armclang Compiler | 33-16 |
| Register ARM Target Hardware with Custom Language Implementation | 33-17 |
| What is Half Precision? | 33-19 |
| Half Precision Applications | 33-19 |
| Benefits of Using Half Precision in Embedded Applications | 33-21 |
| Half Precision in MATLAB | 33-22 |
| Half Precision in Simulink | 33-23 |
| Code Generation with Half Precision | 33-23 |

Fixed-Point Designer for Simulink Models

| | |
|---|--------------|
| Sharing Fixed-Point Models | 34-2 |
| Physical Quantities and Measurement Scales | 34-3 |
| Introduction | 34-3 |
| Selecting a Measurement Scale | 34-3 |
| Select a Measurement Scale for Temperature | 34-5 |
| Why Use Fixed-Point Hardware? | 34-9 |
| Why Use the Fixed-Point Designer Software? | 34-10 |
| Developing and Testing Fixed-Point Systems | 34-11 |
| Supported Data Types | 34-13 |
| Configure Blocks with Fixed-Point Output | 34-14 |
| Specify the Output Data Type and Scaling | 34-14 |
| Specify Fixed-Point Data Types with the Data Type Assistant | 34-15 |
| Rounding | 34-17 |
| Overflow Handling | 34-18 |
| Lock the Output Data Type Setting | 34-18 |
| Real-World Values Versus Stored Integer Values | 34-18 |

| | |
|--|--------------|
| Configure Blocks with Fixed-Point Parameters | 34-20 |
| Specify Fixed-Point Values Directly | 34-20 |
| Specify Fixed-Point Values Via Parameter Objects | 34-20 |
| Pass Fixed-Point Data Between Simulink Models and MATLAB .. | 34-22 |
| Read Fixed-Point Data from the Workspace | 34-22 |
| Write Fixed-Point Data to the Workspace | 34-22 |
| Log Fixed-Point Signals | 34-24 |
| Access Fixed-Point Block Data During Simulation | 34-24 |
| Cast from Doubles to Fixed Point | 34-25 |
| Simulate Using Binary-Point-Only Scaling | 34-25 |
| Simulate Using [Slope Bias] Scaling | 34-27 |

Data Types and Scaling

35

| | |
|---|--------------|
| Data Types and Scaling in Digital Hardware | 35-2 |
| Fixed-Point Data Types | 35-2 |
| Binary Point Interpretation | 35-3 |
| Floating-Point Data Types | 35-3 |
| Scaling | 35-5 |
| Binary-Point-Only Scaling | 35-5 |
| Slope and Bias Scaling | 35-6 |
| Unspecified Scaling | 35-6 |
| Quantization | 35-7 |
| Fixed-Point Format | 35-7 |
| Range and Precision | 35-9 |
| Range | 35-9 |
| Precision | 35-10 |
| Fixed-Point Data Type Parameters | 35-11 |
| Range and Precision of an 8-Bit Fixed-Point Data Type — Binary-Point-Only Scaling | 35-11 |
| Range and Precision of an 8-Bit Fixed-Point Data Type — Slope and Bias Scaling | 35-12 |
| Fixed-Point Numbers in Simulink | 35-13 |
| Fixed-Point Data Type and Scaling Notation | 35-13 |
| Display Port Data Types | 35-15 |
| Scaled Doubles | 35-16 |
| What Are Scaled Doubles? | 35-16 |
| When to Use Scaled Doubles | 35-16 |
| Use Scaled Doubles to Avoid Precision Loss | 35-18 |
| Floating-Point Numbers | 35-20 |
| Floating-Point Numbers | 35-20 |

| | |
|--|-------|
| Scientific Notation | 35-20 |
| IEEE 754 Standard for Floating-Point Numbers | 35-21 |
| Range and Precision | 35-22 |
| Exceptional Arithmetic | 35-24 |

Arithmetic Operations

36

| | |
|--|-------|
| Rounding | 36-2 |
| Choosing a Rounding Method | 36-2 |
| Rounding Modes for Fixed-Point Simulink Blocks | 36-5 |
| Fixed-Point Designer Rounding Modes | 36-5 |
| Rounding Mode: Ceiling | 36-8 |
| Rounding Mode: Convergent | 36-9 |
| Rounding Mode: Floor | 36-10 |
| Rounding Mode: Nearest | 36-11 |
| Rounding Mode: Round | 36-12 |
| Rounding Mode: Simplest | 36-14 |
| Optimize Rounding for Casts | 36-14 |
| Optimize Rounding for High-Level Arithmetic Operations | 36-14 |
| Optimize Rounding for Intermediate Arithmetic Operations | 36-15 |
| Rounding Mode: Zero | 36-17 |
| Rounding to Zero Versus Truncation | 36-17 |
| Maximize Precision | 36-18 |
| Pad with Trailing Zeros | 36-18 |
| Constant Scaling for Best Precision | 36-19 |
| Net Slope and Net Bias Precision | 36-21 |
| What are Net Slope and Net Bias? | 36-21 |
| Detect Net Slope and Net Bias Precision Issues | 36-21 |
| Fixed-Point Constant Underflow | 36-22 |
| Fixed-Point Constant Overflow | 36-22 |
| Fixed-Point Constant Precision Loss | 36-23 |
| Detect Fixed-Point Constant Precision Loss | 36-24 |
| Open the Model | 36-24 |
| Detect Fixed-Point Constant Precision Loss | 36-24 |
| Saturation and Wrapping | 36-25 |
| What Are Saturation and Wrapping? | 36-25 |
| Saturation and Wrapping | 36-25 |
| Guard Bits | 36-28 |

| | |
|---|--------------|
| Determine the Range of Fixed-Point Numbers | 36-29 |
| Handle Overflows in Simulink Models | 36-30 |
| Open the Model | 36-30 |
| Simulate Model with Original Diagnostic Settings | 36-30 |
| Adjust Diagnostic Settings | 36-30 |
| Recommendations for Arithmetic and Scaling | 36-31 |
| Arithmetic Operations and Fixed-Point Scaling | 36-31 |
| Addition | 36-31 |
| Accumulation | 36-33 |
| Multiplication | 36-34 |
| Gain | 36-35 |
| Division | 36-36 |
| Summary | 36-37 |
| Parameter and Signal Conversions | 36-39 |
| Introduction | 36-39 |
| Parameter Conversions | 36-39 |
| Signal Conversions | 36-40 |
| Rules for Arithmetic Operations | 36-42 |
| Computational Units | 36-42 |
| Addition and Subtraction | 36-42 |
| Multiplication | 36-43 |
| Division | 36-47 |
| Shifts | 36-48 |
| The Summation Process | 36-49 |
| The Multiplication Process | 36-51 |
| The Division Process | 36-53 |
| Shifts | 36-54 |
| Shifting Bits to the Right | 36-54 |
| Conversions and Arithmetic Operations | 36-55 |

Realization Structures

37

| | |
|--|-------------|
| Realizing Fixed-Point Digital Filters | 37-2 |
| Introduction | 37-2 |
| Realizations and Data Types | 37-2 |
| Targeting an Embedded Processor | 37-3 |
| Introduction | 37-3 |
| Size Assumptions | 37-3 |
| Operation Assumptions | 37-3 |
| Design Rules | 37-4 |

| | |
|------------------------------|-------------|
| Canonical Forms | 37-5 |
|------------------------------|-------------|

Fixed-Point Advisor

38

| | |
|--|-------------|
| Use the Fixed-Point Tool to Prepare a System for Conversion | 38-2 |
| Preparation Checks | 38-2 |

Fixed-Point Tool

39

| | |
|---|--------------|
| Data Type Conversion Overview | 39-2 |
| Methods for Converting a System to Fixed Point | 39-2 |
| Run Management | 39-5 |
| Convert a Referenced Model to Fixed Point | 39-7 |
| Open ex_mdhref_controller Model | 39-7 |
| View Model Hierarchy in the Fixed-Point Tool | 39-7 |
| Viewing Simulation Ranges for Referenced Models | 39-8 |
| Propose Data Types for a Referenced Model | 39-10 |
| Control Views in the Fixed-Point Tool | 39-13 |
| Filter Results by Run | 39-13 |
| Filter Results by Subsystem | 39-13 |
| Control Column Views | 39-13 |
| Use the Explore Tab to Sort and Filter Results | 39-14 |
| Model Multiple Data Type Behaviors Using a Data Dictionary ... | 39-17 |
| Open the Model | 39-17 |
| Explore How the Data Dictionary is Used in the Model | 39-17 |
| Change Data Types of Model Parameters | 39-19 |
| Compare Numerical Response of Sum Block and Sum in MATLAB Function Block | 39-21 |

Convert Floating-Point Model to Fixed Point

40

| | |
|--|-------------|
| Convert Floating-Point Model to Fixed Point | 40-2 |
| Set up the Model | 40-2 |
| Prepare System for Conversion | 40-3 |
| Collect Ranges | 40-5 |
| Convert Data Types | 40-6 |
| Verify New Settings | 40-7 |

| | |
|---|--------------|
| Replace Unsupported Blocks with a Lookup Table Approximation | 40-8 |
| Verify Behavior of System with Lookup Table Approximation | 40-10 |
| Explore Multiple Floating-Point to Fixed-Point Conversions | 40-11 |
| Set up the Model | 40-11 |
| Convert to Fixed-Point Using Default Proposal Settings | 40-11 |
| Convert Using New Proposal Settings | 40-12 |
| Optimize Fixed-Point Data Types for a System | 40-14 |
| Best Practices for Optimizing Data Types | 40-14 |
| Model Management and Exploration | 40-14 |
| Optimize Fixed-Point Data Types | 40-15 |
| Optimize Data Types Using Multiple Simulation Scenarios | 40-20 |
| Optimize Data Types for an FPGA with DSP Slices | 40-23 |
| Use Data Type Optimization to Minimize Operator Counts | 40-30 |
| Open the Model | 40-30 |
| Define Tolerances and Settings | 40-30 |
| Image Denoising Using Fixed-Point Quantized Restricted Boltzmann Machine Algorithm | 40-33 |
| Optimize the Fixed-Point Data Types of a System Using the Fixed- Point Tool | 40-40 |
| Open Model and Define Simulation Scenarios | 40-40 |
| Prepare System for Conversion | 40-41 |
| Optimize Data Types in the Fixed-Point Tool | 40-41 |
| Examine Results | 40-42 |
| Apply Optimized Data Types to the Model | 40-43 |
| Export Optimization Workflow Steps to a MATLAB Script | 40-44 |
| Perform Data Type Optimization with Custom Behavioral Constraints | 40-46 |

Producing Lookup Table Data

41

| | |
|---|--------------|
| Producing Lookup Table Data | 41-2 |
| Worst-Case Error for a Lookup Table | 41-3 |
| Approximate the Square Root Function | 41-3 |
| Create Lookup Tables for a Sine Function | 41-5 |
| Introduction | 41-5 |
| Set Function Parameters for the Lookup Table | 41-5 |
| Specifying Both errmax and nptsmax | 41-12 |
| Comparison of Example Results | 41-13 |
| Use Lookup Table Approximation Functions | 41-14 |

| | |
|--|--------------|
| Optimize Lookup Tables for Memory-Efficiency | 41-15 |
| Optimize an Existing Lookup Table Using the Lookup Table Optimizer | 41-15 |
| Edit the Optimization Settings and Generate a New Approximate | 41-17 |
| Optimize Lookup Tables for Memory-Efficiency Programmatically | 41-19 |
| Approximate a Function Using a Lookup Table | 41-19 |
| Optimize an Existing Lookup Table | 41-22 |
| Visualize Pareto Front for Memory Optimization Versus Absolute Tolerance | 41-28 |
| Compare Approximations Using On Curve and Off Curve Table Values | 41-30 |
| Generate an Optimized Lookup Table as a MATLAB Function ... | 41-36 |
| Generate an Optimized Lookup Table as a MATLAB Function Programmatically | 41-38 |
| Convert Neural Network Algorithms to Fixed-Point Using fxpopt and Generate HDL Code | 41-41 |
| Convert Neural Network Algorithms to Fixed Point and Generate C Code | 41-52 |
| Effects of Spacing on Speed, Error, and Memory Usage | 41-59 |
| Criteria for Comparing Types of Breakpoint Spacing | 41-59 |
| Model That Illustrates Effects of Breakpoint Spacing | 41-59 |
| Data ROM Required for Each Lookup Table | 41-59 |
| Determination of Out-of-Range Inputs | 41-60 |
| How the Lookup Tables Determine Input Location | 41-60 |
| Interpolation for Each Lookup Table | 41-62 |
| Summary of the Effects of Breakpoint Spacing | 41-63 |
| Approximate Functions with a Direct Lookup Table | 41-65 |
| Generate a Two-Dimensional Direct Lookup Table Approximation | 41-65 |
| Generate a Direct Lookup Table Approximation for a Subsystem .. | 41-67 |
| Convert Digit Recognition Neural Network to Fixed Point and Generate C Code | 41-70 |
| Calculate Complex dB Using a Direct Lookup Table | 41-79 |
| Optimize Lookup Tables for Periodic Functions | 41-82 |
| Replace Fitted Curve with Optimized Lookup Table | 41-89 |

| | |
|---|--------------|
| Choosing a Range Collection Method | 42-2 |
| Best Practices for Fixed-Point Conversion Workflow | 42-5 |
| Enable Signal Logging | 42-5 |
| Back Up Your Simulink Model | 42-5 |
| Convert Individual Subsystems | 42-5 |
| Do Not Use “Save as” on Referenced Models and MATLAB Function blocks | 42-5 |
| Use Lock Output Data Type Setting | 42-5 |
| Save Simulink Signal Objects | 42-6 |
| Do Not Use clear all | 42-6 |
| Models That Might Cause Data Type Propagation Errors | 42-7 |
| Iterative Fixed-Point Conversion Using the Fixed-Point Tool | 42-9 |
| Workflow for Automatic Data Typing | 42-10 |
| Set Up the Model | 42-13 |
| Prepare System for Conversion | 42-14 |
| Set Up the Model | 42-14 |
| Select the System Under Design | 42-14 |
| Set Range Collection Method | 42-15 |
| Specify Simulation Input | 42-16 |
| Edit Signal Tolerances | 42-16 |
| Prepare the System for Conversion | 42-17 |
| Specify Behavioral Constraints | 42-18 |
| Specify Signal Tolerances | 42-18 |
| Use Model Verification Blocks | 42-18 |
| Collect Ranges | 42-20 |
| Collect Ranges | 42-20 |
| Explore Collected Ranges | 42-21 |
| Convert Data Types | 42-22 |
| Edit Proposal Settings | 42-22 |
| Propose Data Types | 42-23 |
| Apply Proposed Data Types | 42-24 |
| Examine Results to Resolve Conflicts | 42-26 |
| Proposed Data Type Summary | 42-28 |
| Needs Attention | 42-28 |
| Range Information | 42-28 |
| Examine the Results and Resolve Conflicts | 42-29 |
| Verify New Settings | 42-30 |
| Simulate Using Embedded Types | 42-30 |
| Examine Visualization | 42-30 |
| Compare Results | 42-31 |

| | |
|--|--------------|
| Explore Additional Data Types | 42-34 |
| Edit Proposal Settings | 42-34 |
| Propose, Apply, Simulate, Compare | 42-34 |
| Iterate | 42-34 |
| Restore Model to Original State | 42-34 |
| Restore Model to Original State | 42-36 |
| Get Proposals for Results with Inherited Types | 42-37 |
| How to Get Proposals for Objects That Use an Inherited Output Data Type | 42-37 |
| When the Fixed-Point Tool Will Not Propose for Inherited Data Types | 42-37 |
| Rescale a Fixed-Point Model | 42-39 |
| About the Feedback Controller Example Model | 42-39 |
| Explore the Numerical Behavior of the Model | 42-43 |
| Propose Fraction Lengths Using Simulation Range Data | 42-45 |
| How the Fixed-Point Tool Proposes Data Types | 42-48 |
| How the Fixed-Point Tool Uses Range Information | 42-48 |
| How the Fixed-Point Tool Uses Target Hardware Information | 42-48 |
| How to Get Proposals for Objects That Use an Inherited Output Data Type | 42-48 |
| When the Fixed-Point Tool Will Not Propose for Inherited Data Types | 42-49 |
| How Hardware Implementation Settings Affect Data Type Proposals | 42-50 |
| Open the Model and Specify Hardware Implementation Settings . | 42-50 |
| Propose Word Lengths Based on Simulation Data | 42-51 |
| Propose Data Types For Merged Simulation Ranges | 42-54 |
| Set up the Model | 42-54 |
| Open the Fixed-Point Tool and Prepare the System for Conversion | 42-55 |
| Collect Ranges and Convert to Fixed-Point | 42-55 |
| Verify Fixed-Point Behavior | 42-56 |
| View Simulation Results | 42-57 |
| Compare Runs | 42-57 |
| Histogram Plot of Signal | 42-58 |
| Fixed-Point Instrumentation and Data Type Override | 42-61 |
| Control Instrumentation Settings | 42-61 |
| Control Data Type Override | 42-61 |
| Instrumentation Settings and Data Type Override for a Model Reference Hierarchy | 42-61 |
| Data Type Override Limitations | 42-62 |
| Model Configuration Changes Made During Data Type Optimization | 42-63 |
| Detect downcast | 42-64 |
| Detect underflow | 42-65 |
| Detect precision loss | 42-65 |

| | |
|-------------------------------------|-------|
| Detect overflow | 42-65 |
| Wrap on overflow | 42-65 |
| Saturate on overflow | 42-66 |
| Signal logging | 42-66 |
| Single simulation output | 42-66 |
| Format | 42-66 |
| Time | 42-66 |
| Output | 42-66 |
| Simulation range checking | 42-67 |
| Show port data types | 42-67 |
| Simulation mode | 42-67 |
| Data type override | 42-67 |

Range Analysis

43

| | |
|--|--------------|
| How Range Analysis Works | 43-2 |
| Analyzing a Model with Range Analysis | 43-2 |
| Automatic Stubbing | 43-4 |
| Model Compatibility with Range Analysis | 43-4 |
| How to Derive Ranges | 43-4 |
| Derive Ranges at the Subsystem Level | 43-6 |
| When to Derive Ranges at the Subsystem Level | 43-6 |
| Derive Ranges at the Subsystem Level | 43-6 |
| Derive Ranges Using Design Ranges | 43-8 |
| Open the Model and View Design Ranges | 43-8 |
| Derive Ranges | 43-8 |
| Derive Ranges Using Block Initial Conditions | 43-10 |
| Open the Model | 43-10 |
| Derive Ranges | 43-10 |
| Derive Ranges for Simulink.Parameter Objects | 43-12 |
| Open the ex_derived_min_max_3 Model | 43-12 |
| Examine Gain Parameters | 43-12 |
| Derive Ranges | 43-13 |
| Insufficient Design Range Information | 43-14 |
| Open the ex_derived_min_max_4 Model | 43-14 |
| Collect Ranges | 43-14 |
| Fix Insufficient Design Ranges | 43-15 |
| Troubleshoot Range Analysis of System Objects | 43-16 |
| Provide More Design Range Information | 43-19 |
| Open Model | 43-19 |
| Collect Ranges in the Fixed-Point Tool | 43-19 |
| Provide Additional Design Range Information | 43-20 |

| | |
|---|--------------|
| Fix Design Range Conflicts | 43-22 |
| Open Model | 43-22 |
| Collect Ranges in the Fixed-Point Tool | 43-22 |
| Fix Design Range Conflicts | 43-23 |
| Intermediate Range Results | 43-24 |
| Open Model | 43-24 |
| Collect Ranges in the Fixed-Point Tool | 43-24 |
| Propose Data Types | 43-25 |
| Inspect Result Details | 43-25 |
| Unsupported Simulink Software Features | 43-27 |
| Simulink Blocks Supported for Range Analysis | 43-28 |
| Overview of Simulink Block Support | 43-28 |
| Limitations of Support for Model Blocks | 43-35 |

Range Collection Workflows

44

| | |
|--|-------------|
| Use the Fixed-Point Tool to Explore Numerical Behavior | 44-2 |
| Open the Fixed-Point Direct Form Filter Model | 44-2 |
| Set Up the Model | 44-3 |
| Open the Fixed-Point Tool and Collect Ranges | 44-4 |
| Explore Fixed-Point Behavior of the Model | 44-6 |
| Use Custom Data Type Override Settings for Range Collection ... | 44-9 |

Working with the MATLAB Function Block

45

| | |
|--|-------------|
| Convert MATLAB Function Block to Fixed Point | 45-2 |
| Best Practices for Working with the MATLAB Function Block in the Fixed-Point Tool | 45-2 |
| Open the Model | 45-2 |
| Prepare for Fixed-Point Conversion | 45-3 |
| Collect Range Information | 45-3 |
| Propose Data Types | 45-3 |
| Inspect Code Using the Code View | 45-4 |
| Apply Proposed Data Types | 45-5 |
| Verify Results | 45-7 |
| Replace Functions in a MATLAB Function Block with a Lookup Table | 45-9 |
| Open the Model | 45-9 |
| Replace Sine Function with Lookup Table Approximation | 45-9 |

| | |
|--|-------|
| Best Practices for Working with the MATLAB Function Block in Automated Fixed-Point Conversion Workflows | 45-12 |
| Unsupported MATLAB Function Block Features | 45-12 |
| Control Data Types and Generate Code with MATLAB Function Block | 45-13 |
| Data Type Override with MATLAB Function Block | 45-13 |
| Fixed-Point Data Types with MATLAB Function Block | 45-14 |
| Share Models Containing Fixed-Point MATLAB Function Blocks .. | 45-17 |
| Specify Fixed-Point Math Properties in MATLAB Function Block | 45-19 |
| Generate Fixed-Point FIR Code Using MATLAB Function Block .. | 45-26 |
| Program the MATLAB Function Block | 45-26 |
| Prepare the Inputs | 45-26 |
| Create the Model | 45-27 |
| Define the fimath Object Using the Model Explorer | 45-28 |
| Run the Simulation | 45-28 |

46 Working with Data Objects in the Fixed-Point Workflow

| | |
|--|------|
| Bus Objects in the Fixed-Point Workflow | 46-2 |
| How Data Type Proposals Are Determined for Bus Objects | 46-2 |
| Bus Naming Conventions with Data Type Override | 46-3 |
| Limitations of Bus Objects in the Fixed-Point Workflow | 46-3 |
| Autoscaling Data Objects Using the Fixed-Point Tool | 46-4 |
| Collecting Ranges for Data Objects | 46-4 |
| Data Type Constraints in Data Objects | 46-4 |
| Autoscale a Model Using Data Objects for Data Type Definitions ... | 46-5 |

47 Command Line Interface for the Fixed-Point Tool

| | |
|--|------|
| The Command-Line Interface for the Fixed-Point Tool | 47-2 |
| Convert a Model to Fixed Point Using the Command Line | 47-4 |

48 Code Generation

| | |
|--|------|
| Fixed-Point Code Generation Support | 48-4 |
| Introduction | 48-4 |
| Languages | 48-4 |

| | |
|--|--------------|
| Data Types | 48-4 |
| Rounding Modes | 48-4 |
| Overflow Handling | 48-4 |
| Blocks | 48-4 |
| Scaling | 48-5 |
| Accelerating Fixed-Point Models | 48-6 |
| Using External Mode or Rapid Simulation Target | 48-7 |
| Introduction | 48-7 |
| External Mode | 48-7 |
| Rapid Simulation Target | 48-7 |
| Net Slope Computation | 48-8 |
| Handle Net Slope Computation | 48-8 |
| Use Division to Handle Net Slope Computation | 48-9 |
| Improve Numerical Accuracy of Simulation Results with Rational Approximations to Handle Net Slope | 48-9 |
| Improve Efficiency of Generated Code with Rational Approximations to Handle Net Slope | 48-12 |
| Use Integer Division to Handle Net Slope Computation | 48-15 |
| Control the Generation of Fixed-Point Utility Functions | 48-16 |
| Optimize Generated Code Using Specified Minimum and Maximum Values | 48-16 |
| Eliminate Unnecessary Utility Functions Using Specified Minimum and Maximum Values | 48-18 |
| Optimize Generated Code with the Model Advisor | 48-21 |
| Identify Blocks that Generate Expensive Fixed-Point and Saturation Code | 48-21 |
| Identify Questionable Fixed-Point Operations | 48-23 |
| Identify Blocks that Generate Expensive Rounding Code | 48-25 |
| Lookup Table Optimization | 48-27 |
| Selecting Data Types for Basic Operations | 48-29 |
| Restrict Data Type Word Lengths | 48-29 |
| Avoid Fixed-Point Scalings with Bias | 48-29 |
| Wrap and Round to Floor or Simplest | 48-29 |
| Limit the Use of Custom Storage Classes | 48-30 |
| Use of Shifts by C Code Generation Products | 48-31 |
| Introduction to Shifts by Code Generation Products | 48-31 |
| Modeling Sources of Shifts | 48-32 |
| Controlling Shifts in Generated Code | 48-32 |
| Use Hardware-Efficient Algorithm to Solve Systems of Complex- Valued Linear Equations | 48-34 |
| Implement Hardware-Efficient QR Decomposition Using CORDIC in a Systolic Array | 48-44 |
| Perform QR Factorization Using CORDIC | 48-52 |

| | |
|--|---------------|
| Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using QR Decomposition | 48-79 |
| Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition | 48-82 |
| Implement Hardware-Efficient Real Burst Matrix Solve Using QR Decomposition | 48-85 |
| Implement Hardware-Efficient Real Burst Matrix Solve Using Q-less QR Decomposition | 48-88 |
| Implement Hardware-Efficient Real Partial-Systolic QR Decomposition | 48-90 |
| Implement Hardware-Efficient Real Partial-Systolic Q-less QR Decomposition | 48-93 |
| Implement Hardware-Efficient Real Burst QR Decomposition ... | 48-96 |
| Implement Hardware-Efficient Real Burst Q-less QR Decomposition | 48-99 |
| Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using QR Decomposition | 48-102 |
| Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition | 48-105 |
| Implement Hardware-Efficient Complex Burst Matrix Solve Using QR Decomposition | 48-108 |
| Implement Hardware-Efficient Complex Burst Matrix Solve Using Q-less QR Decomposition | 48-111 |
| Implement Hardware-Efficient Complex Partial-Systolic QR Decomposition | 48-114 |
| Implement Hardware-Efficient Complex Partial-Systolic Q-less QR Decomposition | 48-118 |
| Implement Hardware-Efficient Complex Burst QR Decomposition | 48-121 |
| Implement Hardware-Efficient Complex Burst Q-less QR Decomposition | 48-124 |
| Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using QR Decomposition with Diagonal Loading | 48-127 |
| Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using QR Decomposition with Diagonal Loading | 48-131 |
| Determine Fixed-Point Types for QR Decomposition | 48-135 |

| | |
|--|---------------|
| Determine Fixed-Point Types for Q-less QR Decomposition | 48-138 |
| Algorithms to Determine Fixed-Point Types for Complex Q-less QR Matrix Solve $A'AX=B$ | 48-140 |
| Determine Fixed-Point Types for Complex Q-less QR Matrix Solve $A'AX=B$ | 48-150 |
| Algorithms to Determine Fixed-Point Types for Complex Least-Squares Matrix Solve $AX=B$ | 48-154 |
| Determine Fixed-Point Types for Complex Least-Squares Matrix Solve $AX=B$ | 48-165 |
| Algorithms to Determine Fixed-Point Types for Real Q-less QR Matrix Solve $A'AX=B$ | 48-169 |
| Determine Fixed-Point Types for Real Q-less QR Matrix Solve $A'AX=B$ | 48-179 |
| Algorithms to Determine Fixed-Point Types for Real Least-Squares Matrix Solve $AX=B$ | 48-183 |
| Determine Fixed-Point Types for Real Least-Squares Matrix Solve $AX=B$ | 48-194 |
| Compute Forgetting Factor Required for Streaming Input Data | 48-198 |
| Estimate Standard Deviation of Quantization Noise of Complex-Valued Signal | 48-200 |
| Estimate Standard Deviation of Quantization Noise of Real-Valued Signal | 48-202 |
| Implement Hardware-Efficient Real Partial-Systolic Q-less QR with Forgetting Factor | 48-204 |
| Implement Hardware-Efficient Complex Partial-Systolic Q-less QR with Forgetting Factor | 48-209 |
| Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Forgetting Factor | 48-214 |
| Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Forgetting Factor | 48-220 |
| Implement Hardware-Efficient Complex Burst Matrix Solve Using QR Decomposition with Tikhonov Regularization | 48-226 |
| Implement Hardware-Efficient Complex Burst Matrix Solve Using Q-less QR Decomposition with Tikhonov Regularization | 48-229 |
| Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using QR Decomposition with Tikhonov Regularization | 48-234 |

| | |
|---|---------------|
| Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Tikhonov Regularization | 48-237 |
| Implement Hardware-Efficient Real Burst Matrix Solve Using QR Decomposition with Tikhonov Regularization | 48-242 |
| Implement Hardware-Efficient Real Burst Matrix Solve Using Q-less QR Decomposition with Tikhonov Regularization | 48-245 |
| Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using QR Decomposition with Tikhonov Regularization | 48-250 |
| Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Tikhonov Regularization | 48-253 |
| Determine Fixed-Point Types for Complex Least-Squares Matrix Solve with Tikhonov Regularization | 48-258 |
| Determine Fixed-Point Types for Complex Q-less QR Matrix Solve with Tikhonov Regularization | 48-262 |
| Determine Fixed-Point Types for Real Least-Squares Matrix Solve with Tikhonov Regularization | 48-266 |
| Determine Fixed-Point Types for Real Q-less QR Matrix Solve with Tikhonov Regularization | 48-270 |
| Implement Hardware-Efficient Real Burst Q-less QR with Forgetting Factor | 48-274 |
| Implement Hardware-Efficient Complex Burst Q-less QR with Forgetting Factor | 48-279 |
| Implement Hardware-Efficient Real Burst Matrix Solve Using Q-less QR Decomposition with Forgetting Factor | 48-285 |
| Implement Hardware-Efficient Complex Burst Matrix Solve Using Q-less QR Decomposition with Forgetting Factor | 48-291 |
| Implement Hardware-Efficient Real Burst Asynchronous Matrix Solve Using Q-less QR Decomposition | 48-297 |
| Implement Hardware-Efficient Complex Burst Asynchronous Matrix Solve Using Q-less QR Decomposition | 48-300 |
| How to Use Square Jacobi SVD HDL Optimized Block | 48-303 |
| Implement HDL Optimized SVD in Feedforward Fashion Without Backpressure | 48-307 |
| Implement HDL Optimized SVD with Backpressure Signal and HDL FIFO Block | 48-312 |

| | |
|--|--------------|
| Fixed-Point Versus Built-in Integer Types | 49-2 |
| Negative Fraction Length | 49-3 |
| Fraction Length Greater Than Word Length | 49-5 |
| fi Constructor Does Not Follow globalfimath Rules | 49-7 |
| Issue | 49-7 |
| Possible Solutions | 49-7 |
| Decide Which Workflow is Right for Your Application | 49-8 |
| Tips for Making Generated Code More Efficient | 49-9 |
| fimath Settings for Efficient Code | 49-9 |
| Replace Functions With More Efficient Fixed-Point Implementations | 49-9 |
| Know When a Function is Supported for Instrumentation and Acceleration | 49-11 |
| Resolve Error: Function is not Supported for Fixed-Point Conversion | 49-12 |
| Issue | 49-12 |
| Possible Solutions | 49-12 |
| Resolve Error: fi*non-fi | 49-14 |
| Issue | 49-14 |
| Possible Solutions | 49-14 |
| Resolve Error: Data Type Mismatch | 49-15 |
| Issue | 49-15 |
| Possible Solutions | 49-15 |
| Resolve Error: Mismatched fimath | 49-16 |
| Issue | 49-16 |
| Possible Solutions | 49-16 |
| Why Does the Fixed-Point Converter App Not Propose Data Types for System Objects? | 49-17 |
| Slow Operations in the Fixed-Point Converter App | 49-18 |
| Blocks That Do Not Support Fixed-Point Data Types | 49-19 |
| Issue | 49-19 |
| Possible Solutions | 49-19 |
| Prevent the Fixed-Point Tool from Overriding Integer Data Types | 49-21 |
| The Fixed-Point Tool did not Propose Data Types | 49-22 |
| Issue | 49-22 |

| | |
|--|--------------|
| Possible Solutions | 49-22 |
| Fraction Lengths and Fixed-Point Numbers | 49-23 |
| Fraction Length Greater Than Word Length | 49-23 |
| Negative Fraction Length | 49-23 |
| Why am I missing data type proposals for MATLAB Function block variables? | 49-24 |
| Data Type Propagation Errors After Applying Proposed Data Types | 49-25 |
| Shared Data Type Groups | 49-25 |
| Block Constraints | 49-26 |
| Internal Block Rules | 49-26 |
| Resolve Range Analysis Issues | 49-27 |
| Issue | 49-27 |
| Possible Solutions | 49-27 |
| Data Type Mismatch and Structure Initial Conditions | 49-28 |
| Specify Bus Signal Initial Conditions Using Simulink.Parameter Objects | 49-28 |
| Data Type Mismatch and Masked Atomic Subsystems | 49-28 |
| Reconversion Using the Fixed-Point Tool | 49-30 |
| Data Type Optimization Not Successful | 49-31 |
| Issue | 49-31 |
| Possible Solutions | 49-31 |
| Compile-Time Recursion Limit Reached | 49-33 |
| Issue | 49-33 |
| Cause | 49-33 |
| Solutions | 49-33 |
| Force Run-Time Recursion | 49-33 |
| Increase the Compile-Time Recursion Limit | 49-35 |
| Output Variable Must Be Assigned Before Run-Time Recursive Call | 49-36 |
| Issue | 49-36 |
| Cause | 49-36 |
| Solution | 49-36 |
| Unable to Determine That Every Element of Cell Array Is Assigned | 49-39 |
| Issue | 49-39 |
| Cause | 49-39 |
| Solution | 49-40 |
| Nonconstant Index into varargin or varargout in a for-Loop | 49-43 |
| Issue | 49-43 |
| Cause | 49-43 |
| Solution | 49-43 |

| | |
|--|------|
| Getting Started with Single Precision Converter | 50-2 |
| Select System Under Design | 50-2 |
| Check Compatibility | 50-2 |
| Convert | 50-3 |
| Verify | 50-3 |
| Convert a System to Single Precision | 50-5 |
| Open Model | 50-5 |
| Convert to Single Precision | 50-5 |

Simulink Half Precision

| | |
|--|-------|
| The Half-Precision Data Type in Simulink | 51-2 |
| Math Operations in Half-Precision | 51-2 |
| Software Features Supported for Half Precision | 51-2 |
| Generate Code for Half Precision Systems | 51-3 |
| Generate Native Half-Precision C Code from Simulink Models ... | 51-5 |
| Generate Native Half-Precision C Code for ARM Cortex-A with Armclang Compiler | 51-5 |
| Register ARM Cortex-A with GCC Compiler | 51-8 |
| Register ARM Target Hardware with Custom Language Implementation | 51-9 |
| Register Other Hardware Targets for Half Precision | 51-9 |
| Half-Precision Field-Oriented Control Algorithm | 51-11 |
| Image Quantization with Half-Precision Data Types | 51-14 |
| Convert Single Precision Lookup Table to Half Precision | 51-15 |
| Digit Classification with Half-Precision Data Types | 51-20 |

Design Cost Estimation

| | |
|--|------|
| Design Cost Model Metrics | 52-2 |
| Data Segment Estimate | 52-2 |
| Operator Count | 52-3 |
| How to Collect Design Cost Metrics | 52-4 |
| Open Project | 52-4 |
| Collect Metric Results | 52-4 |
| Access High-Level Results Programmatically | 52-5 |

| | |
|--|------|
| Generate Report to Access Detailed Results | 52-6 |
| Operator Count | 52-6 |
| Data Segment Table | 52-8 |

Fixed-Point HDL-Optimized Blocks

53

Choose a Block for HDL-Optimized Fixed-Point Matrix Operations

| | |
|--|------|
| | 53-2 |
| Define the Problem to Solve | 53-2 |
| Choose an Architecture | 53-3 |
| Linear System Solvers: Select Synchronous or Asynchronous Operation | 53-4 |
| Data Complexity | 53-5 |
| Hardware Control Signals | 53-5 |

Writing Fixed-Point S-Functions

A

| | |
|---|------|
| Data Type Support | A-2 |
| Supported Data Types | A-2 |
| The Treatment of Integers | A-2 |
| Data Type Override | A-3 |
| Structure of the S-Function | A-4 |
| Storage Containers | A-5 |
| Introduction | A-5 |
| Storage Containers in Simulation | A-5 |
| Storage Containers in Code Generation | A-7 |
| Data Type IDs | A-9 |
| The Assignment of Data Type IDs | A-9 |
| Registering Data Types | A-10 |
| Setting and Getting Data Types | A-11 |
| Getting Information About Data Types | A-11 |
| Converting Data Types | A-13 |
| Overflow Handling and Rounding Methods | A-14 |
| Tokens for Overflow Handling and Rounding Methods | A-14 |
| Overflow Logging Structure | A-14 |
| Create MEX-Files | A-16 |
| Fixed-Point S-Function Examples | A-17 |
| Get the Input Port Data Type | A-18 |
| Set the Output Port Data Type | A-20 |

| | |
|---|-------------|
| Interpret an Input Value | A-21 |
| Write an Output Value | A-23 |
| Determine Output Type Using the Input Type | A-25 |
| API Function Reference | A-26 |

Fixed-Point Designer Examples

54

| | |
|---|---------------|
| Create Fixed-Point Data | 54-2 |
| Perform Fixed-Point Arithmetic | 54-11 |
| Set Fixed-Point Math Attributes | 54-17 |
| View Fixed-Point Number Circles | 54-25 |
| Perform Binary-Point Scaling | 54-35 |
| Compute Quantization Error | 54-39 |
| Detect Limit Cycles in Fixed-Point State-Space Systems | 54-46 |
| Develop Fixed-Point Algorithms | 54-55 |
| Perform QR Factorization Using CORDIC | 54-62 |
| Compute Square Root Using CORDIC | 54-89 |
| Normalize Data for Lookup Tables | 54-96 |
| Implement Fixed-Point Log2 Using Lookup Table | 54-100 |
| Implement Fixed-Point Square Root Using Lookup Table | 54-104 |
| Convert Fast Fourier Transform (FFT) to Fixed Point | 54-108 |
| Set Data Types Using Min/Max Instrumentation | 54-125 |
| Implement FIR Filter Algorithm for Floating-Point and Fixed-Point Types Using cast and zeros | 54-136 |
| Accelerate Fixed-Point Simulation | 54-155 |
| Double to Fixed-Point Conversion | 54-158 |
| Fixed-Point to Fixed-Point Inherited Conversion | 54-160 |
| Fixed-Point Prelookup and Interpolation | 54-161 |

| | |
|--|---------------|
| Sharing Parameters in Prelookup and Interpolation Blocks | 54-163 |
| High Precision Calculations in Interpolation Block | 54-164 |
| Shared Utility Functions for Prelookup Blocks | 54-165 |
| Fixed-Point Multiword Operations In Generated Code | 54-166 |
| Fixed-Point Multiplication Helper Functions in Generated Code | 54-171 |
| Fixed-Point Optimizations Using Specified Minimum and Maximum Values | 54-178 |
| Fixed-Point Function Approximation | 54-182 |
| Fixed-Point Conversion Using Fixed-Point Tool and Derived Range Analysis | 54-191 |
| Fixed-Point Tool | 54-194 |
| Fixed-Point S-Functions: Querying Properties | 54-204 |
| Fixed-Point S-Functions: Arithmetic Shift | 54-205 |
| Fixed-Point S-Functions: Fixed-Point Source | 54-208 |
| Fixed-Point S-Functions: Data Type Propagation | 54-209 |
| Fixed-Point S-Functions: Product and Sum | 54-213 |
| How to Use HDL Optimized Normalized Reciprocal | 54-215 |
| Implement Hardware-Efficient Real Divide HDL Optimized | 54-224 |
| Implement Hardware-Efficient Complex Divide HDL Optimized | 54-226 |
| Implement Hardware-Efficient Hyperbolic Tangent | 54-228 |
| Implement HDL Optimized Modulo By Constant | 54-231 |
| Fixed-Point HDL-Optimized Minimum-Variance Distortionless- Response (MVDR) Beamformer | 54-235 |
| Hardware-Efficient Euler Rotations Using CORDIC | 54-240 |

Simulation Data Inspector

| | |
|---|-------------|
| View Data in the Simulation Data Inspector | 55-2 |
| View Logged Data | 55-2 |

| | |
|---|--------------|
| Import Data from the Workspace or a File | 55-3 |
| View Complex Data | 55-5 |
| View String Data | 55-6 |
| View Frame-Based Data | 55-9 |
| View Event-Based Data | 55-9 |
| Import Data from a CSV File into the Simulation Data Inspector | 55-11 |
| Basic File Format | 55-11 |
| Multiple Time Vectors | 55-11 |
| Signal Metadata | 55-12 |
| Import Data from a CSV File | 55-13 |
| Microsoft Excel Import, Export, and Logging Format | 55-15 |
| Basic File Format | 55-15 |
| Multiple Time Vectors | 55-15 |
| Signal Metadata | 55-16 |
| User-Defined Data Types | 55-18 |
| Complex, Multidimensional, and Bus Signals | 55-20 |
| Function-Call Signals | 55-21 |
| Simulation Parameters | 55-21 |
| Multiple Runs | 55-21 |
| Configure the Simulation Data Inspector | 55-23 |
| Logged Data Size and Location | 55-23 |
| Archive Behavior and Run Limit | 55-24 |
| Incoming Run Names and Location | 55-25 |
| Signal Metadata to Display | 55-26 |
| Signal Selection on the Inspect Pane | 55-27 |
| How Signals Are Aligned for Comparison | 55-27 |
| Colors Used to Display Comparison Results | 55-28 |
| Signal Grouping | 55-28 |
| Data to Stream from Parallel Simulations | 55-29 |
| Options for Saving and Loading Session Files | 55-29 |
| Signal Display Units | 55-29 |
| How the Simulation Data Inspector Compares Data | 55-31 |
| Signal Alignment | 55-31 |
| Synchronization | 55-32 |
| Interpolation | 55-33 |
| Tolerance Specification | 55-33 |
| Limitations | 55-35 |
| Save and Share Simulation Data Inspector Data and Views | 55-36 |
| Save and Load Simulation Data Inspector Sessions | 55-36 |
| Share Simulation Data Inspector Views | 55-37 |
| Share Simulation Data Inspector Plots | 55-37 |
| Create Simulation Data Inspector Report | 55-38 |
| Export Data to the Workspace or a File | 55-39 |
| Export Video Signal to an MP4 File | 55-40 |
| Inspect and Compare Data Programmatically | 55-42 |
| Create a Run and View the Data | 55-42 |
| Compare Two Signals in the Same Run | 55-43 |
| Compare Runs with Global Tolerance | 55-44 |

| | |
|---|--------------|
| Analyze Simulation Data Using Signal Tolerances | 55-45 |
| Limit the Size of Logged Data | 55-48 |
| Limit the Number of Runs Retained in the Simulation Data Inspector Archive | 55-48 |
| Specify a Minimum Disk Space Requirement or Maximum Size for Logged Data | 55-48 |
| View Data Only During Simulation | 55-49 |
| Reduce the Number of Data Points Logged from Simulation | 55-49 |

Fixed-Point Designer for MATLAB Code

Fixed-Point Concepts

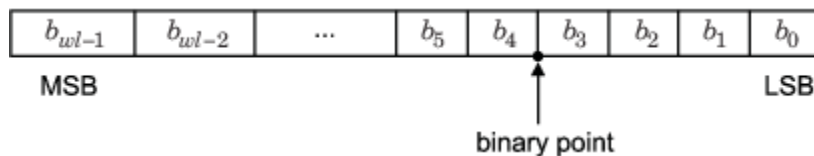
- “Fixed-Point Data Types” on page 1-2
- “Scaling” on page 1-3
- “Compute Slope and Bias” on page 1-4
- “Precision and Range” on page 1-6
- “Arithmetic Operations” on page 1-9
- “fi Objects and C Integer Data Types” on page 1-18

Fixed-Point Data Types

In digital hardware, numbers are stored in binary words. A binary word is a fixed-length sequence of bits (1's and 0's). How hardware components or software functions interpret this sequence of 1's and 0's is defined by the data type. Binary numbers are represented as either fixed-point or floating-point data types.

A fixed-point data type is characterized by the word length in bits, the position of the binary point, and whether it is signed or unsigned. The position of the binary point is the means by which fixed-point values are scaled and interpreted.

For example, a binary representation of a generalized fixed-point number (either signed or unsigned) is shown below:



where

- b_i is the i^{th} binary digit.
- wl is the word length in bits.
- b_{wl-1} is the location of the most significant, or highest, bit (MSB).
- b_0 is the location of the least significant, or lowest, bit (LSB).
- The binary point is shown four places to the left of the LSB. In this example, the number is said to have four fractional bits, or a fraction length of four.

Fixed-point data types can be either signed or unsigned. Whether a fixed-point value is signed or unsigned is usually not encoded explicitly within the binary word; that is, there is no sign bit. Instead, the sign information is implicitly defined within the computer architecture.

Signed binary fixed-point numbers are typically represented in computer hardware in one of three ways:

- Sign/magnitude - One bit of a binary word is always the dedicated sign bit, while the remaining bits of the word encode the magnitude of the number. Negation using sign/magnitude representation consists of flipping the sign bit from 0 (positive) to 1 (negative), or from 1 to 0.
- One's complement - Negating a binary number in one's complement requires a bitwise complement. That is, all 0's are flipped to 1's and all 1's are flipped to 0's. In one's complement notation there are two ways to represent zero. A binary word of all 0's represents "positive" zero, while a binary word of all 1's represents "negative" zero.
- Two's complement - Negation using signed two's complement representation consists of a bit inversion (translation into one's complement) followed by the binary addition of a one. For example, the two's complement of 000101 is 111011.

Two's complement is the most common representation of signed fixed-point numbers and is the only representation used by Fixed-Point Designer documentation.

Scaling

Fixed-point numbers can be encoded according to the scheme

$$\text{real-world value} = (\text{slope} \times \text{integer}) + \text{bias}$$

where the slope can be expressed as

$$\text{slope} = \text{slope adjustment factor} \times 2^{\text{fixed exponent}}$$

The integer is sometimes called the *stored integer*. This is the raw binary number, in which the binary point assumed to be at the far right of the word. In Fixed-Point Designer documentation, the negative of the fixed exponent is often referred to as the *fraction length*.

The slope and bias together represent the scaling of the fixed-point number. In a number with zero bias, only the slope affects the scaling. A fixed-point number that is only scaled by binary point position is equivalent to a number in [Slope Bias] representation that has a bias equal to zero and a slope adjustment factor equal to one. This is referred to as binary point-only scaling or power-of-two scaling:

$$\text{real-world value} = 2^{\text{fixed exponent}} \times \text{integer}$$

or

$$\text{real-world value} = 2^{-\text{fraction length}} \times \text{integer}$$

Fixed-Point Designer software supports both binary point-only scaling and [Slope Bias] scaling.

Note For examples of binary point-only scaling, see the Fixed-Point Designer “Perform Binary-Point Scaling” on page 54-35 example.

For an example of how to compute slope and bias in MATLAB®, see “Compute Slope and Bias” on page 1-4

Compute Slope and Bias

In this section...

“What Is Slope Bias Scaling?” on page 1-4

“Compute Slope and Bias” on page 1-4

What Is Slope Bias Scaling?

With slope bias scaling, you must specify the slope and bias of a number. The real-world value of a slope bias scaled number can be represented by:

$$\text{real-worldvalue} = (\text{slope} \times \text{integer}) + \text{bias}$$

$$\text{slope} = \text{slopedjustmentfactor} \times 2^{\text{fixedexponent}}$$

Compute Slope and Bias

Start with the endpoints that you want, signedness, and word length.

```
lower_bound = 999;
upper_bound = 1000;
is_signed = true;
word_length = 16;
```

To find the range of a `fi` object with a specified word length and signedness, use the `range` function.

```
[Q_min, Q_max] = range(fi([], is_signed, word_length, 0));
```

To find the slope and bias, solve the system of equations:

$$\text{lower_bound} = \text{slope} * \text{Q_min} + \text{bias}$$

$$\text{upper_bound} = \text{slope} * \text{Q_max} + \text{bias}$$

Rewrite these equations in matrix form.

$$\begin{bmatrix} \text{lower_bound} \\ \text{upper_bound} \end{bmatrix} = \begin{bmatrix} \text{Q_min} & 1 \\ \text{Q_max} & 1 \end{bmatrix} \times \begin{bmatrix} \text{slope} \\ \text{bias} \end{bmatrix}$$

Solve for the slope and bias.

```
A = double ([Q_min, 1; Q_max, 1]);
b = double ([lower_bound; upper_bound]);
x = A\b;
format long g
```

To find the slope, or precision, call the first element of the slope-bias vector, `x`.

```
slope = x(1)
```

```
slope =
```

```
1.52590218966964e-05
```

To find the bias, call the second element of vector `x`.

```
bias = x(2)
```

```
bias =
```

```
    999.500007629511
```

Create a `numericType` object with slope bias scaling.

```
T = numericType(is_signed, word_length, slope, bias)
```

```
T =
```

```
    DataTypeMode: Fixed-point: slope and bias scaling
      Signedness: Signed
      WordLength: 16
      Slope: 1.5259021896696368e-5
      Bias: 999.500007629511
```

Create a `fi` object with `numericType` `T`.

```
a = fi(999.255, T)
```

```
a =
```

```
    999.254993514916
```

```
    DataTypeMode: Fixed-point: slope and bias scaling
      Signedness: Signed
      WordLength: 16
      Slope: 1.5259021896696368e-5
      Bias: 999.500007629511
```

Verify that the `fi` object that you created has the correct specifications by finding the range of `a`.

```
range(a)
```

```
ans =
```

```
    999    1000
```

```
    DataTypeMode: Fixed-point: slope and bias scaling
      Signedness: Signed
      WordLength: 16
      Slope: 1.5259021896696368e-5
      Bias: 999.500007629511
```

Precision and Range

In this section...

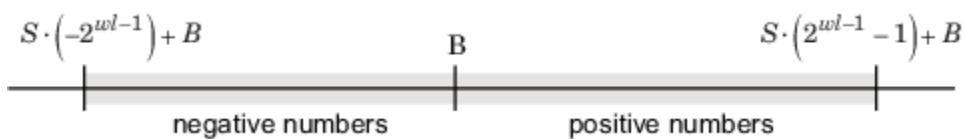
“Range” on page 1-6

“Precision” on page 1-7

Note You must pay attention to the precision and range of the fixed-point data types and scalings you choose in order to know whether rounding methods will be invoked or if overflows or underflows will occur.

Range

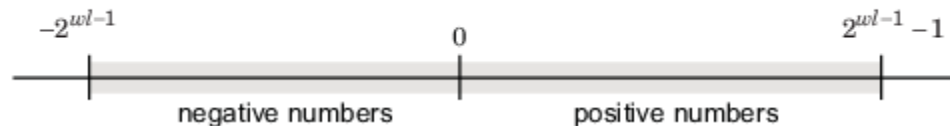
The range is the span of numbers that a fixed-point data type and scaling can represent. The range of representable numbers for a two's complement fixed-point number of word length wl , scaling S and bias B is illustrated below:



For both signed and unsigned fixed-point numbers of any data type, the number of different bit patterns is 2^{wl} .

For example, in two's complement, negative numbers must be represented as well as zero, so the maximum value is $2^{wl-1} - 1$. Because there is only one representation for zero, there are an unequal number of positive and negative numbers. This means there is a representation for -2^{wl-1} but not for 2^{wl-1} :

For slope = 1 and bias = 0:



Overflow Handling

Because a fixed-point data type represents numbers within a finite range, overflows and underflows can occur if the result of an operation is larger or smaller than the numbers in that range.

Fixed-Point Designer software allows you to either *saturate* or *wrap* overflows. Saturation represents positive overflows as the largest positive number in the range being used, and negative overflows as the largest negative number in the range being used. Wrapping uses modulo arithmetic to cast an overflow back into the representable range of the data type.

When you create a `fi` object, any overflows are saturated. The `OverflowAction` property of the default `fimath` is `saturate`. You can log overflows and underflows by setting the `LoggingMode` property of the `fipref` object to `on`.

Precision

The precision of a fixed-point number is the difference between successive values representable by its data type and scaling, which is equal to the value of its least significant bit. The value of the least significant bit, and therefore the precision of the number, is determined by the number of fractional bits. A fixed-point value can be represented to within half of the precision of its data type and scaling.

For example, a fixed-point representation with four bits to the right of the binary point has a precision of 2^{-4} or 0.0625, which is the value of its least significant bit. Any number within the range of this data type and scaling can be represented to within $(2^{-4})/2$ or 0.03125, which is half the precision. This is an example of representing a number with finite precision.

Rounding Methods

When you represent numbers with finite precision, not every number in the available range can be represented exactly. If a number cannot be represented exactly by the specified data type and scaling, a rounding method is used to cast the value to a representable number. Although precision is always lost in the rounding operation, the cost of the operation and the amount of bias that is introduced depends on the rounding method itself. To provide you with greater flexibility in the trade-off between cost and bias, Fixed-Point Designer software currently supports the following rounding methods:

- **Ceiling** rounds to the closest representable number in the direction of positive infinity.
- **Convergent** rounds to the closest representable number. In the case of a tie, **convergent** rounds to the nearest even number. This is the least biased rounding method provided by the toolbox.
- **Zero** rounds to the closest representable number in the direction of zero.
- **Floor**, which is equivalent to two's complement truncation, rounds to the closest representable number in the direction of negative infinity.
- **Nearest** rounds to the closest representable number. In the case of a tie, **nearest** rounds to the closest representable number in the direction of positive infinity. This rounding method is the default for `fi` object creation and `fi` arithmetic.
- **Round** rounds to the closest representable number. In the case of a tie, the **round** method rounds:
 - Positive numbers to the closest representable number in the direction of positive infinity.
 - Negative numbers to the closest representable number in the direction of negative infinity.

Choosing a Rounding Method

Each rounding method has a set of inherent properties. Depending on the requirements of your design, these properties could make the rounding method more or less desirable to you. By knowing the requirements of your design and understanding the properties of each rounding method, you can determine which is the best fit for your needs. The most important properties to consider are:

- **Cost** — Independent of the hardware being used, how much processing expense does the rounding method require?
 - **Low** — The method requires few processing cycles.
 - **Moderate** — The method requires a moderate number of processing cycles.
 - **High** — The method requires more processing cycles.

Note The cost estimates provided here are hardware independent. Some processors have rounding modes built-in, so consider carefully the hardware you are using before calculating the true cost of each rounding mode.

- Bias — What is the expected value of the rounded values minus the original values: $E(\hat{\theta} - \theta)$?
 - $E(\hat{\theta} - \theta) < 0$ — The rounding method introduces a negative bias.
 - $E(\hat{\theta} - \theta) = 0$ — The rounding method is unbiased.
 - $E(\hat{\theta} - \theta) > 0$ — The rounding method introduces a positive bias.

The following table shows a comparison of the different rounding methods available in the Fixed-Point Designer product.

| Fixed-Point Designer Rounding Mode | Cost | Bias |
|------------------------------------|----------|---|
| Ceiling | Low | Large positive |
| Convergent | High | Unbiased |
| Zero | Low | <ul style="list-style-type: none"> • Large positive for negative samples • Unbiased for samples with evenly distributed positive and negative values • Large negative for positive samples |
| Floor | Low | Large negative |
| Nearest | Moderate | Small positive |
| Round | High | <ul style="list-style-type: none"> • Small negative for negative samples • Unbiased for samples with evenly distributed positive and negative values • Small positive for positive samples |
| Simplest (Simulink® only) | Low | Depends on the operation |

Arithmetic Operations

In this section...

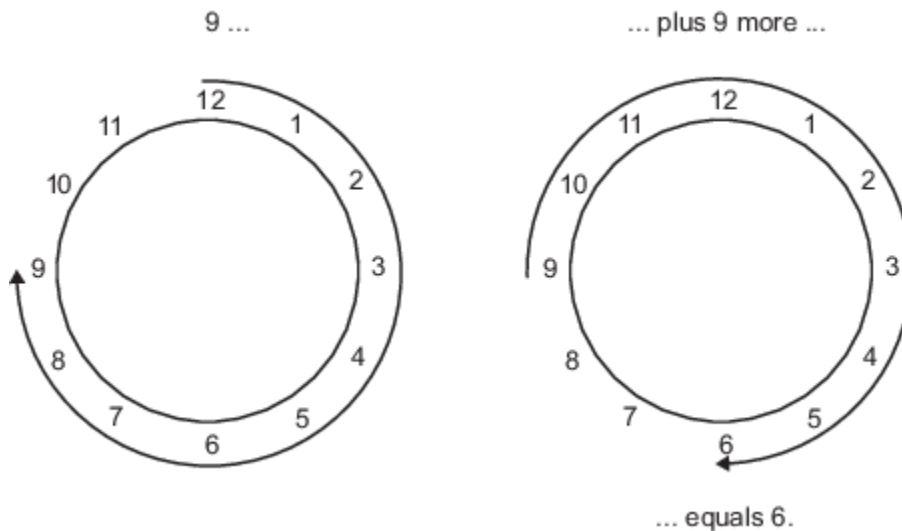
"Modulo Arithmetic" on page 1-9
 "Two's Complement" on page 1-9
 "Addition and Subtraction" on page 1-10
 "Multiplication" on page 1-11
 "Casts" on page 1-15

Note These sections will help you understand what data type and scaling choices result in overflows or a loss of precision.

Modulo Arithmetic

Binary math is based on modulo arithmetic. Modulo arithmetic uses only a finite set of numbers, wrapping the results of any calculations that fall outside the given set back into the set.

For example, the common everyday clock uses modulo 12 arithmetic. Numbers in this system can only be 1 through 12. Therefore, in the "clock" system, 9 plus 9 equals 6. This can be more easily visualized as a number circle:



Similarly, binary math can only use the numbers 0 and 1, and any arithmetic results that fall outside this range are wrapped "around the circle" to either 0 or 1.

Two's Complement

Two's complement is a way to interpret a binary number. In two's complement, positive numbers always start with a 0 and negative numbers always start with a 1. If the leading bit of a two's complement number is 0, the value is obtained by calculating the standard binary value of the

number. If the leading bit of a two's complement number is 1, the value is obtained by assuming that the leftmost bit is negative, and then calculating the binary value of the number. For example,

$$01 = (0 + 2^0) = 1$$
$$11 = ((-2^1) + (2^0)) = (-2 + 1) = -1$$

To compute the negative of a binary number using two's complement,

- 1 Take the one's complement, or “flip the bits.”
- 2 Add a $2^{(-FL)}$ using binary math, where FL is the fraction length.
- 3 Discard any bits carried beyond the original word length.

For example, consider taking the negative of 11010 (-6). First, take the one's complement of the number, or flip the bits:

$$11010 \rightarrow 00101$$

Next, add a 1, wrapping all numbers to 0 or 1:

$$\begin{array}{r} 00101 \\ +1 \\ \hline 00110 \text{ (6)} \end{array}$$

Addition and Subtraction

The addition of fixed-point numbers requires that the binary points of the addends be aligned. The addition is then performed using binary arithmetic so that no number other than 0 or 1 is used.

For example, consider the addition of 010010.1 (18.5) with 0110.110 (6.75):

$$\begin{array}{r} 010010.1 \text{ (18.5)} \\ +0110.110 \text{ (6.75)} \\ \hline 011001.010 \text{ (25.25)} \end{array}$$

Fixed-point subtraction is equivalent to adding while using the two's complement value for any negative values. In subtraction, the addends must be sign-extended to match each other's length. For example, consider subtracting 0110.110 (6.75) from 010010.1 (18.5):

$$\begin{array}{r} 010010.100 \text{ (18.5)} \\ -0110.110 \text{ (6.75)} \\ \hline \end{array}$$

The default `fmath` has a value of 1 (true) for the `CastBeforeSum` property. This casts addends to the sum data type before addition. Therefore, no further shifting is necessary during the addition to line up the binary points.

If `CastBeforeSum` has a value of 0 (false), the addends are added with full precision maintained. After the addition the sum is then quantized.

Multiplication

The multiplication of two's complement fixed-point numbers is directly analogous to regular decimal multiplication, with the exception that the intermediate results must be sign-extended so that their left sides align before you add them together.

For example, consider the multiplication of 10.11 (-1.25) with 011 (3):

$$\begin{array}{r}
 10.11 \text{ (-1.25)} \\
 \underline{011 \text{ (3)}} \\
 11011 \\
 \underline{1011} \\
 1100.01 \text{ (-3.75)}
 \end{array}$$

The extra 1 is the result of necessary sign extension.

The number of fractional bits of the result is the sum of the number of fractional bits of the factors.

Multiplication Data Types

The following diagrams show the data types used for fixed-point multiplication using Fixed-Point Designer software. The diagrams illustrate the differences between the data types used for real-real, complex-real, and complex-complex multiplication.

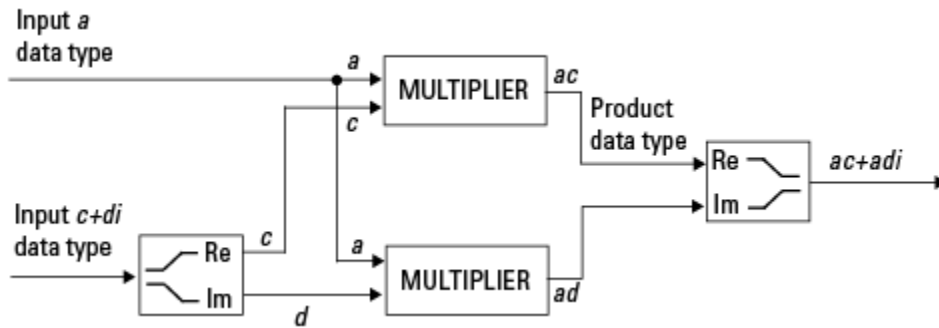
Real-Real Multiplication

The following diagram shows the data types used by the toolbox in the multiplication of two real numbers. The software returns the output of this operation in the product data type, which is governed by the `fimath` object `ProductMode` property.



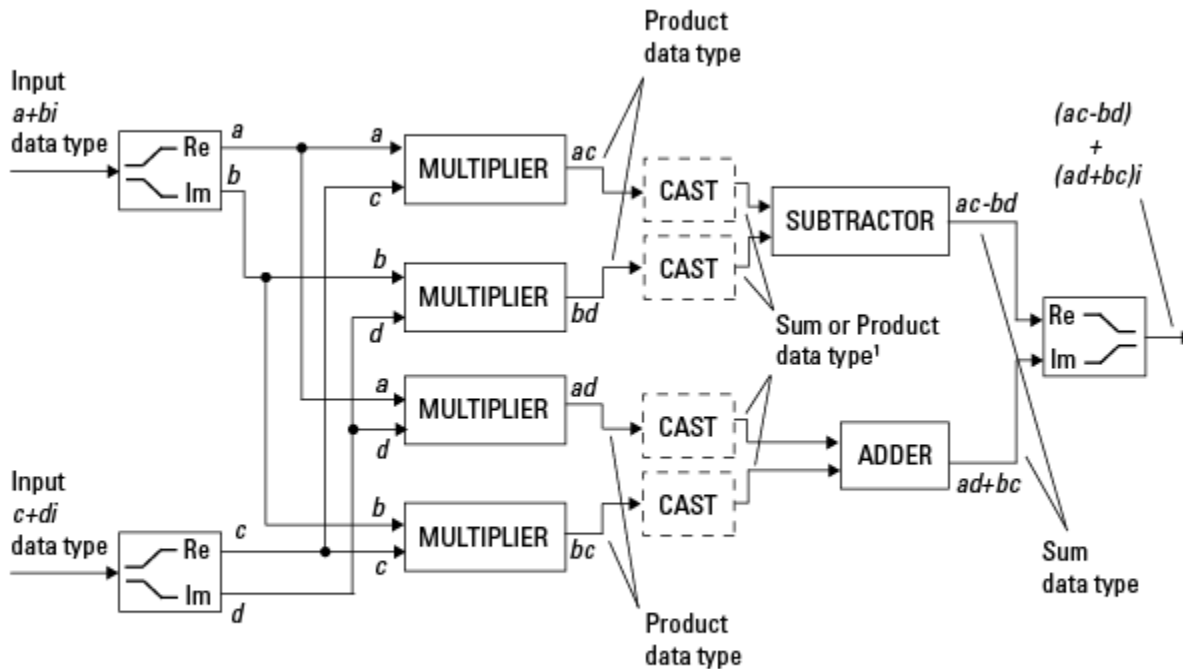
Real-Complex Multiplication

The following diagram shows the data types used by the toolbox in the multiplication of a real and a complex fixed-point number. Real-complex and complex-real multiplication are equivalent. The software returns the output of this operation in the product data type, which is governed by the `fimath` object `ProductMode` property:



Complex-Complex Multiplication

The following diagram shows the multiplication of two complex fixed-point numbers. The software returns the output of this operation in the sum data type, which is governed by the `fimath` object `SumMode` property. The intermediate product data type is determined by the `fimath` object `ProductMode` property.



¹ Sum data type if `CastBeforeSum` is true,
Product data type if `CastBeforeSum` is false

When the `fimath` object `CastBeforeSum` property is true, the casts to the sum data type are present after the multipliers in the preceding diagram. In C code, this is equivalent to

```
acc=ac;
acc-=bd;
```

for the subtractor, and

```
acc=ad;
acc+=bc;
```

for the adder, where *acc* is the accumulator. When the `CastBeforeSum` property is `false`, the casts are not present, and the data remains in the product data type before the subtraction and addition operations.

Multiplication with `fimath`

In the following examples, let

```
F = fimath('ProductMode', 'FullPrecision', ...
'SumMode', 'FullPrecision');
T1 = numerictype('WordLength', 24, 'FractionLength', 20);
T2 = numerictype('WordLength', 16, 'FractionLength', 10);
```

Real*Real

Notice that the word length and fraction length of the result *z* are equal to the sum of the word lengths and fraction lengths, respectively, of the multiplicands. This is because the `fimath` `SumMode` and `ProductMode` properties are set to `FullPrecision`:

```
P = fipref;
P.FimathDisplay = 'none';
x = fi(5, T1, F)

x =

    5

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 24
    FractionLength: 20

y = fi(10, T2, F)

y =

    10

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 10

z = x*y

z =

    50

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 40
    FractionLength: 30
```

Real*Complex

Notice that the word length and fraction length of the result z are equal to the sum of the word lengths and fraction lengths, respectively, of the multiplicands. This is because the `fimath SumMode` and `ProductMode` properties are set to `FullPrecision`:

```
x = fi(5,T1,F)
```

```
x =
```

```
5
```

```
      DataTypeMode: Fixed-point: binary point scaling  
      Signedness: Signed  
      WordLength: 24  
      FractionLength: 20
```

```
y = fi(10+2i,T2,F)
```

```
y =
```

```
10.0000 + 2.0000i
```

```
      DataTypeMode: Fixed-point: binary point scaling  
      Signedness: Signed  
      WordLength: 16  
      FractionLength: 10
```

```
z = x*y
```

```
z =
```

```
50.0000 +10.0000i
```

```
      DataTypeMode: Fixed-point: binary point scaling  
      Signedness: Signed  
      WordLength: 40  
      FractionLength: 30
```

Complex*Complex

Complex-complex multiplication involves an addition as well as multiplication. As a result, the word length of the full-precision result has one more bit than the sum of the word lengths of the multiplicands:

```
x = fi(5+6i,T1,F)
```

```
x =
```

```
5.0000 + 6.0000i
```

```
      DataTypeMode: Fixed-point: binary point scaling  
      Signedness: Signed  
      WordLength: 24  
      FractionLength: 20
```

```
y = fi(10+2i,T2,F)
```


y =

10.0000 + 2.0000i

DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 16
FractionLength: 10

z = x*y

z =

38.0000 + 70.0000i

DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 41
FractionLength: 30

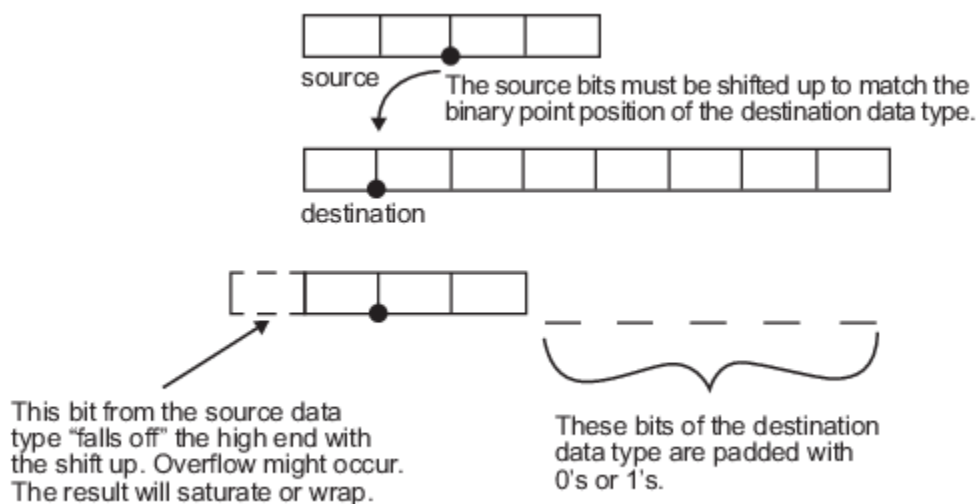
Casts

The `fimath` object allows you to specify the data type and scaling of intermediate sums and products with the `SumMode` and `ProductMode` properties. It is important to keep in mind the ramifications of each cast when you set the `SumMode` and `ProductMode` properties. Depending upon the data types you select, overflow and/or rounding might occur. The following two examples demonstrate cases where overflow and rounding can occur.

Note For more examples of casting, see “Cast fi Objects” on page 2-10.

Casting from a Shorter Data Type to a Longer Data Type

Consider the cast of a nonzero number, represented by a 4-bit data type with two fractional bits, to an 8-bit data type with seven fractional bits:



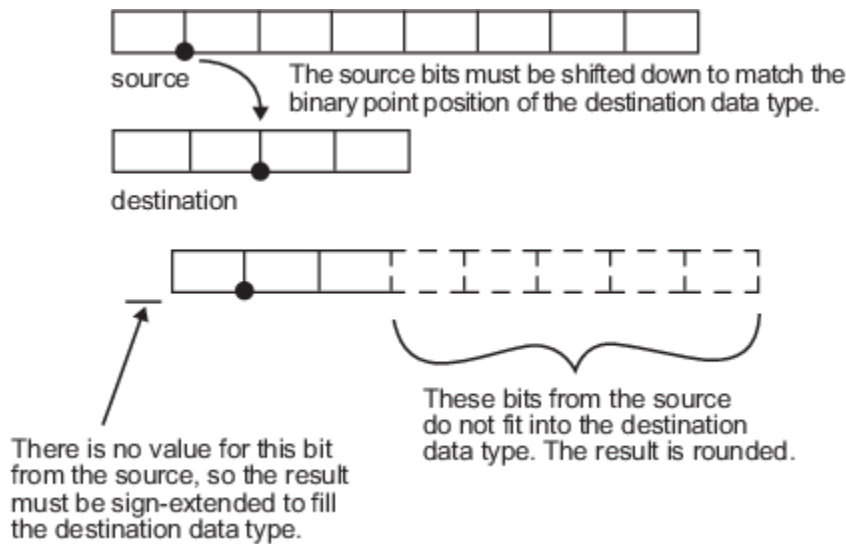
As the diagram shows, the source bits are shifted up so that the binary point matches the destination binary point position. The highest source bit does not fit, so overflow might occur and the result can saturate or wrap. The empty bits at the low end of the destination data type are padded with either 0's or 1's:

- If overflow does not occur, the empty bits are padded with 0's.
- If wrapping occurs, the empty bits are padded with 0's.
- If saturation occurs,
 - The empty bits of a positive number are padded with 1's.
 - The empty bits of a negative number are padded with 0's.

You can see that even with a cast from a shorter data type to a longer data type, overflow can still occur. This can happen when the integer length of the source data type (in this case two) is longer than the integer length of the destination data type (in this case one). Similarly, rounding might be necessary even when casting from a shorter data type to a longer data type, if the destination data type and scaling has fewer fractional bits than the source.

Casting from a Longer Data Type to a Shorter Data Type

Consider the cast of a nonzero number, represented by an 8-bit data type with seven fractional bits, to a 4-bit data type with two fractional bits:



As the diagram shows, the source bits are shifted down so that the binary point matches the destination binary point position. There is no value for the highest bit from the source, so sign extension is used to fill the integer portion of the destination data type. Sign extension is the addition of bits that have the value of the most significant bits to the high end of a two's complement number. Sign extension does not change the value of the binary number. In this example, the bottom five bits of the source do not fit into the fraction length of the destination. Therefore, precision can be lost as the result is rounded.

In this case, even though the cast is from a longer data type to a shorter data type, all the integer bits are maintained. Conversely, full precision can be maintained even if you cast to a shorter data type, as long as the fraction length of the destination data type is the same length or longer than the

fraction length of the source data type. In that case, however, bits are lost from the high end of the result and overflow can occur.

The worst case occurs when both the integer length and the fraction length of the destination data type are shorter than those of the source data type and scaling. In that case, both overflow and a loss of precision can occur.

fi Objects and C Integer Data Types

In this section...

“Integer Data Types” on page 1-18

“Unary Conversions” on page 1-19

“Binary Conversions” on page 1-20

“Overflow Handling” on page 1-21

Note The sections in this topic compare the `fi` object with fixed-point data types and operations in C. In these sections, the information on ANSI[®] C is adapted from Samuel P. Harbison and Guy L. Steele Jr., *C: A Reference Manual*, 3rd ed., Prentice Hall, 1991.

Integer Data Types

This section compares the numerical range of `fi` integer data types to the minimum numerical range of C integer data types, assuming a “Two's Complement” on page 1-9 representation.

C Integer Data Types

Many C compilers support a two's complement representation of signed integer data types. The following table shows the minimum ranges of C integer data types using a two's complement representation. The integer ranges can be larger than or equal to the ranges shown, but cannot be smaller. The range of a `long` must be larger than or equal to the range of an `int`, which must be larger than or equal to the range of a `short`.

In the two's complement representation, a signed integer with n bits has a range from -2^{n-1} to $2^{n-1} - 1$, inclusive. An unsigned integer with n bits has a range from 0 to $2^n - 1$, inclusive. The negative side of the range has one more value than the positive side, and zero is represented uniquely.

| Integer Type | Minimum | Maximum |
|----------------|----------------|---------------|
| signed char | -128 | 127 |
| unsigned char | 0 | 255 |
| short int | -32,768 | 32,767 |
| unsigned short | 0 | 65,535 |
| int | -32,768 | 32,767 |
| unsigned int | 0 | 65,535 |
| long int | -2,147,483,648 | 2,147,483,647 |
| unsigned long | 0 | 4,294,967,295 |

fi Integer Data Types

The following table lists the numerical ranges of the integer data types of the `fi` object, in particular those equivalent to the C integer data types. The ranges are large enough to accommodate the two's complement representation, which is the only signed binary encoding technique supported by Fixed-Point Designer software.

| Constructor | Signed | Word Length | Fraction Length | Minimum | Maximum | Closest ANSI C Equivalent |
|---------------------------|--------|----------------------|-----------------|----------------|---------------|---------------------------|
| <code>fi(x,1,n,0)</code> | Yes | n (2 to 65,535) | 0 | -2^{n-1} | $2^{n-1} - 1$ | Not applicable |
| <code>fi(x,0,n,0)</code> | No | n (2 to 65,535) | 0 | 0 | $2^n - 1$ | Not applicable |
| <code>fi(x,1,8,0)</code> | Yes | 8 | 0 | -128 | 127 | signed char |
| <code>fi(x,0,8,0)</code> | No | 8 | 0 | 0 | 255 | unsigned char |
| <code>fi(x,1,16,0)</code> | Yes | 16 | 0 | -32,768 | 32,767 | short int |
| <code>fi(x,0,16,0)</code> | No | 16 | 0 | 0 | 65,535 | unsigned short |
| <code>fi(x,1,32,0)</code> | Yes | 32 | 0 | -2,147,483,648 | 2,147,483,647 | long int |
| <code>fi(x,0,32,0)</code> | No | 32 | 0 | 0 | 4,294,967,295 | unsigned long |

Unary Conversions

Unary conversions dictate whether and how a single operand is converted before an operation is performed. This section discusses unary conversions in ANSI C and of `fi` objects.

ANSI C Usual Unary Conversions

Unary conversions in ANSI C are automatically applied to the operands of the unary `!`, `-`, `~`, and `*` operators, and of the binary `<<` and `>>` operators, according to the following table:

| Original Operand Type | ANSI C Conversion |
|---------------------------------|----------------------------------|
| char or short | int |
| unsigned char or unsigned short | int or unsigned int ¹ |
| float | float |
| Array of T | Pointer to T |
| Function returning T | Pointer to function returning T |

¹If type `int` cannot represent all the values of the original data type without overflow, the converted type is `unsigned int`.

fi Usual Unary Conversions

The following table shows the `fi` unary conversions:

| C Operator | fi Equivalent | fi Conversion |
|-----------------|--------------------------|---|
| <code>!x</code> | <code>~x = not(x)</code> | Result is logical. |
| <code>~x</code> | <code>bitcmp(x)</code> | Result is same numeric type as operand. |
| <code>*x</code> | No equivalent | Not applicable |

| C Operator | fi Equivalent | fi Conversion |
|-------------------------|--|--|
| <code>x<<n</code> | <code>bitshift(x,n)</code> positive n | Result is same numeric type as operand. Round mode is always floor. Overflow mode is obeyed. 0-valued bits are shifted in on the right. |
| <code>x>>n</code> | <code>bitshift(x,-n)</code> | Result is same numeric type as operand. Round mode is always floor. Overflow mode is obeyed. 0-valued bits are shifted in on the left if the operand is unsigned or signed and positive. 1-valued bits are shifted in on the left if the operand is signed and negative. |
| <code>+x</code> | <code>+x</code> | Result is same numeric type as operand. |
| <code>-x</code> | <code>-x</code> | Result is same numeric type as operand. Overflow mode is obeyed. For example, overflow might occur when you negate an unsigned <code>fi</code> or the most negative value of a signed <code>fi</code> . |

Binary Conversions

This section describes the conversions that occur when the operands of a binary operator are different data types.

ANSI C Usual Binary Conversions

In ANSI C, operands of a binary operator must be of the same type. If they are different, one is converted to the type of the other according to the first applicable conversion in the following table:

| Type of One Operand | Type of Other Operand | ANSI C Conversion |
|---------------------|-----------------------|------------------------------------|
| long double | Any | long double |
| double | Any | double |
| float | Any | float |
| unsigned long | Any | unsigned long |
| long | unsigned | long or unsigned long ¹ |
| long | int | long |
| unsigned | int or unsigned | unsigned |
| int | int | int |

¹Type long is only used if it can represent all values of type unsigned.

fi Usual Binary Conversions

When one of the operands of a binary operator (+, -, *, .*) is a `fi` object and the other is a MATLAB built-in numeric type, then the non-`fi` operand is converted to a `fi` object before the operation is performed, according to the following table:

| Type of One Operand | Type of Other Operand | Properties of Other Operand After Conversion to a fi Object |
|---------------------|-----------------------|---|
| fi | double or single | <ul style="list-style-type: none"> • Signed = same as the original fi operand • WordLength = same as the original fi operand • FractionLength = set to best precision possible |
| fi | int8 | <ul style="list-style-type: none"> • Signed = 1 • WordLength = 8 • FractionLength = 0 |
| fi | uint8 | <ul style="list-style-type: none"> • Signed = 0 • WordLength = 8 • FractionLength = 0 |
| fi | int16 | <ul style="list-style-type: none"> • Signed = 1 • WordLength = 16 • FractionLength = 0 |
| fi | uint16 | <ul style="list-style-type: none"> • Signed = 0 • WordLength = 16 • FractionLength = 0 |
| fi | int32 | <ul style="list-style-type: none"> • Signed = 1 • WordLength = 32 • FractionLength = 0 |
| fi | uint32 | <ul style="list-style-type: none"> • Signed = 0 • WordLength = 32 • FractionLength = 0 |
| fi | int64 | <ul style="list-style-type: none"> • Signed = 1 • WordLength = 64 • FractionLength = 0 |
| fi | uint64 | <ul style="list-style-type: none"> • Signed = 0 • WordLength = 64 • FractionLength = 0 |

Overflow Handling

The following sections compare how ANSI C and Fixed-Point Designer software handle overflows.

ANSI C Overflow Handling

In ANSI C, the result of signed integer operations is whatever value is produced by the machine instruction used to implement the operation. Therefore, ANSI C has no rules for handling signed integer overflow.

The results of unsigned integer overflows wrap in ANSI C.

fi Overflow Handling

Addition and multiplication with `fi` objects yield results that can be exactly represented by a `fi` object, up to word lengths of 65,535 bits or the available memory on your machine. This is not true of division, however, because many ratios result in infinite binary expressions. You can perform division with `fi` objects using the `divide` function, which requires you to explicitly specify the numeric type of the result.

The conditions under which a `fi` object overflows and the results then produced are determined by the associated `fimath` object. You can specify certain overflow characteristics separately for sums (including differences) and products. Refer to the following table:

| fimath Object Properties Related to Overflow Handling | Property Value | Description |
|--|-----------------------|--|
| OverflowAction | 'saturate' | Overflows are saturated to the maximum or minimum value in the range. |
| | 'wrap' | Overflows wrap using modulo arithmetic if unsigned, two's complement wrap if signed. |
| ProductMode | 'FullPrecision' | Full-precision results are kept. Overflow does not occur. An error is thrown if the resulting word length is greater than <code>MaxProductWordLength</code> . The rules for computing the resulting product word and fraction lengths are given in "fimath Object Properties" on page 3-4 in the Property Reference. |
| | 'KeepLSB' | The least significant bits of the product are kept. Full precision is kept, but overflow is possible. This behavior models the C language integer operations. The <code>ProductWordLength</code> property determines the resulting word length. If <code>ProductWordLength</code> is greater than is necessary for the full-precision product, then the result is stored in the least significant bits. If <code>ProductWordLength</code> is less than is necessary for the full-precision product, then overflow occurs. The rule for computing the resulting product fraction length is given in "fimath Object Properties" on page 3-4 in the Property Reference. |

| fimath Object Properties Related to Overflow Handling | Property Value | Description |
|---|--------------------|--|
| | 'KeepMSB' | <p>The most significant bits of the product are kept. Overflow is prevented, but precision may be lost.</p> <p>The ProductWordLength property determines the resulting word length. If ProductWordLength is greater than is necessary for the full-precision product, then the result is stored in the most significant bits. If ProductWordLength is less than is necessary for the full-precision product, then rounding occurs.</p> <p>The rule for computing the resulting product fraction length is given in “fimath Object Properties” on page 3-4 in the Property Reference.</p> |
| | 'SpecifyPrecision' | You can specify both the word length and the fraction length of the resulting product. |
| ProductWordLength | Positive integer | The word length of product results when ProductMode is 'KeepLSB', 'KeepMSB', or 'SpecifyPrecision'. |
| MaxProductWordLength | Positive integer | The maximum product word length allowed when ProductMode is 'FullPrecision'. The default is 65,535 bits. This property can help ensure that your simulation does not exceed your hardware requirements. |
| ProductFractionLength | Integer | The fraction length of product results when ProductMode is 'Specify Precision'. |
| SumMode | 'FullPrecision' | <p>Full-precision results are kept. Overflow does not occur. An error is thrown if the resulting word length is greater than MaxSumWordLength.</p> <p>The rules for computing the resulting sum word and fraction lengths are given in “fimath Object Properties” on page 3-4 in the Property Reference.</p> |

| fimath Object Properties Related to Overflow Handling | Property Value | Description |
|---|--------------------|---|
| | 'KeepLSB' | <p>The least significant bits of the sum are kept. Full precision is kept, but overflow is possible. This behavior models the C language integer operations.</p> <p>The <code>SumWordLength</code> property determines the resulting word length. If <code>SumWordLength</code> is greater than is necessary for the full-precision sum, then the result is stored in the least significant bits. If <code>SumWordLength</code> is less than is necessary for the full-precision sum, then overflow occurs.</p> <p>The rule for computing the resulting sum fraction length is given in “fimath Object Properties” on page 3-4 in the Property Reference.</p> |
| | 'KeepMSB' | <p>The most significant bits of the sum are kept. Overflow is prevented, but precision may be lost.</p> <p>The <code>SumWordLength</code> property determines the resulting word length. If <code>SumWordLength</code> is greater than is necessary for the full-precision sum, then the result is stored in the most significant bits. If <code>SumWordLength</code> is less than is necessary for the full-precision sum, then rounding occurs.</p> <p>The rule for computing the resulting sum fraction length is given in “fimath Object Properties” on page 3-4 in the Property Reference.</p> |
| | 'SpecifyPrecision' | <p>You can specify both the word length and the fraction length of the resulting sum.</p> |
| SumWordLength | Positive integer | <p>The word length of sum results when <code>SumMode</code> is 'KeepLSB', 'KeepMSB', or 'SpecifyPrecision'.</p> |
| MaxSumWordLength | Positive integer | <p>The maximum sum word length allowed when <code>SumMode</code> is 'FullPrecision'. The default is 65,535 bits. This property can help ensure that your simulation does not exceed your hardware requirements.</p> |
| SumFractionLength | Integer | <p>The fraction length of sum results when <code>SumMode</code> is 'SpecifyPrecision'.</p> |

Working with fi Objects

- “Ways to Construct fi Objects” on page 2-2
- “Cast fi Objects” on page 2-10
- “Set fi Object Properties” on page 2-15

Ways to Construct fi Objects

In this section...

“Use fi Constructor to Create fi Objects” on page 2-2

“Use the Insert fi Constructor Dialog to Build fi Object Constructors” on page 2-6

“Determine Property Precedence” on page 2-7

“Create fi Objects For Use in a Types Table” on page 2-8

You can create a `fi` object by using a `fi` constructor function or you can build `fi` object constructors using the **Insert fi Constructor** dialog. You can also use a `fi` constructor function to copy an existing `fi` object.

`numericType` and `fiMath` properties can be specified directly in the `fi` constructor function or you can use an existing `numericType` or `fiMath` object to construct a `fi` object. The value of a property is taken from the last time it is set.

You can write a reusable MATLAB algorithm by keeping the data types of the algorithmic variables in a separate types table.

Use fi Constructor to Create fi Objects

These examples show you several different ways to construct `fi` objects.

Construct fi Object with Default Data Type and Property Values

Create a `fi` object with the default data type and a value of 0.

```
a = fi(0)
```

```
a =
```

```
0
```

```

      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 16
      FractionLength: 15

```

The default `fi` constructor syntax creates a signed `fi` object with a value of 0, word length of 16 bits, and fraction length of 15 bits.

Note The `fi` constructor creates the `fi` object using a `RoundingMethod` of `Nearest` and an `OverflowAction` of `Saturate`. If you construct a `fi` from floating-point values, the default `RoundingMethod` and `OverflowAction` property settings are not used.

For information on the display format of `fi` objects, refer to “View Fixed-Point Data”.

Copy a fi Object

To copy a `fi` object, use assignment.

```

a = fi(pi)
a =
    3.1416

        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 16
        FractionLength: 13

b = a
b =
    3.1416

        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 16
        FractionLength: 13

```

Construct fi Object with Property Name-Value Pair Arguments

You can use property name-value pair arguments to set `fi` object properties in the `fi` constructor. The `fi` object has three types of properties:

- `fi` Object Data Properties
- `fimath` Object Properties
- `numericType` Object Properties

For example, specify the `fimath` object properties for the rounding method and overflow action to use when performing fixed-point arithmetic.

```

a = fi(pi, 'RoundingMethod', 'Floor', ...
        'OverflowAction', 'Wrap')
a =
    3.1415

        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 16
        FractionLength: 13

        RoundingMethod: Floor
        OverflowAction: Wrap
        ProductMode: FullPrecision
        SumMode: FullPrecision

```

If you specify at least one `fimath` object property in the `fi` constructor, the `fi` object has a local `fimath` object. The `fi` object uses default values for the remaining unspecified `fimath` object properties.

If you do not specify any `fimath` object properties in the `fi` object constructor, the `fi` object uses default `fimath` values and has no local `fimath`.

You can use the `isfimathlocal` function to determine whether a `fi` object has a local `fimath` associated with it.

Construct fi Object Using numericType Object

You can create a `fi` object using a `numericType` object. The “numericType Object Properties” define the data type and scaling attributes of a `fi` object.

Create a `numericType` object with default property values.

```
T = numericType
```

```
T =
```

```
      DataTypeMode: Fixed-point: binary point scaling  
      Signedness: Signed  
      WordLength: 16  
      FractionLength: 15
```

Create a `fi` object from the `numericType` object `T`.

```
a = fi(pi,T)
```

```
a =
```

```
1.0000
```

```
      DataTypeMode: Fixed-point: binary point scaling  
      Signedness: Signed  
      WordLength: 16  
      FractionLength: 15
```

You can use a `fimath` object and a `numericType` object in the `fi` constructor.

```
F = fimath('RoundingMethod','Nearest',...  
          'OverflowAction','Saturate',...  
          'ProductMode','FullPrecision',...  
          'SumMode','FullPrecision')
```

```
F =
```

```
      RoundingMethod: Nearest  
      OverflowAction: Saturate  
      ProductMode: FullPrecision  
      SumMode: FullPrecision
```

```
a = fi(pi,T,F)
```

```
a =
```

```
1.0000
```

```
      DataTypeMode: Fixed-point: binary point scaling  
      Signedness: Signed  
      WordLength: 16  
      FractionLength: 15
```

```

RoundingMethod: Nearest
OverflowAction: Saturate
ProductMode: FullPrecision
SumMode: FullPrecision

```

Note The syntax `a = fi(pi,T,F)` is equivalent to `a = fi(pi,F,T)`. You can use both statements to define a `fi` object using a `fimath` object and a `numericType` object.

Construct fi Object Using fimath Object

You can create a `fi` object using a specific `fimath` object. When you do so, a local `fimath` object is assigned to the `fi` object you create. If you do not specify any `numericType` object properties, the word length of the `fi` object defaults to 16 bits. The fraction length is determined by best precision scaling.

For example, create a `fimath` object that specifies the rounding method, overflow action, product mode, and sum mode to use.

```

F = fimath('RoundingMethod','Nearest',...
          'OverflowAction','Saturate',...
          'ProductMode','FullPrecision',...
          'SumMode','FullPrecision')

```

F =

```

          RoundingMethod: Nearest
OverflowAction: Saturate
          ProductMode: FullPrecision
          SumMode: FullPrecision

```

Use dot notation to change the overflow action of the `fimath` object `F`.

```

F.OverflowAction = 'Wrap'

```

F =

```

          RoundingMethod: Nearest
OverflowAction: Wrap
          ProductMode: FullPrecision
          SumMode: FullPrecision

```

Create a `fi` object using the `fimath` object `F`.

```

a = fi(pi,F)

```

a =

```

3.1416

```

```

          DataTypeMode: Fixed-point: binary point scaling
          Signedness: Signed
          WordLength: 16
          FractionLength: 13

```

```
RoundingMethod: Nearest
OverflowAction: Wrap
ProductMode: FullPrecision
SumMode: FullPrecision
```

You can also create `fi` objects using a `fimath` object while specifying various `numericType` properties at creation time. For example, create an unsigned `fi` object with a value of `pi`, word length of 8 bits, fraction length of 6 bits, and `fimath F`.

```
b = fi(pi,0,8,6,F)
```

```
b =
```

```
3.1406
```


```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Unsigned
WordLength: 8
FractionLength: 6
```

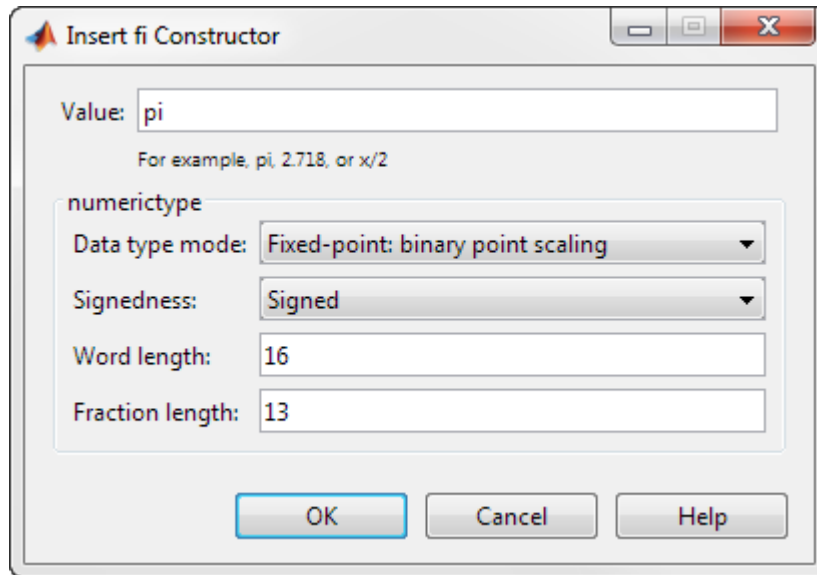
```
RoundingMethod: Nearest
OverflowAction: Wrap
ProductMode: FullPrecision
SumMode: FullPrecision
```

Use the Insert fi Constructor Dialog to Build fi Object Constructors

You can build `fi` object constructors in MATLAB by using the **Insert fi Constructor** dialog box. After specifying the value and properties of the `fi` object in the dialog box, you can insert the prepopulated `fi` object constructor at a specific location in your file.

For example, create a signed `fi` object with a value of `pi`, a word length of 16 bits and a fraction length of 13 bits.

- 1 On the MATLAB **Home** tab, in the **File** section, click **New Script**.
- 2 On the **Editor** tab, in the **Code** section, click the **Specify fixed-point data button**  button arrow. Click **Insert fi** to open the **Insert fi Constructor** dialog box.
- 3 Use the edit boxes and drop-down menus to specify the following properties of the `fi` object:
 - **Value** = `pi`
 - **Data type mode** = Fixed-point: binary point scaling
 - **Signedness** = Signed
 - **Word length** = 16
 - **Fraction length** = 13



- 4 To insert the `fi` object constructor in your file, place your cursor at the desired location in the file, then click **OK** on the **Insert fi Constructor** dialog box. Clicking **OK** closes the **Insert fi Constructor** dialog box and automatically populates the `fi` object constructor in your file.

```
fi(pi, 1, 16, 13)
```

Determine Property Precedence

The value of a property of a `fi` object is taken from the last time it is set. For example, create a `numericity` object with the `Signed` set to `true` and a fraction length of 14.

```
T = numericity('Signed',true,...
    'FractionLength',14)
```

```
T =
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 14
```

Create a `fi` object that specifies the `numericity` property `T` *after* the `Signed` property. The resulting `fi` object is signed.

```
a = fi(pi, 'Signed', false, ...
    'numericity', T)
```

```
a =
```

```
1.9999
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 14
```

Create a second `fi` object that specifies the `numericType` `T` *before* the `Signed` property. The resulting `fi` object is unsigned.

```
b = fi(pi, 'numericType', T, ...
      'Signed', false)

b =

    3.1416

      DataTypeMode: Fixed-point: binary point scaling
      Signedness:   Unsigned
      WordLength:   16
      FractionLength: 14
```

Create fi Objects For Use in a Types Table

You can write a reusable MATLAB algorithm by keeping the data types of the algorithmic variables in a separate types table. For example,

```
function T = mytypes(dt)
    switch dt
        case 'double'
            T.b = double([]);
            T.x = double([]);
            T.y = double([]);

        case 'fixed16'
            T.b = fi([], 1, 16, 15);
            T.x = fi([], 1, 16, 15);
            T.y = fi([], 1, 16, 14);
    end
end
```

Cast the variables in the algorithm to the data types in the types table as described in “Manual Fixed-Point Conversion Best Practices” on page 11-3.

```
function [y,z]=myfilter(b,x,z,T)
    y = zeros(size(x), 'like', T.y);
    for n=1:length(x)
        z(:) = [x(n); z(1:end-1)];
        y(n) = b * z;
    end
end
```

In a separate test file, set up input data to feed into your algorithm, and specify the data types of the inputs.

```
% Test inputs
b = fir1(11, 0.25);
t = linspace(0, 10*pi, 256)';
x = sin((pi/16)*t.^2);
% Linear chirp

% Cast inputs
T=mytypes('fixed16');
b=cast(b, 'like', T.b);
```

```
x=cast(x,'like',T.x);  
z=zeros(size(b),'like',T.x);
```

```
% Run  
[y,z] = myfilter(b,x,z,T);
```

See Also

fi | fimath | fipref | numerictype

Related Examples

- “View Fixed-Point Data”
- “Cast fi Objects” on page 2-10
- “Manual Fixed-Point Conversion Best Practices” on page 11-3

Cast fi Objects

In this section...

“Overwriting by Assignment” on page 2-10

“Ways to Cast with MATLAB Software” on page 2-10

Overwriting by Assignment

Because MATLAB software does not have type declarations, an assignment like `A = B` replaces the type and content of `A` with the type and content of `B`. If `A` does not exist at the time of the assignment, MATLAB creates the variable `A` and assigns it the same type and value as `B`. Such assignment happens with all types in MATLAB — objects and built-in types alike — including `fi`, `double`, `single`, `int8`, `uint8`, `int16`, etc.

For example, the following code overwrites the value and `int8` type of `A` with the value and `int16` type of `B`:

```
A = int8(0);
B = int16(32767);
A = B
```

```
A =
```

```
    int16
```

```
    32767
```

```
class(A)
```

```
ans =
```

```
    'int16'
```

Ways to Cast with MATLAB Software

You may find it useful to cast data into another type—for example, when you are casting data from an accumulator to memory. There are several ways to cast data in MATLAB. The following sections provide examples of four different methods:

- Casting by Subscripted Assignment
- Casting by Conversion Function
- Casting with the Fixed-Point Designer `reinterpretcast` function
- Casting with the `cast` Function

Casting by Subscripted Assignment

The following subscripted assignment statement retains the type of `A` and saturates the value of `B` to an `int8`:

```
A = int8(0);
B = int16(32767);
A(:) = B
```

```
A =
    int8
    127
class(A)
ans =
    'int8'
```

The same is true for `fi` objects:

```
fipref('NumericTypeDisplay', 'short');
A = fi(0, 1, 8, 0);
B = fi(32767, 1, 16, 0);
A(:) = B
A =
    127
    numerictype(1,8,0)
```

Note For more information on subscripted assignments, see the `subsasgn` function.

Casting by Conversion Function

You can convert from one data type to another by using a conversion function. In this example, `A` does not have to be predefined because it is overwritten.

```
B = int16(32767);
A = int8(B)
A =
    int8
    127
class(A)
ans =
    'int8'
```

The same is true for `fi` objects:

```
B = fi(32767,1,16,0)
A = fi(B,1,8,0)
B =
    32767
    numerictype(1,16,0)
A =
```

```
127
    numerictype(1,8,0)
```

Using a `numerictype` Object in the `fi` Conversion Function

Often a specific `numerictype` is used in many places, and it is convenient to predefine `numerictype` objects for use in the conversion functions. Predefining these objects is a good practice because it also puts the data type specification in one place.

```
T8 = numerictype(1,8,0)
T16 = numerictype(1,16,0)
```

```
T8 =
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 8
    FractionLength: 0
```

```
T16 =
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 0
```

```
B = fi(32767,T16)
A = fi(B,T8)
```

```
B =
```

```
    32767
    numerictype(1,16,0)
```

```
A =
```

```
127
    numerictype(1,8,0)
```

Casting with the `reinterpretpcast` Function

You can convert fixed-point and built-in data types without changing the underlying data. The Fixed-Point Designer `reinterpretpcast` function performs this type of conversion.

In the following example, `B` is an unsigned `fi` object with a word length of 8 bits and a fraction length of 5 bits. The `reinterpretpcast` function converts `B` into a signed `fi` object `A` with a word length of 8 bits and a fraction length of 1 bit. The real-world values of `A` and `B` differ, but their binary representations are the same.

```
B = fi([pi/4 1 pi/2 4],0,8,5)
T = numerictype(1,8,1);
A = reinterpretpcast(B,T)
```

```
B =
```

```

0.7813    1.0000    1.5625    4.0000
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Unsigned
      WordLength: 8
      FractionLength: 5

```

A =

```

12.5000    16.0000    25.0000   -64.0000
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 8
      FractionLength: 1

```

To verify that the underlying data has not changed, compare the binary representations of A and B:

```

binary_B = bin(B)
binary_A = bin(A)

binary_B =
      '00011001    00100000    00110010    10000000'

binary_A =
      '00011001    00100000    00110010    10000000'

```

Casting with the cast Function

Using the cast function, you can convert the value of a variable to the same numerictype, complexity, and fimath as another variable.

In the following example, a is cast to the data type of b. The output, c, has the same numerictype and fimath properties as b, and the value of a.

```

a = pi;
b = fi([],1,16,13,'RoundingMethod','Floor');
c = cast(a,'like',b)

c =
      3.1415
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 16
      FractionLength: 13
      RoundingMethod: Floor
      OverflowAction: Saturate
      ProductMode: FullPrecision
      SumMode: FullPrecision

```

Using this syntax allows you to specify data types separately from your algorithmic code as described in “Manual Fixed-Point Conversion Best Practices” on page 11-3.

See Also

cast | fi | numerictype | reinterpretcast | subsasgn

Set fi Object Properties

In this section...

“Set Fixed-Point Properties at Object Creation” on page 2-15
 “Use Subscripted Assignment to Set Real-World Value of fi Object” on page 2-15
 “Direct Property Referencing to Read fi Object Properties” on page 2-16
 “Best Practices for Code Generation” on page 2-16
 “Remove Local fimath Properties from fi Object” on page 2-18

Set fi object properties when you create the fi object.

Set Fixed-Point Properties at Object Creation

You can set properties of fi objects at the time of object creation by including properties after the arguments of the fi constructor function. For example, to set the overflow action to Wrap and the rounding method to Convergent in the fi constructor.

```
a = fi(pi, 'OverflowAction', 'Wrap', ...
      'RoundingMethod', 'Convergent')
a =
    3.1416

      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 16
      FractionLength: 13

      RoundingMethod: Convergent
      OverflowAction: Wrap
      ProductMode: FullPrecision
      SumMode: FullPrecision
```

To set the stored integer value of a fi object, use the parameter name-value pair arguments for the 'int' property when you create the object. For example, create a signed fi object with a stored integer value of 4, 16-bit word length, and 15-bit fraction length.

```
x = fi(0,1,16,15, 'int',4);
```

Verify that the fi object has the expected integer setting.

```
x.int
ans =
    int16
     4
```

Use Subscripted Assignment to Set Real-World Value of fi Object

You can set the real-world value of a fi object via subscripted assignment.

```
a = fi(pi);
a(:) = 2

a =

    2

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 13
```

Direct Property Referencing to Read fi Object Properties

You can read `fi` object property values using MATLAB structure-like referencing. You do this by using a period to index into a property by name.

For example, to get the `WordLength` of `a`.

```
a.WordLength

ans =

    16
```

Best Practices for Code Generation

The following methods for setting `fi` object properties are recommended for compatibility with code generation.

First, define the `fi` object `a`.

```
a = fi(pi, 'OverflowAction', 'Wrap', ...
    'RoundingMethod', 'Convergent')

a =

    3.1416

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 13

    RoundingMethod: Convergent
    OverflowAction: Wrap
    ProductMode: FullPrecision
    SumMode: FullPrecision
```

You can get the `fimath` using dot indexing, use the `fimath` constructor to change the `fimath` settings, then use `setfimath` to set the local `fimath` object back into `fi` object `a`.

```
F = fimath(a.fimath, 'OverflowAction', 'Saturate');
a = setfimath(a,F)

a =
```

```
3.1416
```

```

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 13

    RoundingMethod: Convergent
    OverflowAction: Saturate
    ProductMode: FullPrecision
    SumMode: FullPrecision

```

The `setfimath` function is useful for changing out the `fimath` altogether. For example:

```

a = fi(pi);
F = fixed.fimathLike(a);
a = setfimath(a,F)

```

```
a =
```

```
3.1416
```

```

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 13

    RoundingMethod: Floor
    OverflowAction: Wrap
    ProductMode: SpecifyPrecision
    ProductWordLength: 16
    ProductFractionLength: 13
    SumMode: SpecifyPrecision
    SumWordLength: 16
    SumFractionLength: 13
    CastBeforeSum: true

```

Alternatively, you can call the `fi` object constructor with the value input set to the original `fi` object, then add new `fimath` parameters directly in the `fi` object constructor. For example:

```

a = fi(pi,'OverflowAction','Wrap',...
    'RoundingMethod','Convergent')

```

```
a =
```

```
3.1416
```

```

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 13

    RoundingMethod: Convergent
    OverflowAction: Wrap
    ProductMode: FullPrecision
    SumMode: FullPrecision

```

```
a = fi(a,'OverflowAction','Saturate')
```

```
a =  
3.1416  
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 16  
    FractionLength: 13  
  
    RoundingMethod: Convergent  
    OverflowAction: Saturate  
    ProductMode: FullPrecision  
    SumMode: FullPrecision
```

Note that using dot indexing to write `fi`math and `numeric`type object properties is *not* compatible with code generation. For example:

```
a.OverflowAction = 'Saturate' % Works in interpreted MATLAB only
```

Remove Local `fi`math Properties from `fi` Object

If you have a `fi` object `b` with a local `fi`math object, you can remove the local `fi`math object and force `b` to use default `fi`math values.

```
b = fi(pi,1,'RoundingMethod','Floor')  
b =  
3.1415  
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 16  
    FractionLength: 13  
  
    RoundingMethod: Floor  
    OverflowAction: Saturate  
    ProductMode: FullPrecision  
    SumMode: FullPrecision  
  
b = removefimath(b)  
b =  
3.1415  
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 16  
    FractionLength: 13  
  
isfimathlocal(b)  
ans =  
logical  
0
```

See Also

[fi](#) | [fimath](#) | [numerictype](#) | [fi Object Data Properties](#) | [fimath Object Properties](#) | [numerictype Object Properties](#)

Working with fimath Objects

- “fimath Object Construction” on page 3-2
- “fimath Object Properties” on page 3-4
- “fimath Properties Usage for Fixed-Point Arithmetic” on page 3-10
- “fimath for Rounding and Overflow Modes” on page 3-16
- “fimath for Sharing Arithmetic Rules” on page 3-17
- “fimath ProductMode and SumMode” on page 3-19
- “How Functions Use fimath” on page 3-24

fimath Object Construction

In this section...

“fimath Object Syntaxes” on page 3-2

“Building fimath Object Constructors in a GUI” on page 3-3

fimath Object Syntaxes

The arithmetic attributes of a `fi` object are defined by a local `fimath` object, which is attached to that `fi` object. If a `fi` object has no local `fimath`, the following default `fimath` values are used:

```
RoundingMethod: Nearest
OverflowAction: Wrap
ProductMode: FullPrecision
SumMode: FullPrecision
```

You can create `fimath` objects in Fixed-Point Designer software in one of two ways:

- You can use the `fimath` constructor function to create new `fimath` objects.
- You can use the `fimath` constructor function to copy an existing `fimath` object.

To get started, type

```
F = fimath
```

to create a `fimath` object.

```
F =
```

```
RoundingMethod: Nearest
OverflowAction: Saturate
ProductMode: FullPrecision
SumMode: FullPrecision
```

To copy a `fimath` object, simply use assignment as in the following example:

```
F = fimath;
G = F;
isequal(F,G)
```

```
ans =
```

```
logical
```

```
1
```

The syntax


```
F = fimath(...'PropertyName',PropertyValue...)
```

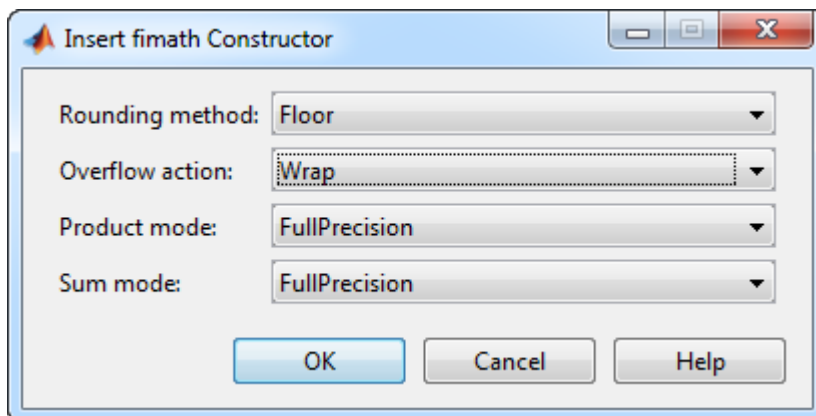
allows you to set properties for a `fimath` object at object creation with property name/property value pairs. Refer to “Setting `fimath` Properties at Object Creation” on page 3-8.

Building fimath Object Constructors in a GUI

When you are working with files in MATLAB, you can build your `fimath` object constructors using the **Insert fimath Constructor** dialog box. After specifying the properties of the `fimath` object in the dialog box, you can insert the prepopulated `fimath` object constructor at a specific location in your file.

For example, to create a `fimath` object that uses convergent rounding and wraps on overflow, perform the following steps:

- 1 On the **Home** tab, in the **File** section, click **New > Script** to open the MATLAB Editor
- 2 On the **Editor** tab, in the **Edit** section, click  in the **Insert** button group. Click the **Insert fimath...** to open the **Insert fimath Constructor** dialog box.
- 3 Use the edit boxes and drop-down menus to specify the following properties of the `fimath` object:
 - **Rounding method** = Floor
 - **Overflow action** = Wrap
 - **Product mode** = FullPrecision
 - **Sum mode** = FullPrecision



- 4 To insert the `fimath` object constructor in your file, place your cursor at the desired location in the file. Then click **OK** on the **Insert fimath Constructor** dialog box. Clicking **OK** closes the **Insert fimath Constructor** dialog box and automatically populates the `fimath` object constructor in your file:

```

1  fimath('RoundingMethod', 'Floor', ...
2      'OverflowAction', 'Wrap', ...
3      'ProductMode', 'FullPrecision', ...
4      'SumMode', 'FullPrecision')
```

fimath Object Properties

In this section...

“Math, Rounding, and Overflow Properties” on page 3-4

“How Properties are Related” on page 3-7

“Setting fimath Object Properties” on page 3-8

Math, Rounding, and Overflow Properties

You can always write to the following properties of fimath objects:

| Property | Description | Valid Values |
|-----------------------|---|--|
| CastBeforeSum | Whether both operands are cast to the sum data type before addition | <ul style="list-style-type: none"> 0 (default) — do not cast before sum 1 — cast before sum <p>Note This property is hidden when the SumMode is set to FullPrecision.</p> |
| MaxProductWordLength | Maximum allowable word length for the product data type | <ul style="list-style-type: none"> 65535 (default) Any positive integer |
| MaxSumWordLength | Maximum allowable word length for the sum data type | <ul style="list-style-type: none"> 65535 (default) Any positive integer |
| OverflowAction | Action to take on overflow | <ul style="list-style-type: none"> Saturate (default) — Saturate to maximum or minimum value of the fixed-point range on overflow. Wrap — Wrap on overflow. This mode is also known as two's complement overflow. |
| ProductBias | Bias of the product data type | <ul style="list-style-type: none"> 0 (default) Any floating-point number |
| ProductFixedExponent | Fixed exponent of the product data type | <ul style="list-style-type: none"> -30 (default) Any positive or negative integer <p>Note The ProductFractionLength is the negative of the ProductFixedExponent. Changing one property changes the other.</p> |
| ProductFractionLength | Fraction length, in bits, of the product data type | <ul style="list-style-type: none"> 30 (default) Any positive or negative integer <p>Note The ProductFractionLength is the negative of the ProductFixedExponent. Changing one property changes the other.</p> |

| Property | Description | Valid Values |
|------------------------------|--|---|
| ProductMode | Defines how the product data type is determined | <ul style="list-style-type: none"> • FullPrecision (default) — The full precision of the result is kept. • KeepLSB— Keep least significant bits. Specify the product word length, while the fraction length is set to maintain the least significant bits of the product. • KeepMSB — Keep most significant bits. Specify the product word length, while the fraction length is set to maintain the most significant bits of the product. • SpecifyPrecision— specify the word and fraction lengths or slope and bias of the product. |
| ProductSlope | Slope of the product data type | <ul style="list-style-type: none"> • 9.3132e-010 (default) • Any floating-point number <p>Note</p> $ProductSlope = ProductSlopeAdjustmentFactor \times 2^{ProductFixedExponent}$ <p>Changing one of these properties affects the others.</p> |
| ProductSlopeAdjustmentFactor | Slope adjustment factor of the product data type | <ul style="list-style-type: none"> • 1 (default) • Any floating-point number greater than or equal to 1 and less than 2 <p>Note</p> $ProductSlope = ProductSlopeAdjustmentFactor \times 2^{ProductFixedExponent}$ <p>Changing one of these properties affects the others.</p> |
| ProductWordLength | Word length, in bits, of the product data type | <ul style="list-style-type: none"> • 32 (default) • Any positive integer |
| RoundingMethod | Rounding method | <ul style="list-style-type: none"> • Nearest (default) — Round toward nearest. Ties round toward positive infinity. • Ceiling — Round toward positive infinity. • Convergent — Round toward nearest. Ties round to the nearest even stored integer (least biased). • Zero — Round toward zero. • Floor — Round toward negative infinity. • Round — Round toward nearest. Ties round toward negative infinity for negative numbers, and toward positive infinity for positive numbers. |

| Property | Description | Valid Values |
|--------------------------|--|---|
| SumBias | Bias of the sum data type | <ul style="list-style-type: none"> • 0 (default) • Any floating-point number |
| SumFixedExponent | Fixed exponent of the sum data type | <ul style="list-style-type: none"> • -30 (default) • Any positive or negative integer <p>Note The SumFractionLength is the negative of the SumFixedExponent. Changing one property changes the other.</p> |
| SumFractionLength | Fraction length, in bits, of the sum data type | <ul style="list-style-type: none"> • 30 (default) • Any positive or negative integer <p>Note The SumFractionLength is the negative of the SumFixedExponent. Changing one property changes the other.</p> |
| SumMode | Defines how the sum data type is determined | <ul style="list-style-type: none"> • FullPrecision (default) — The full precision of the result is kept. • KeepLSB — Keep least significant bits. Specify the sum data type word length, while the fraction length is set to maintain the least significant bits of the sum. • KeepMSB — Keep most significant bits. Specify the sum data type word length, while the fraction length is set to maintain the most significant bits of the sum and no more fractional bits than necessary • SpecifyPrecision — Specify the word and fraction lengths or the slope and bias of the sum data type. |
| SumSlope | Slope of the sum data type | <ul style="list-style-type: none"> • 9.3132e-010 (default) • Any floating-point number <p>Note</p> $SumSlope = SumSlopeAdjustmentFactor \times 2^{SumFixedExponent}$ <p>Changing one of these properties affects the others.</p> |
| SumSlopeAdjustmentFactor | Slope adjustment factor of the sum data type | <ul style="list-style-type: none"> • 1 (default) • Any floating-point number greater than or equal to 1 and less than 2 <p>Note</p> $SumSlope = SumSlopeAdjustmentFactor \times 2^{SumFixedExponent}$ <p>Changing one of these properties affects the others.</p> |

| Property | Description | Valid Values |
|---------------|--|--|
| SumWordLength | Word length, in bits, of the sum data type | <ul style="list-style-type: none"> • 32 (default) • Any positive integer |

For details about these properties, refer to the “Set fi Object Properties” on page 2-15. To learn how to specify properties for fimath objects in Fixed-Point Designer software, refer to “Setting fimath Object Properties” on page 3-8.

How Properties are Related

Sum data type properties

The slope of the sum of two fi objects is related to the SumSlopeAdjustmentFactor and SumFixedExponent properties by

$$\text{SumSlope} = \text{SumSlopeAdjustmentFactor} \times 2^{\text{SumFixedExponent}}$$

If any of these properties are updated, the others are modified accordingly.

In a FullPrecision sum, the resulting word length is represented by

$$W_s = \text{integer length} + F_s$$

where

$$\text{integer length} = \max(W_a - F_a, W_b - F_b) + \text{ceil}(\log_2(\text{NumberOfSummands}))$$

and

$$F_s = \max(F_a, F_b)$$

When the SumMode is set to KeepLSB, the resulting word length and fraction length is determined by

$$W_s = \text{specified in the SumWordLength property}$$

$$F_s = \max(F_a, F_b)$$

When the SumMode is set to KeepMSB, the resulting word length and fraction length is determined by

$$W_s = \text{specified in the SumWordLength property}$$

$$F_s = W_s - \text{integer length}$$

where

$$\text{integer length} = \max(W_a - F_a, W_b - F_b) + \text{ceil}(\log_2(\text{NumberOfSummands}))$$

When the SumMode is set to SpecifyPrecision, you specify both the word and fraction length or slope and bias of the sum data type with the SumWordLength and SumFractionLength, or SumSlope and SumBias properties respectively.

Product data type properties

The slope of the product of two fi objects is related to the ProductSlopeAdjustmentFactor and ProductFixedExponent properties by

$$ProductSlope = ProductSlopeAdjustmentFactor \times 2^{ProductFixedExponent}$$

If any of these properties are updated, the others are modified accordingly.

In a `FullPrecision` multiply, the resulting word length and fraction length are represented by

$$W_p = W_a + W_b$$

$$F_p = F_a + F_b$$

When the `ProductMode` is `KeepLSB` the word length and fraction length are determined by

$$W_p = \text{specified in the ProductWordLength property}$$

$$F_p = F_a + F_b$$

When the `ProductMode` is `KeepMSB` the word length and fraction length are

$$W_p = \text{specified in the ProductWordLength property}$$

$$F_p = W_p - \text{integer length}$$

where

$$\text{integer length} = (W_a + W_b) - (F_a + F_b)$$

When the `ProductMode` is set to `SpecifyPrecision`, you specify both the word and fraction length or slope and bias of the product data type with the `ProductWordLength` and `ProductFractionLength`, or `ProductSlope` and `ProductBias` properties respectively.

For more information about how certain functions use the `fimath` properties, see

Setting fimath Object Properties

- “Setting fimath Properties at Object Creation” on page 3-8
- “Using Direct Property Referencing with fimath” on page 3-9

Setting fimath Properties at Object Creation

You can set properties of `fimath` objects at the time of object creation by including properties after the arguments of the `fimath` constructor function.

For example, to set the overflow action to `Saturate` and the rounding method to `Convergent`,

```
F = fimath('OverflowAction', 'Saturate', 'RoundingMethod', 'Convergent')
```

```
F =
```

```

RoundingMethod: Convergent
OverflowAction: Saturate
ProductMode: FullPrecision
SumMode: FullPrecision

```

In addition to creating a `fimath` object at the command line, you can also set `fimath` properties using the **Insert fimath Constructor** dialog box. For an example of this approach, see “Building fimath Object Constructors in a GUI” on page 3-3.

Using Direct Property Referencing with fimath

You can reference directly into a property for setting or retrieving `fimath` object property values using MATLAB structure-like referencing. You do so by using a period to index into a property by name.

For example, to get the `RoundingMethod` of `F`,

```
F.RoundingMethod
```

```
ans =
```

```
    'Convergent'
```

To set the `OverflowAction` of `F`,

```
F.OverflowAction = 'Wrap'
```

```
F =
```

```
    RoundingMethod: Convergent  
    OverflowAction: Wrap  
    ProductMode: FullPrecision  
    SumMode: FullPrecision
```

fimath Properties Usage for Fixed-Point Arithmetic

In this section...

“fimath Rules for Fixed-Point Arithmetic” on page 3-10

“Binary-Point Arithmetic” on page 3-11

“[Slope Bias] Arithmetic” on page 3-14

fimath Rules for Fixed-Point Arithmetic

`fimath` properties define the rules for performing arithmetic operations on `fi` objects. The `fimath` properties that govern fixed-point arithmetic operations can come from a local `fimath` object or the `fimath` default values.

To determine whether a `fi` object has a local `fimath` object, use the `isfimathlocal` function.

The following sections discuss how `fi` objects with local `fimath` objects interact with `fi` objects without local `fimath`.

Binary Operations

In binary fixed-point operations such as $c = a + b$, the following rules apply:

- If both `a` and `b` have no local `fimath`, the operation uses default `fimath` values to perform the fixed-point arithmetic. The output `fi` object `c` also has no local `fimath`.
- If either `a` or `b` has a local `fimath` object, the operation uses that `fimath` object to perform the fixed-point arithmetic. The output `fi` object `c` has the same local `fimath` object as the input.

Unary Operations

In unary fixed-point operations such as $b = \text{abs}(a)$, the following rules apply:

- If `a` has no local `fimath`, the operation uses default `fimath` values to perform the fixed-point arithmetic. The output `fi` object `b` has no local `fimath`.
- If `a` has a local `fimath` object, the operation uses that `fimath` object to perform the fixed-point arithmetic. The output `fi` object `b` has the same local `fimath` object as the input `a`.

When you specify a `fimath` object in the function call of a unary fixed-point operation, the operation uses the `fimath` object you specify to perform the fixed-point arithmetic. For example, when you use a syntax such as $b = \text{abs}(a, F)$ or $b = \text{sqrt}(a, F)$, the `abs` and `sqrt` operations use the `fimath` object `F` to compute intermediate quantities. The output `fi` object `b` always has no local `fimath`.

Concatenation Operations

In fixed-point concatenation operations such as $c = [a \ b]$, $c = [a; b]$ and $c = \text{bitconcat}(a, b)$, the following rule applies:

- The `fimath` properties of the leftmost `fi` object in the operation determine the `fimath` properties of the output `fi` object `c`.

For example, consider the following scenarios for the operation $d = [a \ b \ c]$:

- If `a` is a `fi` object with no local `fimath`, the output `fi` object `d` also has no local `fimath`.
- If `a` has a local `fimath` object, the output `fi` object `d` has the same local `fimath` object.
- If `a` is not a `fi` object, the output `fi` object `d` inherits the `fimath` properties of the next leftmost `fi` object. For example, if `b` is a `fi` object with a local `fimath` object, the output `fi` object `d` has the same local `fimath` object as the input `fi` object `b`.

fimath Object Operations: add, mpy, sub

The output of the `fimath` object operations `add`, `mpy`, and `sub` always have no local `fimath`. The operations use the `fimath` object you specify in the function call, but the output `fi` object never has a local `fimath` object.

MATLAB Function Block Operations

Fixed-point operations performed with the MATLAB Function block use the same rules as fixed-point operations performed in MATLAB.

All input signals to the MATLAB Function block that you treat as `fi` objects associate with whatever you specify for the **MATLAB Function block `fimath`** parameter. When you set this parameter to `Same as MATLAB`, your `fi` objects do not have local `fimath`. When you set the **MATLAB Function block `fimath`** parameter to `Specify other`, you can define your own set of `fimath` properties for all `fi` objects in the MATLAB Function block to associate with. You can choose to treat only fixed-point input signals as `fi` objects or both fixed-point and integer input signals as `fi` objects. See “Use `fimath` Objects in MATLAB Function Blocks” on page 45-16.

Binary-Point Arithmetic

The `fimath` object encapsulates the math properties of Fixed-Point Designer software.

`fi` objects only have a local `fimath` object when you explicitly specify `fimath` properties in the `fi` constructor. When you use the `sfi` or `ufi` constructor or do not specify any `fimath` properties in the `fi` constructor, the resulting `fi` object does not have any local `fimath` and uses default `fimath` values.

```
a = fi(pi)
a =
    3.1416
        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 16
        FractionLength: 13

a.fimath
isfimathlocal(a)

ans =
        RoundingMethod: Nearest
        OverflowAction: Saturate
        ProductMode: FullPrecision
        SumMode: FullPrecision
```

```
ans =  
    logical  
    0
```

To perform arithmetic with `+`, `-`, `.*`, or `*` on two `fi` operands with local `fimath` objects, the local `fimath` objects must be identical. If one of the `fi` operands does not have a local `fimath`, the `fimath` properties of the two operands need not be identical. See “`fimath` Rules for Fixed-Point Arithmetic” on page 3-10 for more information.

```
a = fi(pi);  
b = fi(8);  
isequal(a.fimath, b.fimath)
```

```
ans =  
    logical  
    1
```

```
a + b
```

```
ans =  
    11.1416
```

```
        DataTypeMode: Fixed-point: binary point scaling  
        Signedness: Signed  
        WordLength: 19  
        FractionLength: 13
```

To perform arithmetic with `+`, `-`, `.*`, or `*`, two `fi` operands must also have the same data type. For example, you can add two `fi` objects with data type `double`, but you cannot add an object with data type `double` and one with data type `single`:

```
a = fi(3, 'DataType', 'double')
```

```
a =  
    3
```

```
        DataTypeMode: Double
```

```
b = fi(27, 'DataType', 'double')
```

```
b =  
    27
```

```
        DataTypeMode: Double
```

```
a + b
```

```
ans =  
    30
```

```
        DataTypeMode: Double
```

```
c = fi(12, 'DataType', 'single')
```

```
c =
```

```
12
```

```
    DataTypeMode: Single
```

```
a + c
```

```
Error using + (line 24)
```

```
Math operations are not allowed on fi objects with different data types.
```

Fixed-point `fi` object operands do not have to have the same scaling. You can perform binary math operations on a `fi` object with a fixed-point data type and a `fi` object with a scaled doubles data type. In this sense, the scaled double data type acts as a fixed-point data type:

```
a = fi(pi)
```

```
a =
```

```
3.1416
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 13
```

```
b = fi(magic(2), ...
'DataTypeMode', 'Scaled double: binary point scaling')
```

```
b =
```

```
1 3
4 2
```

```
    DataTypeMode: Scaled double: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 12
```

```
a + b
```

```
ans =
```

```
4.1416 6.1416
7.1416 5.1416
```

```
    DataTypeMode: Scaled double: binary point scaling
    Signedness: Signed
    WordLength: 18
    FractionLength: 13
```

Use the `divide` function to perform division with doubles, singles, or binary point-only scaling `fi` objects.

[Slope Bias] Arithmetic

Fixed-Point Designer software supports fixed-point arithmetic using the local `fimath` object or default `fimath` for all binary point-only signals. The toolbox also supports arithmetic for [Slope Bias] signals with the following restrictions:

- [Slope Bias] signals must be real.
- You must set the `SumMode` and `ProductMode` properties of the governing `fimath` to `'SpecifyPrecision'` for sum and multiply operations, respectively.
- You must set the `CastBeforeSum` property of the governing `fimath` to `'true'`.
- Fixed-Point Designer does not support the `divide` function for [Slope Bias] signals.

```
f = fimath('SumMode', 'SpecifyPrecision', ...
          'SumFractionLength', 16)
```

```
f =
```

```

    RoundingMethod: Nearest
    OverflowAction: Saturate
      ProductMode: FullPrecision
        SumMode: SpecifyPrecision
    SumWordLength: 32
    SumFractionLength: 16
    CastBeforeSum: true
```

```
a = fi(pi, 'fimath', f)
```

```
a =
```

```
3.1416
```

```

    DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 16
    FractionLength: 13
```

```

    RoundingMethod: Nearest
    OverflowAction: Saturate
      ProductMode: FullPrecision
        SumMode: SpecifyPrecision
    SumWordLength: 32
    SumFractionLength: 16
    CastBeforeSum: true
```

```
b = fi(22, true, 16, 2^-8, 3, 'fimath', f)
```

```
b =
```

```
22
```

```

    DataTypeMode: Fixed-point: slope and bias scaling
      Signedness: Signed
      WordLength: 16
      Slope: 0.00390625
      Bias: 3
```

```
RoundingMethod: Nearest
OverflowAction: Saturate
  ProductMode: FullPrecision
    SumMode: SpecifyPrecision
  SumWordLength: 32
SumFractionLength: 16
CastBeforeSum: true
```

a + b

ans =

25.1416

```
DataTypeMode: Fixed-point: binary point scaling
  Signedness: Signed
  WordLength: 32
FractionLength: 16
```

```
RoundingMethod: Nearest
OverflowAction: Saturate
  ProductMode: FullPrecision
    SumMode: SpecifyPrecision
  SumWordLength: 32
SumFractionLength: 16
CastBeforeSum: true
```

Setting the `SumMode` and `ProductMode` properties to `SpecifyPrecision` are mutually exclusive except when performing the `*` operation between matrices. In this case, you must set both the `SumMode` and `ProductMode` properties to `SpecifyPrecision` for [Slope Bias] signals. Doing so is necessary because the `*` operation performs both sum and multiply operations to calculate the result.

fimath for Rounding and Overflow Modes

Only rounding methods and overflow actions set prior to an operation with `fi` objects affect the outcome of those operations. Once you create a `fi` object in MATLAB, changing its rounding or overflow settings does not affect its value. For example, consider the `fi` objects `a` and `b`:

```
p = fipref('NumberDisplay', 'RealWorldValue', ...  
         'NumericTypeDisplay', 'none', 'FimathDisplay', 'none');  
T = numericitytype('WordLength',8,'FractionLength',7);  
F = fimath('RoundingMethod', 'Floor', 'OverflowAction', 'Wrap');  
a = fi(1,T,F)
```

```
a =  
  
    -1
```

```
b = fi(1,T)
```

```
b =  
  
    0.9922
```

Because you create `a` with a `fimath` object `F` that has `OverflowAction` set to `Wrap`, the value of `a` wraps to `-1`. Conversely, because you create `b` with the default `OverflowAction` value of `Saturate`, its value saturates to `0.9922`.

Now, assign the `fimath` object `F` to `b`:

```
b.fimath = F
```

```
b =  
  
    0.9922
```

Because the assignment operation and corresponding overflow and saturation happened when you created `b`, its value does not change when you assign it the new `fimath` object `F`.

Note `fi` objects with no local `fimath` and created from a floating-point value always get constructed with a `RoundingMethod` of `Nearest` and an `OverflowAction` of `Saturate`. To construct `fi` objects with different `RoundingMethod` and `OverflowAction` properties, specify the desired `RoundingMethod` and `OverflowAction` properties in the `fi` constructor.

For more information about the `fimath` object and its properties, see “`fimath` Object Properties” on page 3-4

fimath for Sharing Arithmetic Rules

There are two ways of sharing `fimath` properties in Fixed-Point Designer software:

- “Default `fimath` Usage to Share Arithmetic Rules” on page 3-17
- “Local `fimath` Usage to Share Arithmetic Rules” on page 3-17

Sharing `fimath` properties across `fi` objects ensures that the `fi` objects are using the same arithmetic rules and helps you avoid “mismatched `fimath`” errors.

Default `fimath` Usage to Share Arithmetic Rules

You can ensure that your `fi` objects are all using the same `fimath` properties by not specifying any local `fimath`. To assure no local `fimath` is associated with a `fi` object, you can:

- Create a `fi` object using the `fi` constructor without specifying any `fimath` properties in the constructor call. For example:

```
a = fi(pi)
```

- Create a `fi` object using the `sfi` or `ufi` constructor. All `fi` objects created with these constructors have no local `fimath`.

```
b = sfi(pi)
```

- Use `removefimath` to remove a local `fimath` object from an existing `fi` object.

Local `fimath` Usage to Share Arithmetic Rules

You can also use a `fimath` object to define common arithmetic rules that you would like to use for multiple `fi` objects. You can then create your `fi` objects, using the same `fimath` object for each. To do so, you must also create a `numericType` object to define a common data type and scaling. Refer to “`numericType` Object Construction” on page 5-2 for more information on `numericType` objects. The following example shows the creation of a `numericType` object and `fimath` object, and then uses those objects to create two `fi` objects with the same `numericType` and `fimath` attributes:

```
T = numericType('WordLength', 32, 'FractionLength', 30)
```

```
T =
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 32
    FractionLength: 30
```

```
F = fimath('RoundingMethod', 'Floor', 'OverflowAction', 'Wrap')
```

```
F =
```

```
    RoundingMethod: Floor
    OverflowAction: Wrap
    ProductMode: FullPrecision
    SumMode: FullPrecision
```

```
a = fi(pi, T, F)
```

a =

-0.8584

 DataTypeMode: Fixed-point: binary point scaling
 Signedness: Signed
 WordLength: 32
 FractionLength: 30

 RoundingMethod: Floor
 OverflowAction: Wrap
 ProductMode: FullPrecision
 SumMode: FullPrecision

b = fi(pi/2, T, F)

b =

1.5708

 DataTypeMode: Fixed-point: binary point scaling
 Signedness: Signed
 WordLength: 32
 FractionLength: 30

 RoundingMethod: Floor
 OverflowAction: Wrap
 ProductMode: FullPrecision
 SumMode: FullPrecision

fimath ProductMode and SumMode

In this section...

“Example Setup” on page 3-19
 “FullPrecision” on page 3-19
 “KeepLSB” on page 3-20
 “KeepMSB” on page 3-21
 “SpecifyPrecision” on page 3-22

Example Setup

The examples in the sections of this topic show the differences among the four settings of the ProductMode and SumMode properties:

- FullPrecision
- KeepLSB
- KeepMSB
- SpecifyPrecision

To follow along, first set the following preferences:

```
p = fipref;
p.NumericTypeDisplay = 'short';
p.FimathDisplay = 'none';
p.LoggingMode = 'on';
F = fimath('OverflowAction','Wrap','RoundingMethod','Floor',...
          'CastBeforeSum',false);
warning off
format compact
```

Next, define `fi` objects `a` and `b`. Both have signed 8-bit data types. The fraction length gets chosen automatically for each `fi` object to yield the best possible precision:

```
a = fi(pi, true, 8)
a =
    3.1562
    numerictype(1,8,5)
b = fi(exp(1), true, 8)
b =
    2.7188
    numerictype(1,8,5)
```

FullPrecision

Now, set ProductMode and SumMode for `a` and `b` to FullPrecision and look at some results:

```
F.ProductMode = 'FullPrecision';
F.SumMode = 'FullPrecision';
a.fimath = F;
```

```
b.fimath = F;
a
a =
    3.1562
    numerictype(1,8,5)
b
b =
    2.7188
    numerictype(1,8,5)
a*b
ans =
    8.5811
    numerictype(1,16,10)
a+b
ans =
    5.8750
    numerictype(1,9,5)
```

In `FullPrecision` mode, the product word length grows to the sum of the word lengths of the operands. In this case, each operand has 8 bits, so the product word length is 16 bits. The product fraction length is the sum of the fraction lengths of the operands, in this case $5 + 5 = 10$ bits.

The sum word length grows by one most significant bit to accommodate the possibility of a carry bit. The sum fraction length aligns with the fraction lengths of the operands, and all fractional bits are kept for full precision. In this case, both operands have 5 fractional bits, so the sum has 5 fractional bits.

KeepLSB

Now, set `ProductMode` and `SumMode` for `a` and `b` to `KeepLSB` and look at some results:

```
F.ProductMode = 'KeepLSB';
F.ProductWordLength = 12;
F.SumMode = 'KeepLSB';
F.SumWordLength = 12;
a.fimath = F;
b.fimath = F;
a
a =
    3.1562
    numerictype(1,8,5)
b
b =
    2.7188
    numerictype(1,8,5)
a*b
```

```

ans =
    0.5811
    numerictype(1,12,10)

a+b

ans =
    5.8750
    numerictype(1,12,5)

```

In **KeepLSB** mode, you specify the word lengths and the least significant bits of results are automatically kept. This mode models the behavior of integer operations in the C language.

The product fraction length is the sum of the fraction lengths of the operands. In this case, each operand has 5 fractional bits, so the product fraction length is 10 bits. In this mode, all 10 fractional bits are kept. Overflow occurs because the full-precision result requires 6 integer bits, and only 2 integer bits remain in the product.

The sum fraction length aligns with the fraction lengths of the operands, and in this model all least significant bits are kept. In this case, both operands had 5 fractional bits, so the sum has 5 fractional bits. The full-precision result requires 4 integer bits, and 7 integer bits remain in the sum, so no overflow occurs in the sum.

KeepMSB

Now, set ProductMode and SumMode for a and b to KeepMSB and look at some results:

```

F.ProductMode = 'KeepMSB';
F.ProductWordLength = 12;
F.SumMode = 'KeepMSB';
F.SumWordLength = 12;
a.fimath = F;
b.fimath = F;
a

a =
    3.1562
    numerictype(1,8,5)

b

b =
    2.7188
    numerictype(1,8,5)

a*b

ans =
    8.5781
    numerictype(1,12,6)

a+b

ans =
    5.8750
    numerictype(1,12,8)

```

In **KeepMSB** mode, you specify the word lengths and the most significant bits of sum and product results are automatically kept. This mode models the behavior of many DSP devices where the

product and sum are kept in double-wide registers, and the programmer chooses to transfer the most significant bits from the registers to memory after each operation.

The full-precision product requires 6 integer bits, and the fraction length of the product is adjusted to accommodate all 6 integer bits in this mode. No overflow occurs. However, the full-precision product requires 10 fractional bits, and only 6 are available. Therefore, precision is lost.

The full-precision sum requires 4 integer bits, and the fraction length of the sum is adjusted to accommodate all 4 integer bits in this mode. The full-precision sum requires only 5 fractional bits; in this case there are 8, so there is no loss of precision.

This example shows that, in KeepMSB mode the fraction length changes regardless of whether an overflow occurs. The fraction length is set to the amount needed to represent the product in case both terms use the maximum possible value ($18+18-16=20$ in this example).

```
F = fimath('SumMode','KeepMSB','ProductMode','KeepMSB',...
          'ProductWordLength',16,'SumWordLength',16);
a = fi(100,1,16,-2,'fimath',F);
a*a

ans =
    0
    numerictype(1,16,-20)
```

SpecifyPrecision

Now set ProductMode and SumMode for a and b to SpecifyPrecision and look at some results:

```
F.ProductMode = 'SpecifyPrecision';
F.ProductWordLength = 8;
F.ProductFractionLength = 7;
F.SumMode = 'SpecifyPrecision';
F.SumWordLength = 8;
F.SumFractionLength = 7;
a.fimath = F;
b.fimath = F;
a

a =
    3.1562
    numerictype(1,8,5)

b

b =
    2.7188
    numerictype(1,8,5)

a*b

ans =
    0.5781
    numerictype(1,8,7)

a+b
```

```
ans =  
-0.1250  
    numerictype(1,8,7)
```

In **SpecifyPrecision** mode, you must specify both word length and fraction length for sums and products. This example unwisely uses fractional formats for the products and sums, with 8-bit word lengths and 7-bit fraction lengths.

The full-precision product requires 6 integer bits, and the example specifies only 1, so the product overflows. The full-precision product requires 10 fractional bits, and the example only specifies 7, so there is precision loss in the product.

The full-precision sum requires 4 integer bits, and the example specifies only 1, so the sum overflows. The full-precision sum requires 5 fractional bits, and the example specifies 7, so there is no loss of precision in the sum.

For more information about the `fimath` object and its properties, see “[fimath Object Properties](#)” on page 3-4

How Functions Use fimath

| In this section... |
|--|
| “Functions that use then discard attached fimath” on page 3-24 |
| “Functions that ignore and discard attached fimath” on page 3-24 |
| “Functions that do not perform math” on page 3-24 |

Functions that use then discard attached fimath

| Functions | Note |
|--------------|-----------------------------------|
| conv, filter | Error if attached fimaths differ. |
| mean, median | — |

Functions that ignore and discard attached fimath

| Functions | Note |
|--|---|
| accumneg, accumpos | <ul style="list-style-type: none"> By default, use Floor rounding method and Wrap overflow |
| add, sub, mpy | <ul style="list-style-type: none"> Override and discard any fimath objects attached to the input fi objects Uses the fimath from input, F, as in add(F, a, b) |
| CORDIC functions (see “CORDIC Algorithms in MATLAB”) | CORDIC functions use their own internal fimath: <ul style="list-style-type: none"> Rounding Mode - Floor Overflow Action - Wrap |
| mod | — |
| qr | — |
| quantize | Uses the math settings on the quantizer object, ignores and discards any fimath settings on the input |
| Trigonometric functions — atan2, cos, sin | — |

Functions that do not perform math

| Functions | Note |
|--|--|
| Built-in Types—int32, int64, int8, uint16, uint32, uint64, uint8 | Ignore any fimath settings on the input. Always use the rounding method Round when casting to the new data type. The output is not a fi object so it has no attached fimath. |
| bitsll, bitsra, bitsrl | OverflowAction and RoundingMethod are ignored — bits drop off the end. |

| Functions | Note |
|------------------|---|
| bitshift | RoundingMethod is ignored, but OverflowAction property is obeyed. |

Working with fipref Objects

- “Set fi Object Display Preferences Using fipref” on page 4-2
- “Underflow and Overflow Logging Using fipref” on page 4-3
- “Data Type Override Preferences Using fipref” on page 4-7

Set fi Object Display Preferences Using fipref

You use the `fipref` object to specify three aspects of the display of `fi` objects: the object values, the local `fimath` properties, and the `numericType` properties.

For example, the following code shows the default `fipref` display for a `fi` object with a local `fimath` object:

```
a = fi(pi, 'RoundingMethod', 'Floor', 'OverflowAction', 'Wrap')
a =
    3.1415

    DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 16
  FractionLength: 13

  RoundingMethod: Floor
  OverflowAction: Wrap
    ProductMode: FullPrecision
      SumMode: FullPrecision
```

The default `fipref` display for a `fi` object with no local `fimath` is as follows:

```
a = fi(pi)
a =
    3.1416

    DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 16
  FractionLength: 13
```

Next, change the `fipref` display properties:

```
P = fipref;
P.NumberDisplay = 'bin';
P.NumericTypeDisplay = 'short';
P.FimathDisplay = 'none'

P =
    NumberDisplay: 'bin'
  NumericTypeDisplay: 'short'
    FimathDisplay: 'none'
      LoggingMode: 'Off'
  DataTypeOverride: 'ForceOff'

a

a =
0110010010001000
    numericType(1,16,13)
```

Underflow and Overflow Logging Using fipref

In this section...

“Logging Overflows and Underflows as Warnings” on page 4-3

“Accessing Logged Information with Functions” on page 4-5

Logging Overflows and Underflows as Warnings

Overflows and underflows are logged as warnings for all assignment, plus, minus, and multiplication operations when the `fipref` `LoggingMode` property is set to `on`. For example, try the following:

- 1 Create a signed `fi` object that is a vector of values from 1 to 5, with 8-bit word length and 6-bit fraction length.


```
a = fi(1:5,1,8,6);
```
- 2 Define the `fimath` object associated with `a`, and indicate that you will specify the sum and product word and fraction lengths.


```
F = a.fimath;
F.SumMode = 'SpecifyPrecision';
F.ProductMode = 'SpecifyPrecision';
a.fimath = F;
```
- 3 Define the `fipref` object and turn on overflow and underflow logging.


```
P = fipref;
P.LoggingMode = 'on';
```
- 4 Suppress the `numericType` and `fimath` displays.


```
P.NumericTypeDisplay = 'none';
P.FimathDisplay = 'none';
```
- 5 Specify the sum and product word and fraction lengths.


```
a.SumWordLength = 16;
a.SumFractionLength = 15;
a.ProductWordLength = 16;
a.ProductFractionLength = 15;
```
- 6 Warnings are displayed for overflows and underflows in assignment operations. For example, try:


```
a(1) = pi
```

Warning: 1 overflow(s) occurred in the `fi` assignment operation.

```
a =
    1.9844    1.9844    1.9844    1.9844    1.9844
```

```
a(1) = double(eps(a))/10
```

Warning: 1 underflow(s) occurred in the `fi` assignment operation.

```
a =
     0    1.9844    1.9844    1.9844    1.9844
```

- 7** Warnings are displayed for overflows and underflows in addition and subtraction operations. For example, try:

```
a+a
```

```
Warning: 12 overflow(s) occurred in the fi + operation.  
> In + (line 24)
```

```
ans =
```

```
0 1.0000 1.0000 1.0000 1.0000
```

```
a-a
```

```
Warning: 8 overflow(s) occurred in the fi - operation.  
> In - (line 24)
```

```
ans =
```

```
0 0 0 0 0
```

- 8** Warnings are displayed for overflows and underflows in multiplication operations. For example, try:

```
a.*a
```

```
Warning: 4 product overflow(s) occurred in the fi .* operation.  
> In .* (line 24)
```

```
ans =
```

```
0 1.0000 1.0000 1.0000 1.0000
```

```
a*a'
```

```
Warning: 4 product overflow(s) occurred in the fi * operation.  
> In * (line 25)  
Warning: 3 sum overflow(s) occurred in the fi * operation.  
> In * (line 25)
```

```
ans =
```

```
1.0000
```

The final example above is a complex multiplication that requires both multiplication and addition operations. The warnings inform you of overflows and underflows in both.

Because overflows and underflows are logged as warnings, you can use the `dbstop` MATLAB function with the syntax

```
dbstop if warning
```

to find the exact lines in a file that are causing overflows or underflows.

Use

```
dbstop if warning fi:underflow
```

to stop only on lines that cause an underflow. Use

```
dbstop if warning fi:overflow
```

to stop only on lines that cause an overflow.

Accessing Logged Information with Functions

When the `fipref` `LoggingMode` property is set to `on`, you can use the following functions to return logged information about assignment and creation operations to the MATLAB command line:

- `maxlog` — Returns the maximum real-world value
- `minlog` — Returns the minimum value
- `noverflows` — Returns the number of overflows
- `nunderflows` — Returns the number of underflows

`LoggingMode` must be set to `on` before you perform any operation in order to log information about it. To clear the log, use the function `resetlog`.

For example, consider the following. First turn logging on, then perform operations, and then finally get information about the operations:

```
fipref('LoggingMode','on');
x = fi([-1.5 eps 0.5], true, 16, 15);
x(1) = 3.0;
```

```
maxlog(x)
```

```
ans =
```

```
1.0000
```

```
minlog(x)
```

```
ans =
```

```
-1
```

```
noverflows(x)
```

```
ans =
```

```
2
```

```
nunderflows(x)
```

```
ans =
```

```
1
```

Next, reset the log and request the same information again. Note that the functions return empty `[]`, because logging has been reset since the operations were run:

```
resetlog(x)
```

```
maxlog(x)
```

```
Warning: Logging is turned on in 'maxlog'. However,
no values have been logged for this variable yet.
```

```
ans =
```

```
    []
```

```
minlog(x)
```

```
Warning: Logging is turned on in 'minlog'. However,  
no values have been logged for this variable yet.
```

```
ans =
```

```
    []
```

```
noverflows(x)
```

```
Warning: Logging is turned on in 'noverflows'. However,  
no values have been logged for this variable yet.
```

```
ans =
```

```
    []
```

```
nunderflows(x)
```

```
Warning: Logging is turned on in 'nunderflows'. However,  
no values have been logged for this variable yet.
```

```
ans =
```

```
    []
```

Data Type Override Preferences Using fipref

In this section...

“Overriding the Data Type of fi Objects” on page 4-7

“Data Type Override for Fixed-Point Scaling” on page 4-8

Overriding the Data Type of fi Objects

Use the `fipref` `DataTypeOverride` property to override `fi` objects with singles, doubles, or scaled doubles. Data type override only occurs when the `fi` constructor function is called. Objects that are created while data type override is on have the overridden data type. They maintain that data type when data type override is later turned off. To obtain an object with a data type that is not the override data type, you must create an object when data type override is off:

```
p = fipref('DataTypeOverride', 'TrueDoubles')
```

```
p =
```

```
    NumberDisplay: 'RealWorldValue'
  NumericTypeDisplay: 'full'
    FimathDisplay: 'full'
      LoggingMode: 'Off'
  DataTypeOverride: 'TrueDoubles'
```

```
a = fi(pi)
```

```
a =
```

```
    3.1416
```

```
    DataTypeMode: Double
```

```
p = fipref('DataTypeOverride', 'ForceOff')
```

```
p =
```

```
    NumberDisplay: 'RealWorldValue'
  NumericTypeDisplay: 'full'
    FimathDisplay: 'full'
      LoggingMode: 'Off'
  DataTypeOverride: 'ForceOff'
```

```
a
```

```
a =
```

```
    3.1416
```

```
    DataTypeMode: Double
```

```
b = fi(pi)
```

```
b =
```

3.1416

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 16
FractionLength: 13
```

Tip To reset the `fipref` object to its default values use `reset(fipref)` or `reset(p)`, where `p` is a `fipref` object. This is useful to ensure that data type override and logging are off.

Data Type Override for Fixed-Point Scaling

Choosing the scaling for the fixed-point variables in your algorithms can be difficult. In Fixed-Point Designer software, you can use a combination of data type override and min/max logging to help you discover the numerical ranges that your fixed-point data types need to cover. These ranges dictate the appropriate scalings for your fixed-point data types. In general, the procedure is

- 1 Implement your algorithm using fixed-point `fi` objects, using initial “best guesses” for word lengths and scalings.
- 2 Set the `fipref` `DataTypeOverride` property to `ScaledDoubles`, `TrueSingles`, or `TrueDoubles`.
- 3 Set the `fipref` `LoggingMode` property to `on`.
- 4 Use the `maxlog` and `minlog` functions to log the maximum and minimum values achieved by the variables in your algorithm in floating-point mode.
- 5 Set the `fipref` `DataTypeOverride` property to `ForceOff`.
- 6 Use the information obtained in step 4 to set the fixed-point scaling for each variable in your algorithm such that the full numerical range of each variable is representable by its data type and scaling.

A detailed example of this process is shown in the Fixed-Point Designer “Set Data Types Using Min/Max Instrumentation” on page 54-125 example.

Working with numerictype Objects

- “numerictype Object Construction” on page 5-2
- “numerictype Object Properties” on page 5-5
- “numerictype of Fixed-Point Objects” on page 5-9
- “numerictype Objects Usage to Share Data Type and Scaling Settings of fi objects” on page 5-12

numerictype Object Construction

In this section...

“numerictype Object Syntaxes” on page 5-2

“Example: Construct a numerictype Object with Property Name and Property Value Pairs” on page 5-2

“Example: Copy a numerictype Object” on page 5-3

“Example: Build numerictype Object Constructors in a GUI” on page 5-4

numerictype Object Syntaxes

numerictype objects define the data type and scaling attributes of `fi` objects, as well as Simulink signals and model parameters. You can create numerictype objects in Fixed-Point Designer software in one of two ways:

- You can use the numerictype constructor function to create a new object.
- You can use the numerictype constructor function to copy an existing numerictype object.

To create a default numerictype object, type

```
T = numerictype
```

```
T =
```

```

    DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 16
    FractionLength: 15

```

To see all of the numerictype object syntaxes, refer to the numerictype constructor function reference page.

The following examples show different ways of constructing numerictype objects. For more examples of constructing numerictype objects, see the “Examples” on the numerictype constructor function reference page.

Example: Construct a numerictype Object with Property Name and Property Value Pairs

When you create a numerictype object using property name and property value pairs, Fixed-Point Designer software first creates a default numerictype object, and then, for each property name you specify in the constructor, assigns the corresponding value.

This behavior differs from the behavior that occurs when you use a syntax such as `T = numerictype(s,w)`, where you only specify the property values in the constructor. Using such a syntax results in no default numerictype object being created, and the numerictype object receives only the assigned property values that are specified in the constructor.

The following example shows how the property name/property value syntax creates a slightly different `numerictype` object than the property values syntax, even when you specify the same property values in both constructors.

To demonstrate this difference, suppose you want to create an unsigned `numerictype` object with a word length of 32 bits.

First, create the `numerictype` object using property name/property value pairs.

```
T1 = numerictype('Signed',0,'WordLength',32)
```

```
T1 =
```

```

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Unsigned
    WordLength: 32
    FractionLength: 15

```

The `numerictype` object `T1` has the same `DataTypeMode` and `FractionLength` as a default `numerictype` object, but the `WordLength` and `Signed` properties are overwritten with the values you specified.

Now, create another unsigned 32 bit `numerictype` object, but this time specify only property values in the constructor.

```
T2 = numerictype(0,32)
```

```
T2 =
```

```

    DataTypeMode: Fixed-point: unspecified scaling
    Signedness: Unsigned
    WordLength: 32

```

Unlike `T1`, `T2` only has the property values you specified. The `DataTypeMode` of `T2` is `Fixed-Point: unspecified scaling`, so no fraction length is assigned.

`fi` objects cannot have unspecified `numerictype` properties. Thus, all unspecified `numerictype` object properties become specified at the time of `fi` object creation.

Example: Copy a numerictype Object

To copy a `numerictype` object, use assignment:

```
T = numerictype;
U = T;
isequal(T,U)
```

```
ans =
```


```
logical
```

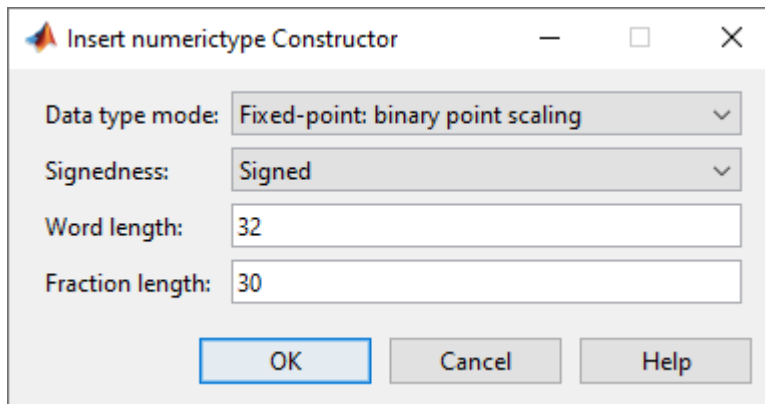
```
1
```

Example: Build numerictype Object Constructors in a GUI

When you are working with files in MATLAB, you can build your `numerictype` object constructors using the **Insert numerictype Constructor** dialog box. After specifying the properties of the `numerictype` object in the dialog box, you can insert the pre-populated `numerictype` object constructor at a specific location in your file.

For example, create a signed `numerictype` object with binary-point scaling, a word length of 32 bits and a fraction length of 30 bits:

- 1 On the **Home** tab, in the **File** section, click **New > Script** to open the MATLAB Editor
- 2 On the **Editor** tab, in the **Edit** section of the toolbar, click  in the **Insert** button group. Click the **Insert numerictype** to open the **Insert numerictype Constructor** dialog box.
- 3 Use the edit boxes and drop-down menus to specify the following properties of the `numerictype` object:
 - **Data type mode:** Fixed-point: binary point scaling
 - **Signedness:** Signed
 - **Word length:** 32
 - **Fraction length:** 30



- 4 To insert the `numerictype` object constructor in your file, place your cursor at the desired location in the file, and click **OK** on the **Insert numerictype Constructor** dialog box. Clicking **OK** closes the **Insert numerictype Constructor** dialog box and automatically populates the `numerictype` object constructor in your file:

```
numerictype(1, 32, 30)
```

numerictype Object Properties

| In this section... |
|---|
| “Data Type and Scaling Properties” on page 5-5 |
| “How Properties are Related” on page 5-7 |
| “Set numerictype Object Properties” on page 5-8 |

Data Type and Scaling Properties

All properties of a numerictype object are writable. However, the numerictype properties of a fi object become read only after the fi object has been created. Any numerictype properties of a fi object that are unspecified at the time of fi object creation are automatically set to their default values. The properties of a numerictype object are:

| Property | Description | Valid Values |
|--------------|---|---|
| Bias | Bias associated with the object. Along with the slope, the bias forms the scaling of a fixed-point number. | <ul style="list-style-type: none"> Any floating-point number |
| DataType | Data type category | <ul style="list-style-type: none"> Fixed (default) — Fixed-point or integer data type boolean — Built-in MATLAB boolean data type double — Built-in MATLAB double data type ScaledDouble — Scaled double data type single — Built-in MATLAB single data type |
| DataTypeMode | Data type and scaling associated with the object | <ul style="list-style-type: none"> Fixed-point: binary point scaling (default) — Fixed-point data type and scaling defined by the word length and fraction length Boolean — Built-in boolean Double — Built-in double Fixed-point: slope and bias scaling — Fixed-point data type and scaling defined by the slope and bias Fixed-point: unspecified scaling — Fixed-point data type with unspecified scaling Scaled double: binary point scaling — Double data type with fixed-point word length and fraction length information retained Scaled double: slope and bias scaling — Double data type with fixed-point slope and bias information retained Scaled double: unspecified scaling — Double data type with unspecified fixed-point scaling Single — Built-in single |

| Property | Description | Valid Values |
|----------------|---|--|
| FixedExponent | Fixed-point exponent associated with the object | <ul style="list-style-type: none"> Any integer <p>Note The FixedExponent property is the negative of the FractionLength. Changing one property changes the other.</p> |
| FractionLength | Fraction length of the stored integer value, in bits | <ul style="list-style-type: none"> Best precision fraction length based on value of the object and the word length (default) Any integer <p>Note The FractionLength property is the negative of the FixedExponent. Changing one property changes the other.</p> |
| Scaling | Scaling mode of the object | <ul style="list-style-type: none"> BinaryPoint (default) — Scaling for the fi object is defined by the fraction length. SlopeBias — Scaling for the fi object is defined by the slope and bias. Unspecified — A temporary setting that is only allowed at fi object creation, to allow for the automatic assignment of a binary point best-precision scaling. |
| Signed | Whether the object is signed | <ul style="list-style-type: none"> true (default) — signed false — unsigned 1 — signed 0 — unsigned [] — auto <p>Note Although the Signed property is still supported, the Signedness property always appears in the numeric type object display. If you choose to change or set the signedness of your numeric type objects using the Signed property, MATLAB updates the corresponding value of the Signedness property.</p> |
| Signedness | Whether the object is signed, unsigned, or has an unspecified sign | <ul style="list-style-type: none"> Signed (default) Unsigned Auto — unspecified sign |
| Slope | Slope associated with the object | <ul style="list-style-type: none"> Any finite floating-point number greater than zero <p>Note</p> $\text{slope} = \text{slopeadjustmentfactor} \times 2^{\text{fixedexponent}}$ <p>Changing one of these properties changes the other.</p> |
| | Along with the bias, the slope forms the scaling of a fixed-point number. | |

| Property | Description | Valid Values |
|------------------------|---|--|
| Slope AdjustmentFactor | Slope adjustment associated with the object | <ul style="list-style-type: none"> Any number greater than or equal to 1 and less than 2 |
| | The slope adjustment is equivalent to the fractional slope of a fixed-point number. | <p>Note</p> $slope = slopeadjustmentfactor \times 2^{fixedexponent}$ <p>Changing one of these properties changes the other.</p> |
| WordLength | Word length of the stored integer value, in bits | <ul style="list-style-type: none"> 16 (default) Any positive integer if Signedness is Unsigned or unspecified Any integer greater than one if Signedness is set to Signed |

These properties are described in detail in the “Set fi Object Properties” on page 2-15. To learn how to specify properties for numerictype objects in Fixed-Point Designer software, refer to “Set numerictype Object Properties” on page 5-8.

How Properties are Related

Properties that affect the slope

The **Slope** field of the numerictype object is related to the SlopeAdjustmentFactor and FixedExponent properties by

$$slope = slopeadjustmentfactor \times 2^{fixedexponent}$$

The FixedExponent and FractionLength properties are related by

$$fixedexponent = -fractionlength$$

If you set the SlopeAdjustmentFactor, FixedExponent, or FractionLength property, the **Slope** field is modified.

Stored integer value and real world value

In binary point scaling the numerictype StoredIntegerValue and RealWorldValue properties are related according to

$$real-worldvalue = storedintegervalue \times 2^{-fractionlength}$$

In [Slope Bias] scaling the RealWorldValue can be represented by

$$real-worldvalue = storedintegervalue \times (slopeadjustmentfactor \times 2^{fixedexponent}) + bias$$

which is equivalent to

$$real-worldvalue = (slope \times storedinteger) + bias$$

If any of these properties are updated, the others are modified accordingly.

Set numerictype Object Properties

Setting numerictype Properties at Object Creation

You can set properties of numerictype objects at the time of object creation by including properties after the arguments of the numerictype constructor function.

For example, to set the word length to 32 bits and the fraction length to 30 bits,

```
T = numerictype('WordLength',32,'FractionLength',30)
```

```
T =
```

```
      DataTypeMode: Fixed-point: binary point scaling  
      Signedness: Signed  
      WordLength: 32  
      FractionLength: 30
```

In addition to creating a numerictype object at the command line, you can also set numerictype properties using the **Insert numerictype Constructor** dialog box. For an example of this approach, see “Example: Build numerictype Object Constructors in a GUI” on page 5-4.

Use Direct Property Referencing with numerictype Objects

You can reference directly into a property for setting or retrieving numerictype object property values using MATLAB structure-like referencing. You do this by using a period to index into a property by name.

For example, to get the word length of T,

```
T.WordLength
```

```
ans =
```

```
32
```

To set the fraction length of T,

```
T.FractionLength = 31
```

```
T =
```

```
      DataTypeMode: Fixed-point: binary point scaling  
      Signedness: Signed  
      WordLength: 32  
      FractionLength: 31
```


numerictype of Fixed-Point Objects

| |
|--|
| In this section... |
| “Valid Values for numerictype Object Properties” on page 5-9 |
| “Properties That Affect the Slope” on page 5-10 |
| “Stored Integer Value and Real World Value” on page 5-10 |

Valid Values for numerictype Object Properties

The numerictype object contains all the data type and scaling attributes of a fixed-point object. The numerictype object behaves like any MATLAB object, except that it only lets you set valid values for defined fields. The following table shows the possible settings of each field of the object.

Note When you create a `fi` object, any unspecified field of the numerictype object reverts to its default value. Thus, if the `DataTypeMode` is set to `unspecified` `scaling`, it defaults to `binary point scaling` when the `fi` object is created. If the `Signedness` property of the numerictype object is set to `Auto`, it defaults to `Signed` when the `fi` object is created.

| DataTypeMode | DataType | Scaling | Signedness | Word-Length | Fraction-Length | Slope | Bias |
|---|-----------------|----------------|----------------------------|--|------------------------------------|---|--------------------------------------|
| <i>Fixed-point data types</i> | | | | | | | |
| Fixed-point: binary point scaling | Fixed | BinaryPoint | Signed Unsigned Auto | Positive integer from 1 to 65,535 | Positive or negative integer | 2 ^(- fraction length) | 0 |
| Fixed-point: slope and bias scaling | Fixed | SlopeBias | Signed Unsigned Auto | Positive integer from 1 to 65,535 | N/A | Any floating- point number greater than zero | Any floating - point number |
| Fixed-point: unspecified scaling | Fixed | Unspecified | Signed Unsigned Auto | Positive integer from 1 to 65,535 | N/A | N/A | N/A |
| <i>Scaled double data types</i> | | | | | | | |
| Scaled double: binary point scaling | ScaledDouble | BinaryPoint | Signed Unsigned Auto | Positive integer from 1 to 65,535 | Positive or negative integer | 2 ^(- fraction length) | 0 |

| DataTypeMode | DataType | Scaling | Signedness | Word-Length | Fraction-Length | Slope | Bias |
|---------------------------------------|--------------|-------------|----------------------------|-----------------------------------|-----------------|---|---------------------------|
| Scaled double: slope and bias scaling | ScaledDouble | SlopeBias | Signed Unsigned Auto | Positive integer from 1 to 65,535 | N/A | Any floating-point number greater than zero | Any floating-point number |
| Scaled double: unspecified scaling | ScaledDouble | Unspecified | Signed Unsigned Auto | Positive integer from 1 to 65,535 | N/A | N/A | N/A |
| <i>Built-in data types</i> | | | | | | | |
| Double | double | N/A | 1 true | 64 | 0 | 1 | 0 |
| Single | single | N/A | 1 true | 32 | 0 | 1 | 0 |
| Boolean | boolean | N/A | 0 false | 1 | 0 | 1 | 0 |

You cannot change the numerictype properties of a fi object after fi object creation.

Properties That Affect the Slope

The **Slope** field of the numerictype object is related to the SlopeAdjustmentFactor and FixedExponent properties by

$$\text{slope} = \text{slopeadjustmentfactor} \times 2^{\text{fixedexponent}}$$

The FixedExponent and FractionLength properties are related by

$$\text{fixedexponent} = -\text{fractionlength}$$

If you set the SlopeAdjustmentFactor, FixedExponent, or FractionLength property, the **Slope** field is modified.

Stored Integer Value and Real World Value

In binary point scaling the numerictype StoredIntegerValue and RealWorldValue properties are related according to

$$\text{real-worldvalue} = \text{storedintegervalue} \times 2^{-\text{fractionlength}}$$

In [Slope Bias] scaling the RealWorldValue can be represented by

$$\text{real-worldvalue} = \text{storedintegervalue} \times (\text{slopeadjustmentfactor} \times 2^{\text{fixedexponent}}) + \text{bias}$$

which is equivalent to

$$\text{real-worldvalue} = (\text{slope} \times \text{storedinteger}) + \text{bias}$$

If any of these properties are updated, the others are modified accordingly.

See Also

More About

- “numeric type Object Properties” on page 5-5
- “Scaling” on page 1-3

numerictype Objects Usage to Share Data Type and Scaling Settings of fi objects

You can use a numerictype object to define common data type and scaling rules that you would like to use for many fi objects. You can then create multiple fi objects, using the same numerictype object for each.

Example 1

In the following example, you create a numerictype object T with word length 32 and fraction length 28. Next, to ensure that your fi objects have the same numerictype attributes, create fi objects a and b using your numerictype object T.

```
format long g
T = numerictype('WordLength',32,'FractionLength',28)

T =

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 32
    FractionLength: 28

a = fi(pi,T)

a =

    3.1415926553309

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 32
    FractionLength: 28

b = fi(pi/2,T)

b =

    1.5707963258028

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 32
    FractionLength: 28
```

Example 2

In this example, start by creating a numerictype object T with [Slope Bias] scaling. Next, use that object to create two fi objects, c and d with the same numerictype attributes:

```
T = numerictype('Scaling','slopebias','Slope',2^2,'Bias',0)
```

```

T =
    DataTypeMode: Fixed-point: slope and bias scaling
    Signedness: Signed
    WordLength: 16
    Slope: 2^2
    Bias: 0

c = fi(pi,T)
c =
    4
    DataTypeMode: Fixed-point: slope and bias scaling
    Signedness: Signed
    WordLength: 16
    Slope: 2^2
    Bias: 0

d = fi(pi/2,T)
d =
    0
    DataTypeMode: Fixed-point: slope and bias scaling
    Signedness: Signed
    WordLength: 16
    Slope: 2^2
    Bias: 0

```

For more detail on the properties of numerictype objects see “numerictype Object Properties” on page 5-5.

Working with quantizer Objects

Transformations for Quantized Data

You can convert data values from numeric to hexadecimal or binary according to the specifications of a quantizer object.

- Use `num2bin` to convert data to binary
- Use `num2hex` to convert data to hexadecimal
- Use `hex2num` to convert hexadecimal data to numeric
- Use `bin2num` to convert binary data to numeric

For example,

```
q = quantizer([3 2]);  
x = [0.75  -0.25  
     0.50  -0.50  
     0.25  -0.75  
      0    -1  ];  
b = num2bin(q,x)
```

```
b = 8x3 char array  
    '011'  
    '010'  
    '001'  
    '000'  
    '111'  
    '110'  
    '101'  
    '100'
```

produces all two's complement fractional representations of 3-bit fixed-point numbers.

See Also

`num2bin` | `num2hex` | `hex2num` | `bin2num`

Automated Fixed-Point Conversion

- “Fixed-Point Conversion Workflows” on page 7-2
- “Automated Fixed-Point Conversion” on page 7-4
- “Debug Numerical Issues in Fixed-Point Conversion Using Variable Logging” on page 7-23
- “MATLAB Language Features Supported for Automated Fixed-Point Conversion” on page 7-35
- “Generated Fixed-Point Code” on page 7-37
- “Fixed-Point Code for MATLAB Classes” on page 7-42
- “Automated Fixed-Point Conversion Best Practices” on page 7-44
- “Replacing Functions Using Lookup Table Approximations” on page 7-50
- “Custom Plot Functions” on page 7-51
- “Generate Fixed-Point MATLAB Code for Multiple Entry-Point Functions” on page 7-52
- “Convert Code Containing Global Data to Fixed Point” on page 7-56
- “Convert Code Containing Global Variables to Fixed-Point” on page 7-60
- “Convert Code Containing Structures to Fixed Point” on page 7-64
- “Convert Identical Functions Called with Different Data Types” on page 7-67
- “Data Type Issues in Generated Code” on page 7-71
- “System Objects Supported by Fixed-Point Converter App” on page 7-73
- “Convert dsp.FIRFilter Object to Fixed-Point Using the Fixed-Point Converter App” on page 7-74

Fixed-Point Conversion Workflows

In this section...

“Choosing a Conversion Workflow” on page 7-2

“Automated Workflow” on page 7-2

“Manual Workflow” on page 7-2

Choosing a Conversion Workflow

MathWorks® provides a number of solutions for fixed-point conversion. Which conversion method you use depends on your end goal and your level of fixed-point expertise.

| Goal | Conversion Method | See Also |
|---|---|--|
| Use generated fixed-point MATLAB code for simulation purposes. | If you are new to fixed-point modeling, use the Fixed-Point Converter app. | “Automated Workflow” on page 7-2 |
| | If you are familiar with fixed-point modeling, and want to quickly explore design tradeoffs, convert your code manually. | “Manual Workflow” on page 7-2 |
| Generate fixed-point C code (requires MATLAB Coder™) | MATLAB Coder Fixed-Point Conversion tool | “Convert MATLAB Code to Fixed-Point C Code” (MATLAB Coder) |
| Generated HDL code (requires HDL Coder™) | HDL Coder Workflow Advisor | “Floating-Point to Fixed-Point Conversion” (HDL Coder) |
| Integrate fixed-point MATLAB code in larger applications for system-level simulation. | Generate a MEX function from the fixed-point algorithm and call the MEX function instead of the original MATLAB function. | “Propose Data Types Based on Simulation Ranges” on page 8-13 and “Propose Data Types Based on Derived Ranges” on page 8-24 |

Automated Workflow

If you are new to fixed-point modeling and you are looking for a direct path from floating-point MATLAB to fixed-point MATLAB code, use the automated workflow. Using this automated workflow, you can obtain data type proposals based on simulation ranges, static ranges, or both. For more information, see “Automated Fixed-Point Conversion” on page 7-4, “Propose Data Types Based on Simulation Ranges” on page 8-13, and “Propose Data Types Based on Derived Ranges” on page 8-24.

Manual Workflow

If you have a baseline understanding of fixed-point implementation details and an interest in exploring design tradeoffs to achieve optimized results, use the separate algorithm/data type workflow. Separating algorithmic code from data type specifications allows you to quickly explore design tradeoffs. This approach provides readable, portable fixed-point code that you can easily integrate into other projects. For more information, see “Manual Fixed-Point Conversion Workflow”

on page 11-2 and “Implement FIR Filter Algorithm for Floating-Point and Fixed-Point Types Using cast and zeros” on page 54-136.

Automated Fixed-Point Conversion

In this section...

“Automated Fixed-Point Conversion Capabilities” on page 7-4
“Code Coverage” on page 7-5
“Proposing Data Types” on page 7-7
“Locking Proposed Data Types” on page 7-10
“Viewing Functions” on page 7-10
“Viewing Variables” on page 7-17
“Log Data for Histogram” on page 7-19
“Function Replacements” on page 7-21
“Validating Types” on page 7-21
“Testing Numerics” on page 7-22
“Detecting Overflows” on page 7-22

Automated Fixed-Point Conversion Capabilities

You can convert floating-point MATLAB code to fixed-point code using the Fixed-Point Converter app or at the command line using the `fiaccel` function `-float2fixed` option. You can choose to propose data types based on simulation range data, derived (also known as static) range data, or both.

You can manually enter static ranges. These manually entered ranges take precedence over simulation ranges and the app uses them when proposing data types. In addition, you can modify and lock the proposed type so that the app cannot change it. For more information, see “Locking Proposed Data Types” on page 7-10.

For a list of supported MATLAB features and functions, see “MATLAB Language Features Supported for Automated Fixed-Point Conversion” on page 7-35.

During fixed-point conversion, you can:

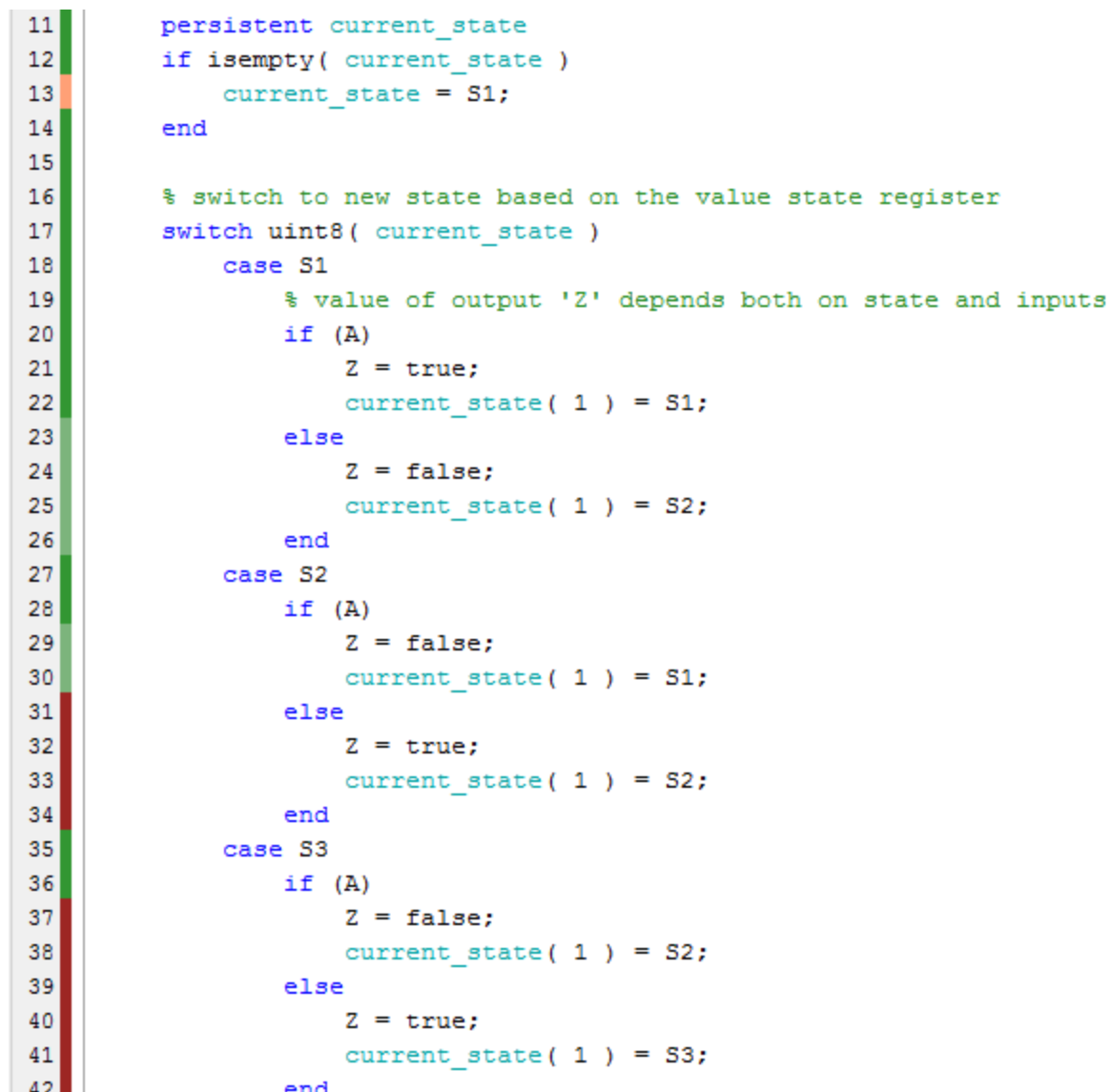
- Verify that your test files cover the full intended operating range of your algorithm using code coverage results.
- Propose fraction lengths based on default word lengths.
- Propose word lengths based on default fraction lengths.
- Optimize whole numbers.
- Specify safety margins for simulation min/max data.
- Validate that you can build your project with the proposed data types.
- Test numerics by running the test file with the fixed-point types applied.
- View a histogram of bits that each variable uses.
- Detect overflows.

Code Coverage

By default, the app shows code coverage results. Your test files must exercise the algorithm over its full operating range so that the simulation ranges are accurate. The quality of the proposed fixed-point data types depends on how well the test files cover the operating range of the algorithm with the accuracy that you want.

Reviewing code coverage results helps you to verify that your test files are exercising the algorithm adequately. If the code coverage is inadequate, modify the test files or add more test files to increase coverage. If you simulate multiple test files in one run, the app displays cumulative coverage. However, if you specify multiple test files, but run them one at a time, the app displays the coverage of the file that ran last.

The app displays a color-coded coverage bar to the left of the code.



```

11 persistent current_state
12 if isempty( current_state )
13     current_state = S1;
14 end
15
16 % switch to new state based on the value state register
17 switch uint8( current_state )
18     case S1
19         % value of output 'Z' depends both on state and inputs
20         if (A)
21             Z = true;
22             current_state( 1 ) = S1;
23         else
24             Z = false;
25             current_state( 1 ) = S2;
26         end
27     case S2
28         if (A)
29             Z = false;
30             current_state( 1 ) = S1;
31         else
32             Z = true;
33             current_state( 1 ) = S2;
34         end
35     case S3
36         if (A)
37             Z = false;
38             current_state( 1 ) = S2;
39         else
40             Z = true;
41             current_state( 1 ) = S3;
42         end

```

This table describes the color coding.

| Coverage Bar Color | Indicates |
|--------------------|---|
| Green | One of the following situations: <ul style="list-style-type: none"> The entry-point function executes multiple times and the code executes more than one time. The entry-point function executes one time and the code executes one time. Different shades of green indicate different ranges of line execution counts. The darkest shade of green indicates the highest range. |
| Orange | The entry-point function executes multiple times, but the code executes one time. |
| Red | Code does not execute. |


When you place your cursor over the coverage bar, the color highlighting extends over the code. For each section of code, the app displays the number of times that the section executes.

| | | |
|----|---|----------|
| 11 | <code>persistent current_state</code> | |
| 12 | <code>if isempty(current_state)</code> | |
| 13 | <code>current_state = S1;</code> | 1 calls |
| 14 | <code>end</code> | 51 calls |
| 15 | | |
| 16 | <code>% switch to new state based on the value state register</code> | |
| 17 | <code>switch uint8(current_state)</code> | |
| 18 | <code>case S1</code> | |
| 19 | <code> % value of output 'Z' depends both on state and inputs</code> | |
| 20 | <code> if (A)</code> | |
| 21 | <code> Z = true;</code> | 37 calls |
| 22 | <code> current_state(1) = S1;</code> | |
| 23 | <code> else</code> | 7 calls |
| 24 | <code> Z = false;</code> | |
| 25 | <code> current_state(1) = S2;</code> | |
| 26 | <code> end</code> | |
| 27 | <code>case S2</code> | 51 calls |
| 28 | <code> if (A)</code> | |
| 29 | <code> Z = false;</code> | 7 calls |
| 30 | <code> current_state(1) = S1;</code> | |
| 31 | <code> else</code> | 0 calls |
| 32 | <code> Z = true;</code> | |
| 33 | <code> current_state(1) = S2;</code> | |
| 34 | <code> end</code> | |
| 35 | <code>case S3</code> | 51 calls |
| 36 | <code> if (A)</code> | |
| 37 | <code> Z = false;</code> | 0 calls |
| 38 | <code> current_state(1) = S2;</code> | |
| 39 | <code> else</code> | |
| 40 | <code> Z = true;</code> | |
| 41 | <code> current_state(1) = S3;</code> | |
| 42 | <code> end</code> | |

To verify that your test files are testing your algorithm over the intended operating range, review the code coverage results.

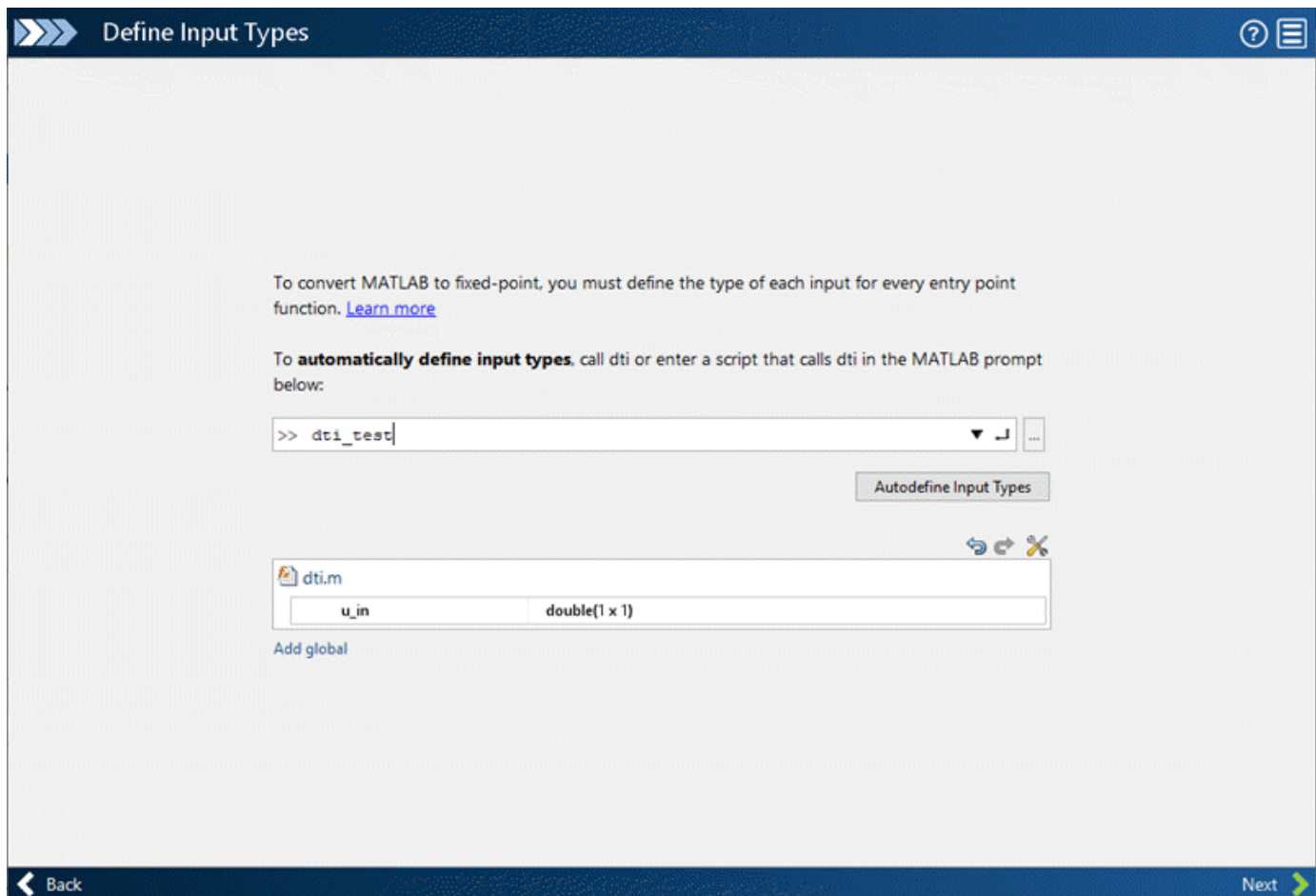
| Coverage Bar Color | Action |
|--------------------|--|
| Green | If you expect sections of code to execute more frequently than the coverage shows, either modify the MATLAB code or the test files. |
| Orange | This behavior is expected for initialization code, for example, the initialization of persistent variables. If you expect the code to execute more than one time, either modify the MATLAB code or the test files. |
| Red | If the code that does not execute is an error condition, this behavior is acceptable. If you expect the code to execute, either modify the MATLAB code or the test files. If the code is written conservatively and has upper and lower boundary limits, and you cannot modify the test files to reach this code, add static minimum and maximum values. See “Computing Derived Ranges” on page 7-9. |

Code coverage is on by default. Turn it off only after you have verified that you have adequate test file coverage. Turning off code coverage can speed up simulation. To turn off code coverage, on the **Convert to Fixed Point** page:

- 1 Click the **Analyze** arrow .
- 2 Clear the **Show code coverage** check box.

Proposing Data Types

In the define input types step, you specify a test file that calls the entry-point function. The app runs the test file to analyze the code and infer the types for entry-point input arguments.



The app proposes fixed-point data types based on computed ranges and the word length or fraction length setting. The computed ranges are based on simulation range data, derived range data (also known as static ranges), or both. If you run a simulation and compute derived ranges, the app merges the simulation and derived ranges.

Note You cannot propose data types based on derived ranges for MATLAB classes.

Derived range analysis is not supported for non-scalar variables.

You can manually enter static ranges. These manually entered ranges take precedence over simulation ranges and the app uses them when proposing data types. If you analyze ranges using derived range analysis alone, you must enter static ranges.

You can modify and lock the proposed type so that the tool cannot change it. For more information, see “Locking Proposed Data Types” on page 7-10.

Running a Simulation

During fixed-point conversion, the app generates an instrumented MEX function for your entry-point MATLAB file. If the build completes without errors, the app displays compiled information (type, size, complexity) for functions and variables in your code. To navigate to local functions, click the

Functions tab. If build errors occur, the app provides error messages that link to the line of code that caused the build issues. You must address these errors before running a simulation. Use the link to navigate to the offending line of code in the MATLAB editor and modify the code to fix the issue. If your code uses functions that are not supported for fixed-point conversion, the app displays them on the **Function Replacements** tab. See “Function Replacements” on page 7-21.

Before running a simulation, specify the test file or files that you want to run. When you run a simulation, the app runs the test file, calling the instrumented MEX function. If you modify the MATLAB design code, the app automatically generates an updated MEX function before running a test file.

If the test file runs successfully, the simulation minimum and maximum values and the proposed types are displayed on the **Variables** tab. If you manually enter static ranges for a variable, the manually entered ranges take precedence over the simulation ranges. If you manually modify the proposed types by typing or using the histogram, the data types are locked so that the app cannot modify them.

If the test file fails, the errors are displayed on the **Output** tab.

Test files must exercise your algorithm over its full operating range. The quality of the proposed fixed-point data types depends on how well the test file covers the operating range of the algorithm with the accuracy that you want. You can add test files and select to run more than one test file during the simulation. If you run multiple test files, the app merges the simulation results.

Optionally, you can select to log data for histograms. After running a simulation, you can view the histogram for each variable. For more information, see “Log Data for Histogram” on page 7-19.

Computing Derived Ranges

The advantage of proposing data types based on derived ranges is that you do not have to provide test files that exercise your algorithm over its full operating range. Running such test files often takes a very long time. The app can compute derived ranges for scalar variables only.

To compute derived ranges and propose data types based on these ranges, provide static minimum and maximum values or proposed data types for all input variables. To improve the analysis, enter as much static range information as possible for other variables. You can manually enter ranges or promote simulation ranges to use as static ranges. Manually entered static ranges always take precedence over simulation ranges.

If you know what data type your hardware target uses, set the proposed data types to match this type. Manually entered data types are locked so that the app cannot modify them. The app uses these data types to calculate the input minimum and maximum values and to derive ranges for other variables. For more information, see “Locking Proposed Data Types” on page 7-10.

When you select **Compute Derived Ranges**, the app runs a derived range analysis to compute static ranges for variables in your MATLAB algorithm. When the analysis is complete, the static ranges are displayed on the **Variables** tab. If the run produces $+/-\text{Inf}$ derived ranges, consider defining ranges for all persistent variables.

Optionally, you can select **Quick derived range analysis**. With this option, the app performs faster static analysis. The computed ranges might be larger than necessary. Select this option in cases where the static analysis takes more time than you can afford.

If the derived range analysis for your project is taking a long time, you can optionally set a timeout. When the timeout is reached, the app aborts the analysis.

Locking Proposed Data Types

You can lock proposed data types against changes by the app using one of the following methods:

- Manually setting a proposed data type in the app.
- Right-clicking a type proposed by the tool and selecting `Lock computed value`.

The app displays locked data types in bold so that they are easy to identify. You can unlock a type using one of the following methods:

- Manually overwriting it.
- Right-clicking it and selecting `Undo changes`. This action unlocks only the selected type.
- Right-clicking and selecting `Undo changes for all variables`. This action unlocks all locked proposed types.

Viewing Functions

During the **Convert to Fixed Point** step of the fixed-point conversion process, you can view a list of functions in your project in the left pane. This list also includes function specializations and class methods. When you select a function from the list, the MATLAB code for that function or class method is displayed in the code window and the variables that they use are displayed on the **Variables** tab.

After conversion, the left pane also displays a list of output files including the fixed-point version of the original algorithm. If your function is not specialized, the app retains the original function name in the fixed-point file name and appends the fixed-point suffix. For example, here the fixed-point version of `ex_2ndOrder_filter.m` is `ex_2ndOrder_filter_fixpt.m`.

The screenshot shows the 'Fixed-Point Converter' application window. The main area displays MATLAB code for a Butterworth filter. The code defines a function `ex_2ndOrder_Filter_fixpt(x)` that takes an input `x` and returns an output `y`. The code uses the `fi` function to convert floating-point values to fixed-point format. The filter coefficients are defined as `b` and `a`. The filter is implemented as a second-order Butterworth filter with a cutoff frequency of 0.25. The code uses the `fi_signed` function to handle signed fixed-point numbers.

The 'Output Files' section shows the generated files, including `ex_2ndOrder_filter_fixpt.m`, `ex_2ndOrder_filter_wrapper_fixpt.r`, `ex_2ndOrder_filter_report.html`, `index.html`, `ex_2ndOrder_filter_fixpt_args.mat`, and `ex_2ndOrder_filter_wrapper_fixpt_`.

The 'Variables' tab shows the conversion results for the input, output, and persistent variables. The table below summarizes the data:

| Variable | Type | Size | Signed | Word Length | Fraction Length |
|-------------------|-------------|---------|--------|-------------|-----------------|
| Input | | | | | |
| x | embedded.fi | 1 x 256 | Yes | 16 | 14 |
| Output | | | | | |
| y | embedded.fi | 1 x 256 | Yes | 16 | 14 |
| Persistent | | | | | |
| z | embedded.fi | 2 x 1 | Yes | 16 | 15 |
| Local | | | | | |

A 'Validation succeeded' message is displayed at the bottom of the application.

Classes

The app displays information for the class and each of its methods. For example, consider a class, `Counter`, that has a static method, `MAX_VALUE`, and a method, `next`.

If you select the class, the app displays the class and its properties on the **Variables** tab.

The screenshot displays a MATLAB class definition for `Counter` with the following code:

```

1 classdef Counter < handle
2     properties
3         Value
4     end
5
6     methods (Static)
7         function t = MAX_VALUE()
8             t = 128;
9         end
10    end
11
12    methods
13        function this = Counter()
14            this.Value = 0;
15        end
16
17        function v = next(this)
18            v = this.Value;
19            if this.Value == this.MAX_VALUE
20                this.Value = 0;
21            else
22                this.Value = this.Value + 1;
23            end
24        end
25    end
26 end
27

```

The 'Type Validation Output' table is shown below the code editor:

| Variable | Type | Sim Min | Sim Max | Static Min | Static Max | Whole Nu... | Proposed Type |
|----------|--------|---------|---------|------------|------------|-------------|----------------------|
| ▲ Output | | | | | | | |
| ▲ this | | | | | | | |
| Value | double | 0 | 128 | | | Yes | numerictype(0, 8, 0) |

If you select a method, the app displays only the variables that the method uses.

The screenshot displays the MATLAB IDE interface. On the left, the 'Source Code' pane shows a file explorer with 'Counter > Counter' selected. The main editor window shows the following MATLAB code:

```

1 classdef Counter < handle
2     properties
3         Value
4     end
5
6     methods (Static)
7         function t = MAX_VALUE()
8             t = 128;
9         end
10    end
11
12    methods
13        function this = Counter()
14            this.Value = 0;
15        end
16
17        function v = next(this)
18            v = this.Value;
19            if this.Value == this.MAX_VALUE
20                this.Value = 0;
21            else
22                this.Value = this.Value + 1;
23            end
24        end
25    end
26 end
27

```

Below the code editor, the 'Type Validation Output' tab is active, showing a table with the following data:

| Variable | Type | Sim Min | Sim Max | Static Min | Static Max | Whole Nu... | Proposed Type |
|----------|--------|---------|---------|------------|------------|-------------|----------------------|
| t | double | 128 | 128 | | | Yes | numerictype(0, 8, 0) |

At the bottom of the IDE, there are 'Back' and 'Next' navigation buttons.

Specializations

If a function is specialized, the app lists each specialization and numbers them sequentially. For example, consider a function, `dut`, that calls subfunctions, `foo` and `bar`, multiple times with different input types.

```

function y = dut(u, v)

tt1 = foo(u);
tt2 = foo([u v]);
tt3 = foo(complex(u,v));

ss1 = bar(u);
ss2 = bar([u v]);
ss3 = bar(complex(u,v));

y = (tt1 + ss1) + sum(tt2 + ss2) + real(tt3) + real(ss3);

end

function y = foo(u)
    y = u * 2;
end

function y = bar(u)

```

```

y = u * 4;
end

```

If you select the top-level function, the app displays all the variables on the **Variables** tab.

The screenshot shows the MATLAB IDE interface. On the left, the 'Source Code' pane lists several functions: 'dut', 'dut > bar > 1', 'dut > bar > 2', 'dut > bar > 3', 'dut > foo > 1', 'dut > foo > 2', and 'dut > foo > 3'. The main editor displays the source code for three functions: 'dut(u, v)', 'foo(u)', and 'bar(u)'. Below the code, the 'Variables' tab is active, showing a table of variables with their types, simulation ranges, and proposed fixed-point types.

| Variable | Type | Sim Min | Sim Max | Whole Number | Proposed Type |
|---------------|--------|---------|---------|--------------|----------------------|
| Input | | | | | |
| u | double | 10 | 10 | Yes | numerictype(0, 4, 0) |
| v | double | 20 | 20 | Yes | numerictype(0, 5, 0) |
| Output | | | | | |
| y | double | 300 | 300 | Yes | numerictype(0, 9, 0) |
| Local | | | | | |
| tt1 | double | 20 | 20 | Yes | numerictype(0, 5, 0) |

If you select the tree view, the app also displays the line numbers for the call to each specialization.

The screenshot displays a software interface for automated fixed-point conversion. On the left, a tree view under 'Source Code' shows a hierarchy of functions: 'dut' (expanded) and its sub-functions 'foo > 3', 'foo > 1', 'foo > 2', 'bar > 1', 'bar > 2', and 'bar > 3'. The main area shows the source code for these functions. Below the code is a table with tabs for 'Variables', 'Function Replacements', and 'Output'. The 'Variables' tab is active, showing a table with columns: Variable, Type, Sim Min, Sim Max, Whole Number, and Proposed Type.

```

1 function y = dut(u, v)
2
3     tt1 = foo(u);
4     tt2 = foo([u v]);
5     tt3 = foo(complex(u,v));
6
7
8     ss1 = bar(u);
9     ss2 = bar([u v]);
10    ss3 = bar(complex(u,v));
11
12    y = (tt1 + ss1) + sum(tt2 + ss2) + real(tt3) + real(ss3);
13 end
14
15 function y = foo(u)
16     y = u * 2;
17 end
18
19 function y = bar(u)
20     y = u * 4;
21 end

```

| Variable | Type | Sim Min | Sim Max | Whole Number | Proposed Type |
|----------|----------------|---------|---------|--------------|----------------------|
| Input | | | | | |
| u | complex double | 10 | 20 | Yes | numerictype(0, 5, 0) |
| Output | | | | | |
| y | complex double | 20 | 40 | Yes | numerictype(0, 6, 0) |

If you select a specialization, the app displays only the variables that the specialization uses.

The screenshot displays the MATLAB IDE interface. On the left, a 'Source Code' pane lists several function calls: `dut`, `dut > bar > 1`, `dut > bar > 2`, `dut > bar > 3`, `dut > foo > 1`, `dut > foo > 2`, and `dut > foo > 3`. The main editor shows the source code for three functions:

```

1 function y = dut(u, v)
2
3     tt1 = foo(u);
4     tt2 = foo([u v]);
5     tt3 = foo(complex(u,v));
6
7     ss1 = bar(u);
8     ss2 = bar([u v]);
9     ss3 = bar(complex(u,v));
10
11    y = (tt1 + ss1) + sum(tt2 + ss2) + real(tt3) + real(ss3);
12
13    end
14
15 function y = foo(u)
16     y = u * 2;
17    end
18
19 function y = bar(u)
20     y = u * 4;
21    end

```

Below the code editor, there is a table with tabs for 'Variables', 'Function Replacements', and 'Output'. The 'Output' tab is active, showing the following data:

| Variable | Type | Sim Min | Sim Max | Whole Number | Proposed Type |
|----------|----------------|---------|---------|--------------|----------------------|
| Input | | | | | |
| u | complex double | 10 | 20 | Yes | numerictype(0, 5, 0) |
| Output | | | | | |
| y | complex double | 20 | 40 | Yes | numerictype(0, 6, 0) |

In the generated fixed-point code, the number of each fixed-point specialization matches the number in the **Source Code** list, which makes it easy to trace between the floating-point and fixed-point versions of your code. For example, the generated fixed-point function for `foo > 1` is named `foo_s1`.

The screenshot displays the MATLAB IDE interface. On the left, the 'Source Code' pane shows the function definition for 'dut_fixpt' and its sub-functions 'foo_s1' and 'foo_s2'. The 'Output Files' pane lists generated files like 'dut_fixpt.m' and 'dut_report.html'. The main editor shows the MATLAB code for the function, which uses 'fi' blocks for fixed-point operations. Below the code, the 'Variables' tab is active, showing a table of variables and their properties.

| Variable | Type | Sim Min | Sim Max | Whole Number | Proposed Type |
|---------------|----------------|---------|---------|--------------|----------------------|
| Input | | | | | |
| u | complex double | 10 | 20 | Yes | numerictype(0, 5, 0) |
| Output | | | | | |
| y | complex double | 40 | 80 | Yes | numerictype(0, 7, 0) |

At the bottom of the IDE, a green notification box indicates 'Validation succeeded'.

Viewing Variables

The **Variables** tab provides the following information for each variable in the function selected in the **Navigation** pane:

- **Type** — The original data type of the variable in the MATLAB algorithm.
- **Sim Min** and **Sim Max** — The minimum and maximum values assigned to the variable during simulation.

You can edit the simulation minimum and maximum values. Edited fields are shown in bold. Editing these fields does not trigger static range analysis, but the tool uses the edited values in subsequent analyses. You can revert to the types proposed by the app.

- **Static Min** and **Static Max** — The static minimum and maximum values.

To compute derived ranges and propose data types based on these ranges, provide static minimum and maximum values for all input variables. To improve the analysis, enter as much static range information as possible for other variables.

When you compute derived ranges, the app runs a static analysis to compute static ranges for variables in your code. When the analysis is complete, the static ranges are displayed. You can edit the computed results. Edited fields are shown in bold. Editing these fields does not trigger static range analysis, but the tool uses the edited values in subsequent analyses. You can revert to the types proposed by the app.

- **Whole Number** — Whether all values assigned to the variable during simulation are integers.

The app determines whether a variable is always a whole number. You can modify this field. Edited fields are shown in bold. Editing these fields does not trigger static range analysis, but the app uses the edited values in subsequent analyses. You can revert to the types proposed by the app.

- The proposed fixed-point data type for the specified word (or fraction) length. Proposed data types use the `numerictype` notation. For example, `numerictype(1,16,12)` denotes a signed fixed-point type with a word length of 16 and a fraction length of 12. `numerictype(0,16,12)` denotes an unsigned fixed-point type with a word length of 16 and a fraction length of 12.

Because the app does not apply data types to expressions, it does not display proposed types for them. Instead, it displays their original data types.

You can also view and edit variable information in the code pane by placing your cursor over a variable name.

You can use `Ctrl+F` to search for variables in the MATLAB code and on the **Variables** tab. The app highlights occurrences in the code and displays only the variable with the specified name on the **Variables** tab.

Viewing Information for MATLAB Classes

The app displays:

- Code for MATLAB classes and code coverage for class methods in the code window. Use the **Source Code** list on the **Convert to Fixed Point** page to select which class or class method to view. If you select a class method, the app highlights the method in the code window.

The screenshot shows the MATLAB IDE with a class definition for `Counter` and its type validation output table. The class definition is as follows:

```

1 classdef Counter < handle
2     properties
3         Value
4     end
5
6     methods (Static)
7         function t = MAX_VALUE ()
8             t = 128;
9         end
10    end
11
12    methods
13        function this = Counter()
14            this.Value = 0;
15        end
16
17        function v = next(this)
18            v = this.Value;
19            if this.Value == this.MAX_VALUE
20                this.Value = 0;
21            else
22                this.Value = this.Value + 1;
23            end
24        end
25    end
26 end
27

```

The type validation output table is shown below:


| Variable | Type | Sim Min | Sim Max | Static Min | Static Max | Whole Nu... | Proposed Type |
|----------|---------|---------|---------|------------|------------|-------------|---------------------|
| Input | | | | | | | |
| this | Counter | Unknown | Unknown | | | | |
| Value | double | 0 | 128 | | | Yes | numericity(0, 8, 0) |
| Output | | | | | | | |
| v | double | 0 | 128 | | | Yes | numericity(0, 8, 0) |

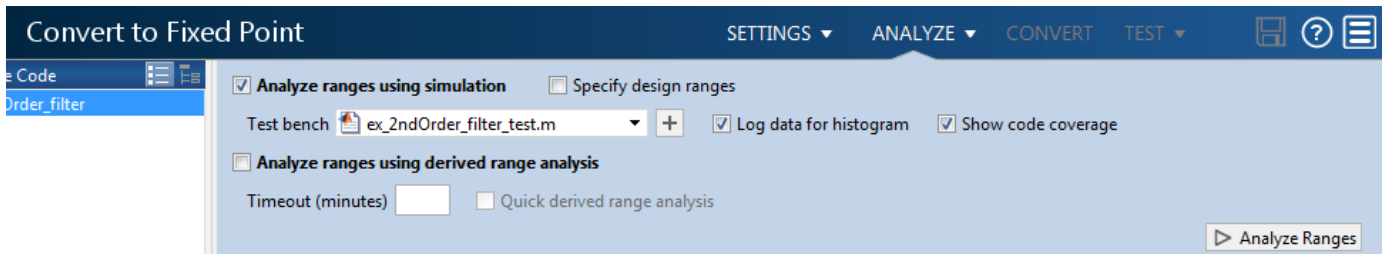
- Information about MATLAB classes on the **Variables** tab.

| Variable | Type | Sim Min | Sim Max | Static Min | Static Max | Whole Num... | Proposed Type |
|----------|---------|---------|---------|------------|------------|--------------|---------------------|
| Input | | | | | | | |
| this | Counter | Unknown | Unknown | | | | |
| Value | double | 0 | 128 | | | Yes | numericity(0, 8, 0) |
| Output | | | | | | | |
| v | double | 0 | 128 | | | Yes | numericity(0, 8, 0) |

Log Data for Histogram

To log data for histograms:

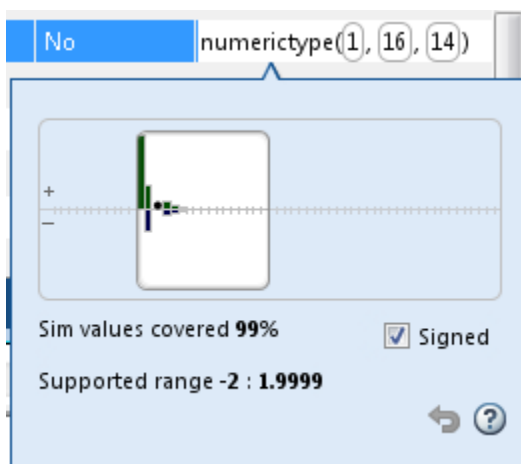
- On the **Convert to Fixed Point** page, click the **Analyze** arrow .
- Select **Log data for histogram**.



- Click **Analyze Ranges**.

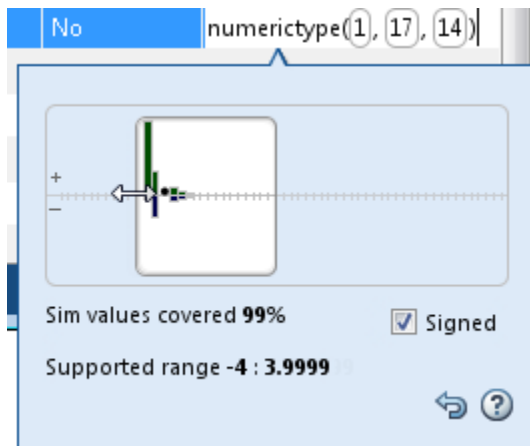
After simulation, to view the histogram for a variable, on the **Variables** tab, click the **Proposed Type** field for that variable.

The histogram provides the range of the proposed data type and the percentage of simulation values that the proposed data type covers. The bit weights are displayed along the X-axis, and the percentage of occurrences along the Y-axis. Each bin in the histogram corresponds to a bit in the binary word. For example, this histogram displays the range for a variable of type `numerictype(1, 16, 14)`.




You can view the effect of changing the proposed data types by:

- Dragging the edges of the bounding box in the histogram window to change the proposed data type.



- Selecting or clearing **Signed**.

To revert to the types proposed by the automatic conversion, in the histogram window, click .

Function Replacements

If your MATLAB code uses functions that do not have fixed-point support, the app lists these functions on the **Function Replacements** tab. You can choose to replace unsupported functions with a custom function replacement or with a lookup table.

You can add and remove function replacements from this list. If you enter a function replacement for a function, the replacement function is used when you build the project. If you do not enter a replacement, the app uses the type specified in the original MATLAB code for the function.

Note Using this table, you can replace the names of the functions but you cannot replace argument patterns.

If code generation readiness screening is disabled, the list of unsupported functions on the **Function Replacements** tab can be incomplete or incorrect. In this case, add the functions manually. See .

Validating Types

Converting the code to fixed point validates the build using the proposed fixed-point data types. If the validation is successful, you are ready to test the numerical behavior of the fixed-point MATLAB algorithm.

If the errors or warnings occur during validation, they are displayed on the **Output** tab. If errors or warning occur:

- On the **Variables** tab, inspect the proposed types and manually modified types to verify that they are valid.
- On the **Function Replacements** tab, verify that you have provided function replacements for unsupported functions.

Testing Numerics

After converting code to fixed point and validating the proposed fixed-point data types, click **Test** to verify the behavior of the fixed-point MATLAB algorithm. By default, if you added a test file to define inputs or run a simulation, the app uses this test file to test numerics. Optionally, you can add test files and select to run more than one test file. The app compares the numerical behavior of the generated fixed-point MATLAB code with the original floating-point MATLAB code. If you select to log inputs and outputs for comparison plots, the app generates an additional plot for each scalar output. This plot shows the floating-point and fixed-point results and the difference between them. For nonscalar outputs, only the error information is shown.










After fixed-point simulation, if the numerical results do not meet the accuracy that you want, modify fixed-point data type settings and repeat the type validation and numerical testing steps. You might have to iterate through these steps multiple times to achieve the results that you want.

Detecting Overflows

When testing numerics, selecting **Use scaled doubles to detect overflows** enables overflow detection. When this option is selected, the conversion app runs the simulation using scaled double versions of the proposed fixed-point types. Because scaled doubles store their data in double-precision floating-point, they carry out arithmetic in full range. They also retain their fixed-point settings, so they are able to report when a computation goes out of the range of the fixed-point type. For more information, see “Scaled Doubles” on page 35-16.

If the app detects overflows, on its **Overflow** tab, it provides:

- A list of variables and expressions that overflowed
- Information on how much each variable overflowed
- A link to the variables or expressions in the code window

| Variables | Function Replacements | Overflows | |
|---|-----------------------|-----------|--|
| | Function | Line | Description |
|  | overflow_fixpt | 7 | Overflow error in expression 'x'. |
|  | overflow_fixpt | 7 | Overflow error in expression 'y'. |
|  | overflow_fixpt | 10 | Overflow error in expression 'z'. |
|  | overflow_fixpt | 10 | Overflow error in expression 'z = fi(x*y, 0, 8, 0, fm)'. |
|  | overflow_fixpt | 10 | Overflow error in expression 'fi(x*y, 0, 8, 0, fm)'. |
|  | overflow_fixpt | 10 | Overflow error in expression 'x'. |
|  | overflow_fixpt | 10 | Overflow error in expression 'x*y'. |
|  | overflow_fixpt | 10 | Overflow error in expression 'y'. |
|  | overflow_fixpt | 11 | Overflow error in expression 'z'. |

If your original algorithm uses scaled doubles, the app also provides overflow information for these expressions.

See Also

“Detect Overflows” on page 8-5

Debug Numerical Issues in Fixed-Point Conversion Using Variable Logging

In this section...

“Prerequisites” on page 7-23
 “Convert to Fixed Point Using Default Configuration” on page 7-26
 “Determine Where Numerical Issues Originated” on page 7-29
 “Adjust fimath Settings” on page 7-30
 “Adjust Word Length Settings” on page 7-31
 “Replace Constant Functions” on page 7-32

This example shows some best practices for debugging your fixed-point code when you need more than out of the box conversion.

Prerequisites

- 1 Create a local working folder, for example, `c:\kalman_filter`.
- 2 In your local working folder, create the following files.

- **kalman_filter function**

This is the main entry-point function for your project.

```

function [y] = kalman_filter(z,N0)
    %#codegen
    A = kalman_stm();

    % Measurement Matrix
    H = [1 0];

    % Process noise variance
    Q = 0;
    % Measurement noise variance
    R = N0 ;

    persistent x_est p_est
    if isempty(x_est)
        % Estimated state
        x_est = [0; 1];
        % Estimated error covariance
        p_est = N0 * eye(2, 2);
    end

    % Kalman algorithm
    % Predicted state and covariance
    x_prd = A * x_est;
    p_prd = A * p_est * A' + Q;

    % Estimation
    S = H * p_prd' * H' + R;
    B = H * p_prd';
    klm_gain = matrix_solve(S,B)';
  
```

```
% Estimated state and covariance
x_est = x_prd + klm_gain * (z - H * x_prd);
p_est = p_prd - klm_gain * H * p_prd;

% Compute the estimated measurements
y = H * x_est;
```

end

- **kalman_stm function**

This function is called by the `kalman_filter` function and computes the state transition matrix.

```
function A = kalman_stm()
    f0 = 200;
    dt = 1/1e4;
    % Kalman filter initialization
    % State transition Matrix
    A = [cos(2*pi*f0*dt), -sin(2*pi*f0*dt);
         sin(2*pi*f0*dt), cos(2*pi*f0*dt)];
```

end

- **matrix_solve function**

This function is a more efficient implementation of a matrix left divide.

```
function x = matrix_solve(a,b)
    %fixed-point conversion friendly matrix solve: a * x = b

    % initialize x
    x = zeros(size(a,1),size(b,2));
    % compute lu decomposition of a
    [l, u] = lu_replacement(a);
    % solve x = a\b for every column of b
    % through forward and backward substitution
    for col = 1:size(b,2)
        bcol = b(:,col);
        y = forward_substitute(l,bcol);
        x(:,col) = back_substitute(u,y);
    end
```

end

- **lu_replacement function**

This function is called by the `matrix_solve` function.

```
function [l,A]=lu_replacement(A)
    N=size(A,1);
    l = eye(N);
    for n=1:N-1
        piv = A(n,n);
        for k=n+1:N
            mult = divide_no_zero(A(k,n),piv);
            A(k,:) = -mult*A(n,:) + A(k,:);
            l(k,n) = mult;
        end
    end
```

end

end

- **forward_substitute function**

This function is called by the `matrix_solve` function.

```
function y = forward_substitute(l,b)
    % forward substitution
    N = size(b,1);
    y = zeros(N,1);
    % forward substitution
    y(1) = divide_no_zero(b(1),l(1,1));
    for n = 2:N
        acc = 0;
        for k = 1:n-1
            acc(:) = acc + y(k)*l(n,k);
        end
        y(n) = divide_no_zero((b(n)-acc),l(n,n));
    end
end
```

- **back_substitute function**

This function is called by the `matrix_solve` function.

```
function x = back_substitute(u,y)
    % backwards substitution
    N = size(u,1);
    x = zeros(N,1);

    % backward substitution
    x(N) = divide_no_zero(y(N),u(N,N));

    for n = (N-1):(-1):(1)
        acc = 0;
        for k = n:(N)
            acc(:) = acc + x(k)*u(n,k);
        end
        x(n) = divide_no_zero((y(n) - acc),u(n,n));
    end
end
```

- **divide_no_zero function**

This function is called by the `lu_replacement`, `forward_substitute` and `back_substitute` functions.

```
function y = divide_no_zero(num, den)
    % Divide and avoid division by zero
    if den == 0
        y = 0;
    else
        y = num/den;
    end
end
```

- **kalman_filter_tb test file**

This script generates a noisy sine wave, and calls the `kalman_filter` function to filter the noisy signal. It then plots the signals for comparison.

```
% KALMAN FILTER EXAMPLE TEST BENCH
clear all
step = ((400*pi)/1000)/10;
TIME_STEPS = 400;
X = 0:step:TIME_STEPS;
rng default;
rng(1);
Orig_Signal = sin(X);
Noisy_Signal = Orig_Signal + randn(size(X));
Clean_Signal = zeros(size(X));
for i = 1:length(X)
Clean_Signal(i) = kalman_filter(Noisy_Signal(i), 1);
end
figure
subplot(5,1,1)
plot(X,rand(size(X)))
axis([1 TIME_STEPS 0 1.25]);
title('Noise')

% Plot Noisy Signal
subplot(5,1,2)
plot(X,Noisy_Signal)
axis([1 TIME_STEPS -4 4]);
title('Noisy Signal')

% Plot Filtered Clean Signal
subplot(5,1,3)
plot(X,Clean_Signal)
axis([1 TIME_STEPS -2 2]);
title('Filtered Signal')

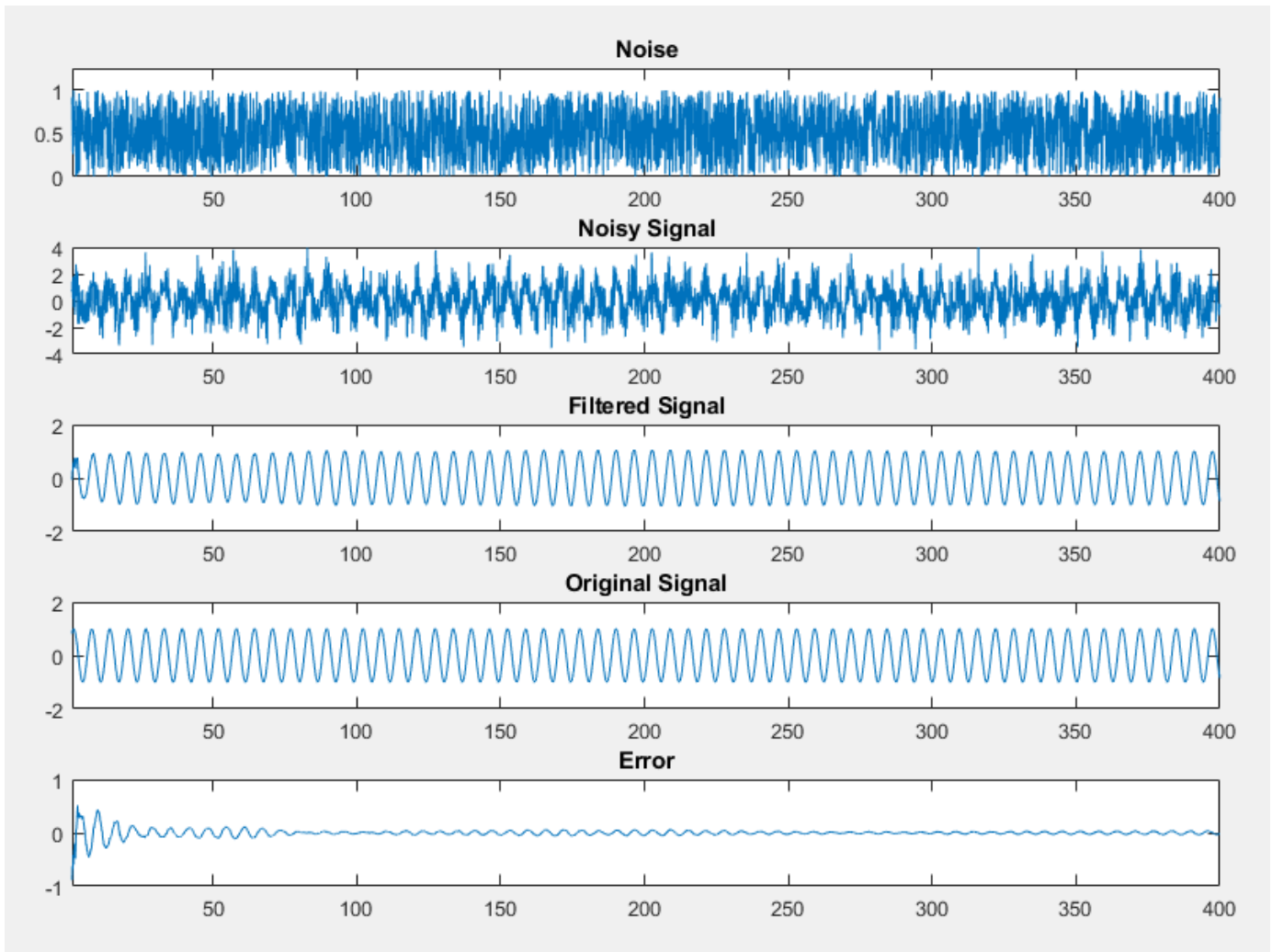
% Plot Original Signal
subplot(5,1,4)
plot(X,Orig_Signal)
axis([1 TIME_STEPS -2 2]);
title('Original Signal')

% Plot Error
subplot(5,1,5)
plot(X, (Clean_Signal - Orig_Signal))
axis([1 TIME_STEPS -1 1]);
title('Error')
figure(gcf)
```

Convert to Fixed Point Using Default Configuration

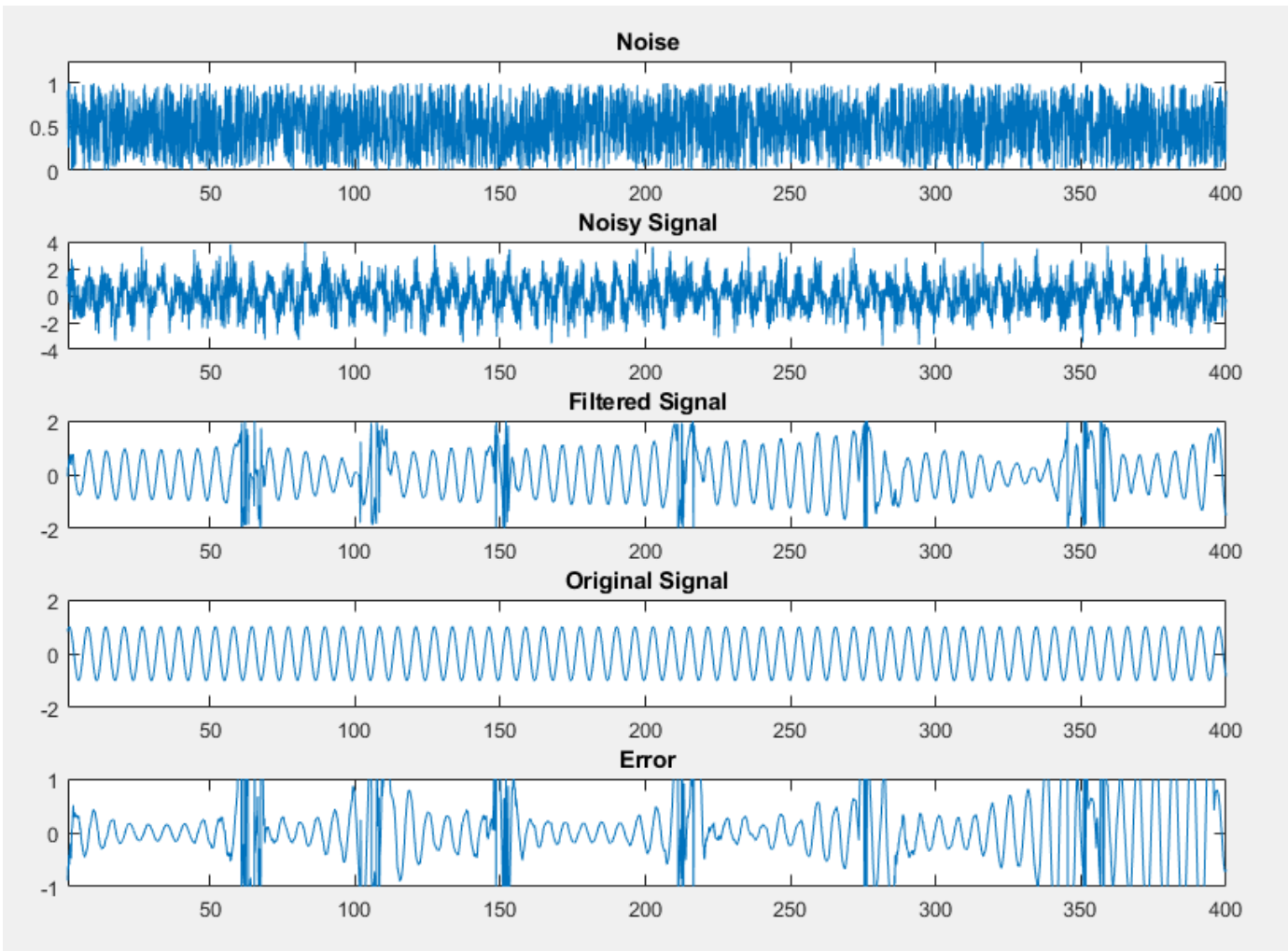
- 1 From the apps gallery, open the Fixed-Point Converter app.
- 2 On the **Select** page, browse to the `kalman_filter.m` file and click **Open**.
- 3 Click **Next**. On the **Define Input Types** page, browse to the `kalman_filter_tb` file. Click **Autodefine Input Types**.

The test file runs and plots the input noisy signal, the filtered signal, the ideal filtered signal, and the difference between the filtered and the ideal filtered signal.

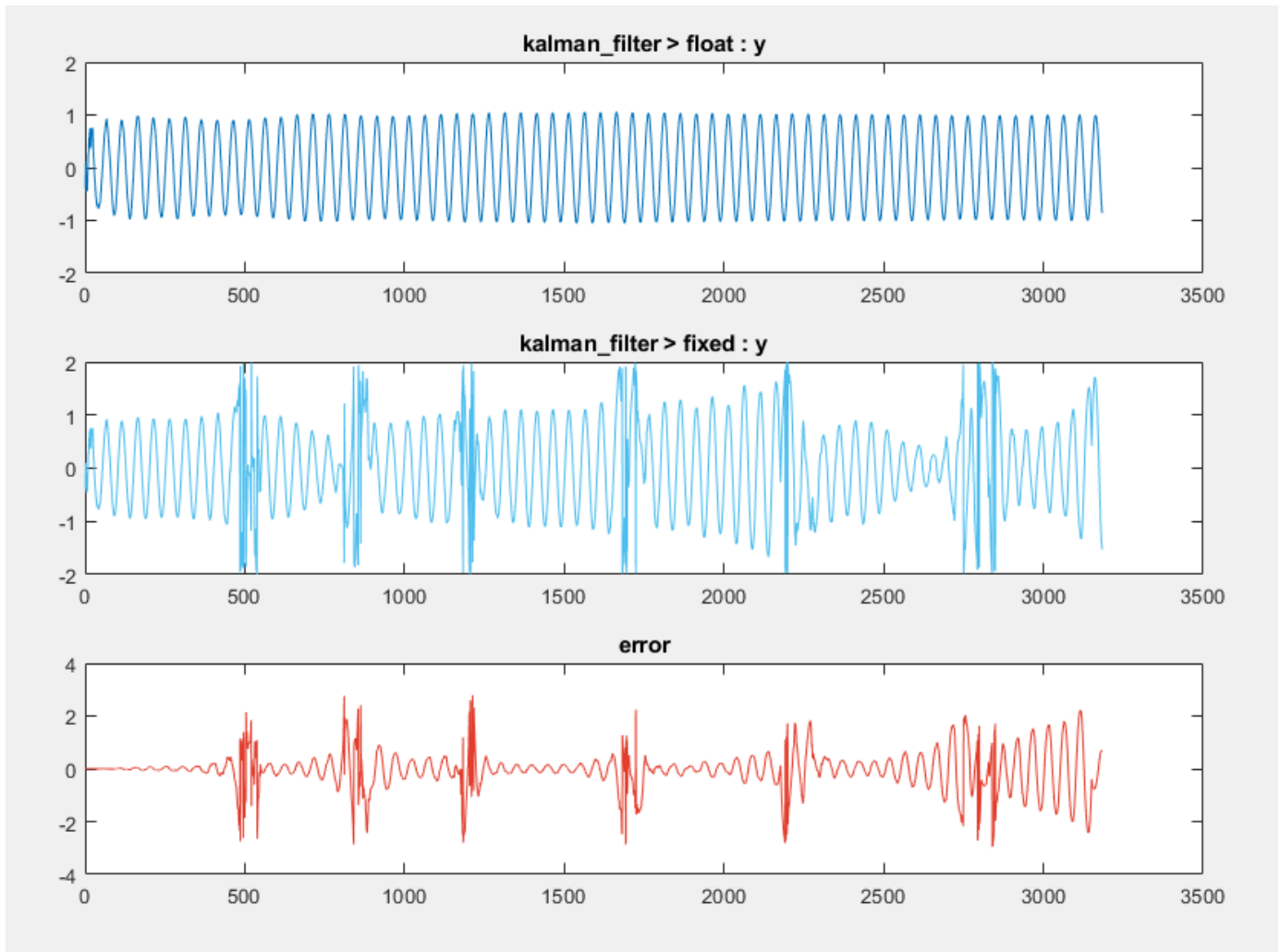


- 4** Click **Next**. On the **Convert to Fixed Point** page, click **Analyze** to gather range information and data type proposals using the default settings.
- 5** Click **Convert** to apply the proposed data types.
- 6** Click the **Test** arrow and select the **Log inputs and outputs for comparison plots** check box. Click **Test**. The Fixed-Point Converter runs the test file `kalman_filter_tb.m` to test the generated fixed-point code. Floating-point and fixed-point simulations are run, with errors calculated for the input and output variables.

The generated plots show that the current fixed-point implementation does not produce good results.



The error for the output variable y is extremely high, at over 282 percent.



Determine Where Numerical Issues Originated

Log any function input and output variables that you suspect are the cause of numerical issues to the output arguments of the top-level function.

- 1 Click `kalman_filter` in the **Source Code** pane to return to the floating-point code.

When you select the **Log inputs and outputs for comparison plots** option during the **Test** phase, the input and output variables of the top level-function, `kalman_filter` are automatically logged for plotting.

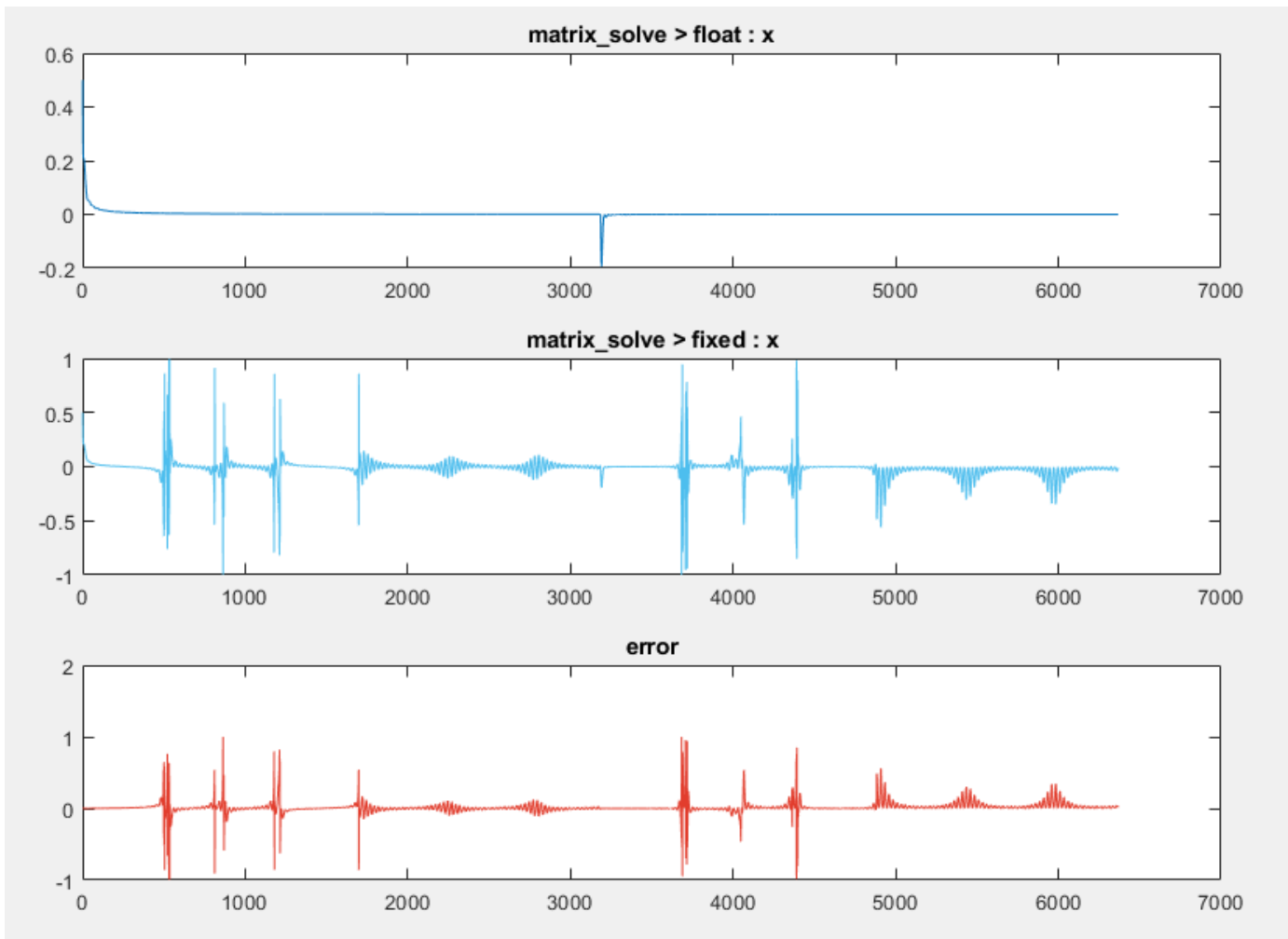
- 2 The `kalman_filter` function calls the `matrix_solve` function, which contains calls to several other functions. To investigate whether numerical issues are originating in the `matrix_solve` function, select `kalman_filter > matrix_solve` in the **Source Code** pane.

In the **Log Data** column, select the function input and output variables that you want to log. In this example, select all three, `a`, `b`, and `x`.

| Variables | Function Replacements | Output | Errors | Verification Output | | | | | |
|------------|-----------------------|---------|---------|---------------------|------------------------|----------|----------|--|--|
| Variable | Type | Sim Min | Sim Max | Whole ... | Proposed Type | Log Data | Max Diff | | |
| [-] Input | | | | | | | | | |
| a | double | 1 | 2 | No | numerictype(0, 16, 14) | ✓ | | | |
| b | 1 x 2 double | -0.25 | 1 | No | numerictype(1, 16, 14) | ✓ | | | |
| [-] Output | | | | | | | | | |
| x | 1 x 2 double | -0.19 | 0.5 | No | numerictype(1, 16, 15) | ✓ | | | |
| [-] Local | | | | | | | | | |
| l | double | 1 | 1 | Yes | numerictype(0, 1, 0) | | | | |

3 Click **Test**.

The generated plot shows a large error for the output variable of the `matrix_solve` function.

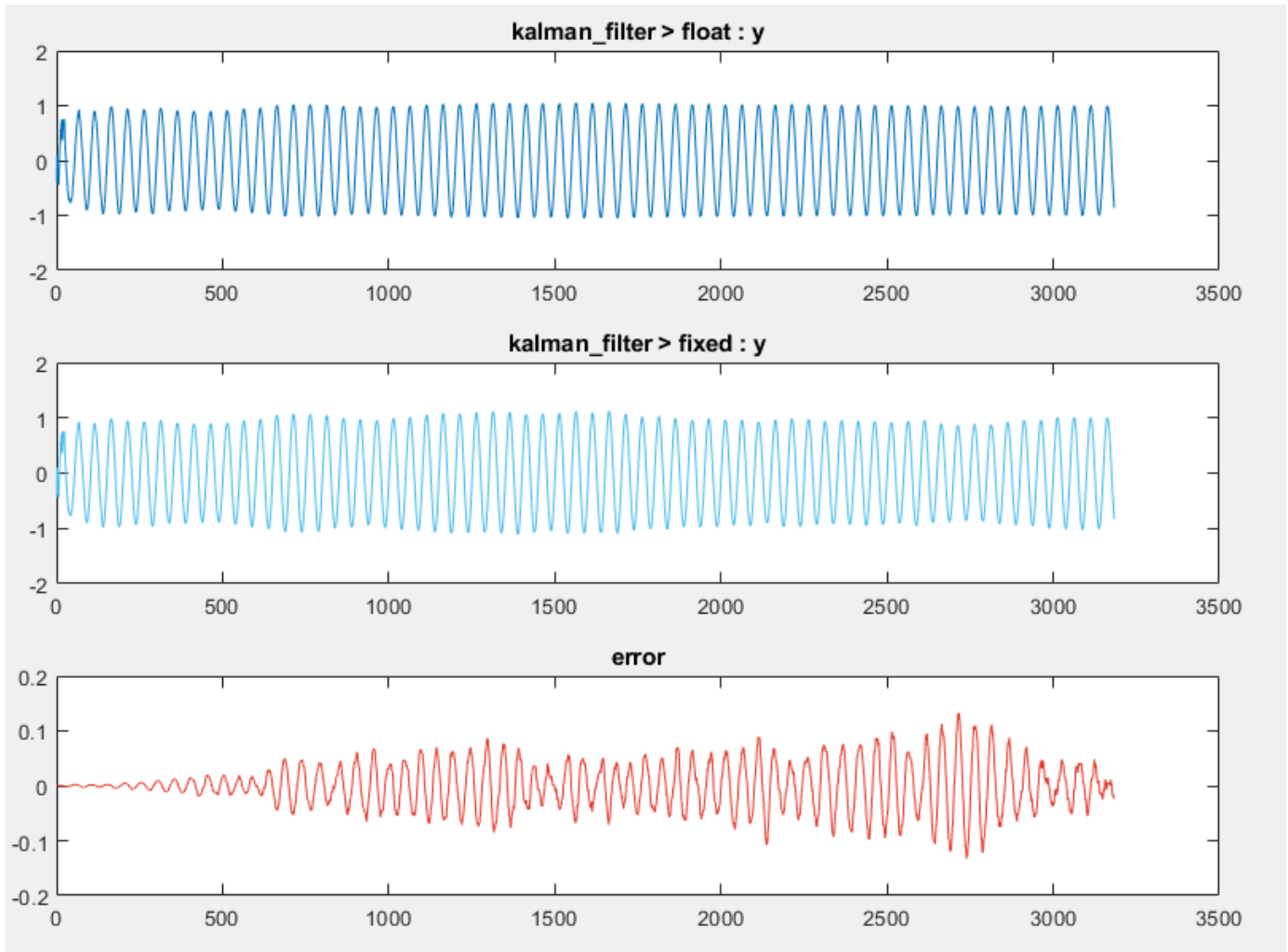


Adjust fimath Settings

1 On the **Convert to Fixed Point** page, click **Settings**.

- Under **fimath**, set the **Rounding method** to Nearest. Set the **Overflow action** to Saturate.
- 2** Click **Convert** to apply the new settings.
- 3** Click the arrow next to **Test** and ensure that **Log inputs and outputs for comparison plots** is selected. Enable logging for any function input or output variables. Click **Test**.

Examine the plot for top-level function output variable, *y*.

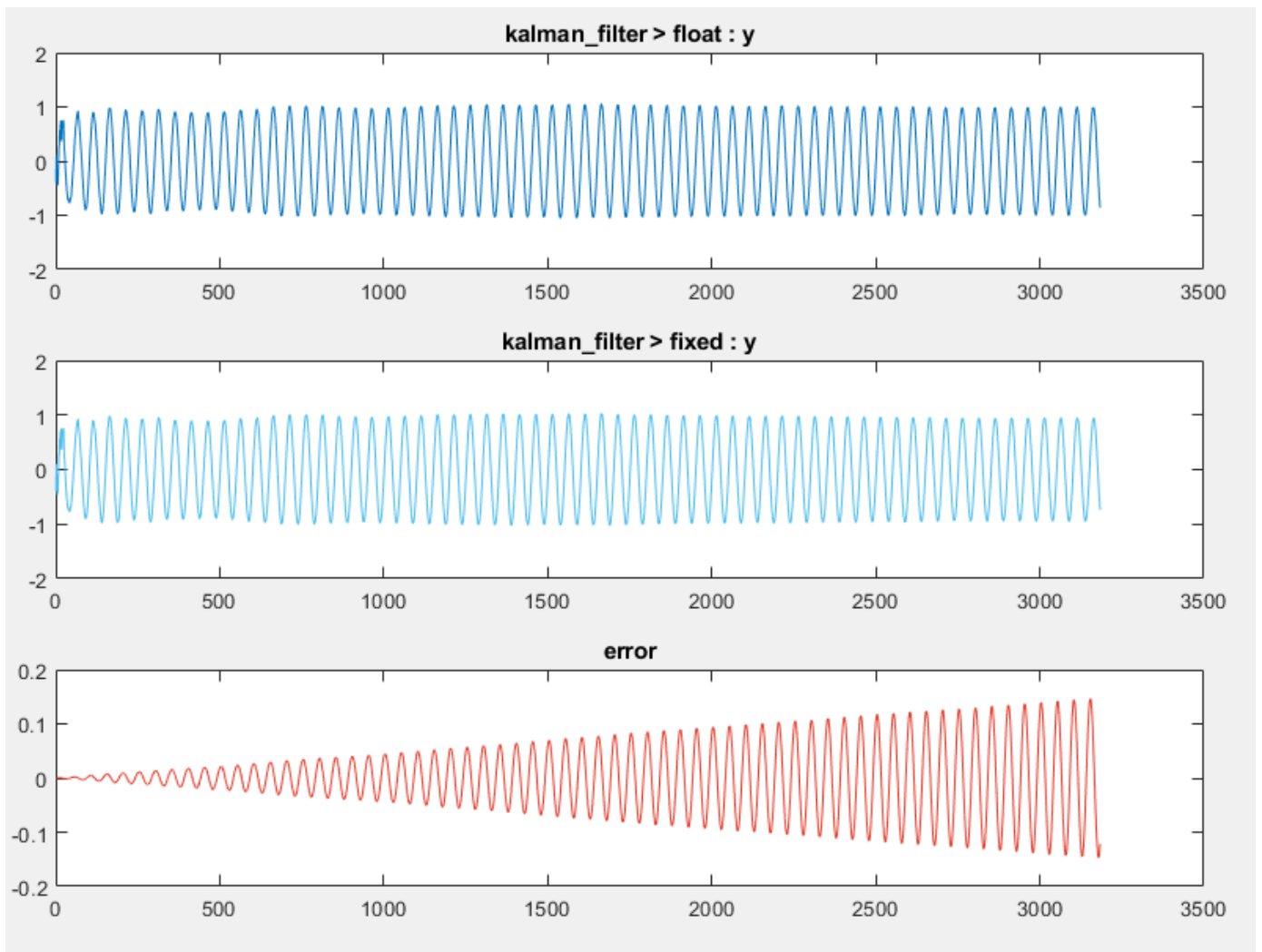


The new `fimath` settings improve the output, but some error still remains.

Adjust Word Length Settings

Adjusting the default word length improves the accuracy of the generated fixed-point design.

- 1** Click **Settings** and change the default word length to 32. Click **Convert** to apply the new settings.
- 2** Click **Test**. The error for the output variable *y* is accumulating.



- 3 Close the Fixed-Point Converter and plot window.

Replace Constant Functions

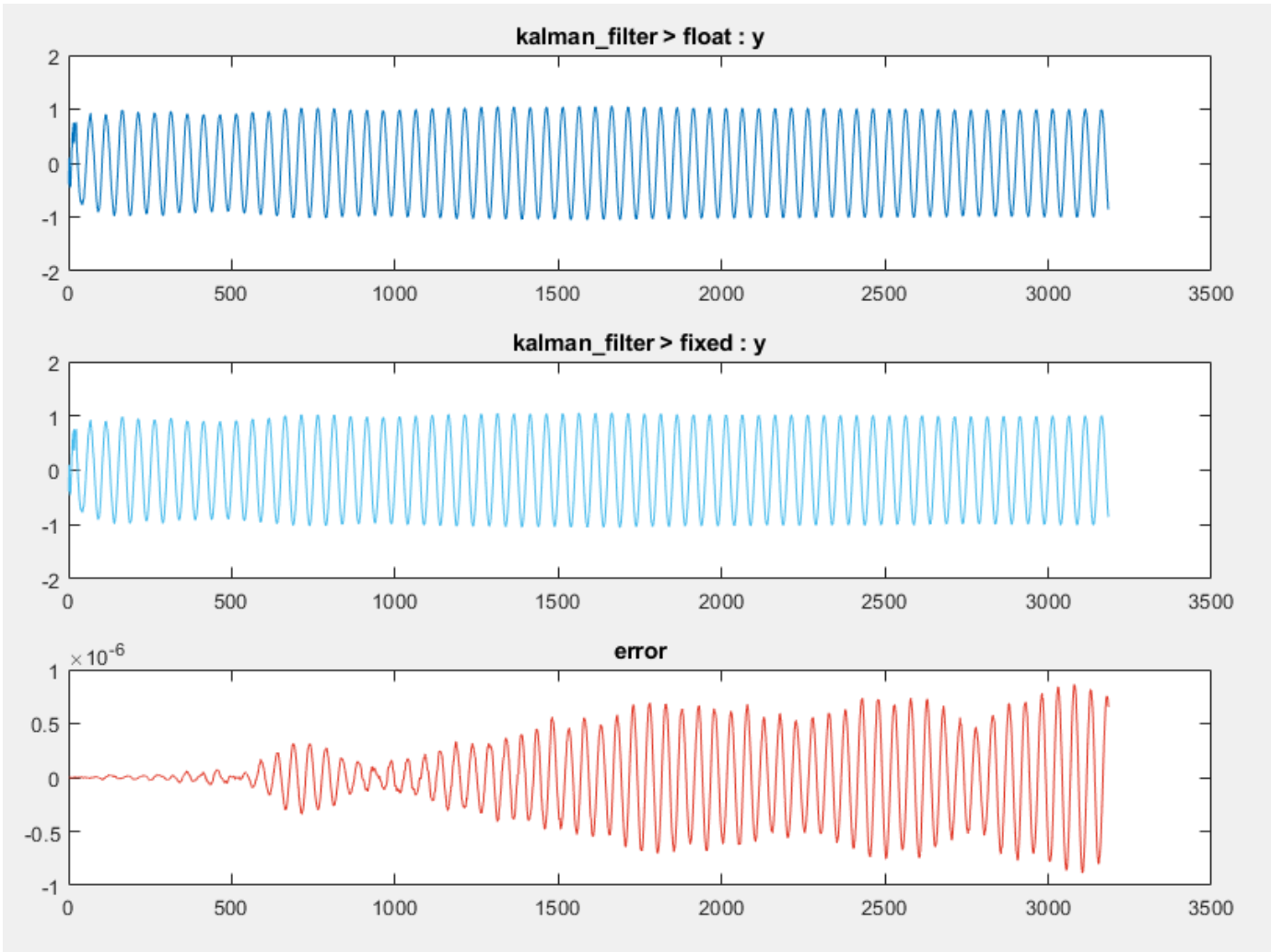
The `kalman_stm` function computes the state transition matrix, which is a constant. You do not need to convert this function to fixed point. To avoid unnecessary quantization through computation, replace this function with a double-precision constant. By replacing the function with a constant, the state transition matrix undergoes quantization only once.

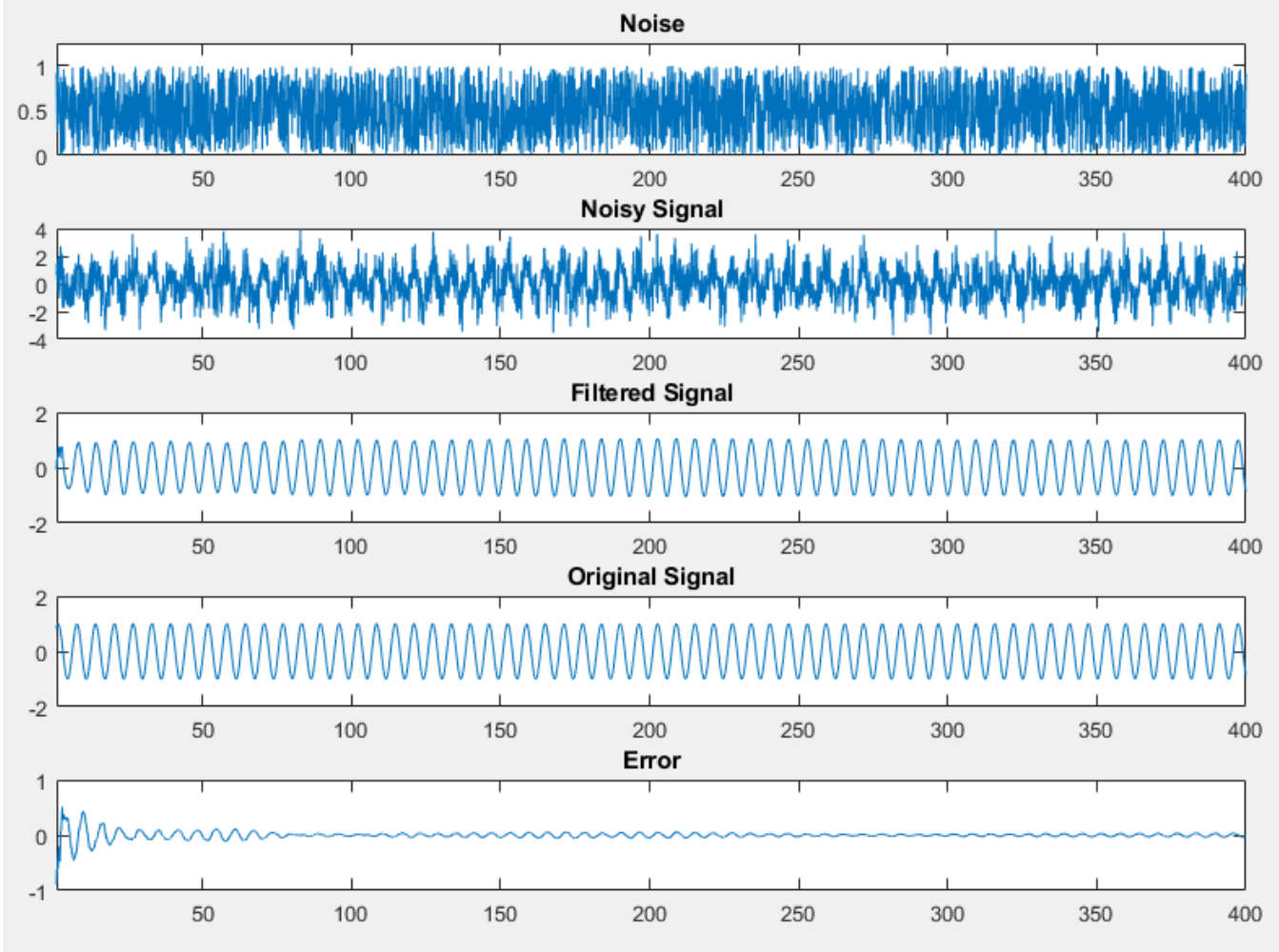
- 1 Click the `kalman_filter` function in the **Source Code** pane. Edit the `kalman_filter` function. Replace the call to the `kalman_stm` function with the equivalent double constant.

```
A = [0.992114701314478, -0.125333233564304;
     0.125333233564304, 0.992114701314478];
```

Save the changes.

- 2 Click **Analyze** to refresh the proposals.
- 3 Click **Convert** to apply the new proposals.
- 4 Click **Test**. The error on the plot for the functions output `y` is now on the order of 10^{-6} , which is acceptable for this design.





MATLAB Language Features Supported for Automated Fixed-Point Conversion

In this section...

“MATLAB Language Features Supported for Automated Fixed-Point Conversion” on page 7-35

“MATLAB Language Features Not Supported for Automated Fixed-Point Conversion” on page 7-36

MATLAB Language Features Supported for Automated Fixed-Point Conversion

Fixed-Point Designer supports the following MATLAB language features in automated fixed-point conversion:

- N-dimensional arrays
- Matrix operations, including deletion of rows and columns
- Variable-sized data (see “Generate Code for Variable-Size Data” (MATLAB Coder)). Range computation for variable-sized data is supported via simulation mode only. Variable-sized data is not supported for comparison plotting.
- Subscripting (see “Incompatibility with MATLAB in Matrix Indexing Operations for Code Generation” on page 29-22)
- Complex numbers (see “Code Generation for Complex Data” on page 16-8)
- Numeric classes (see “Supported Variable Types” on page 18-13)
- Double-precision, single-precision, and integer math
- Fixed-point arithmetic (see “Code Acceleration and Code Generation from MATLAB” on page 12-2)
- Program control statements `if`, `switch`, `for`, `while`, and `break`
- Arithmetic, relational, and logical operators
- Local functions
- Global variables
- Persistent variables
- Structures, including arrays of structures. Range computation for structures is supported via simulation mode only.
- Characters

The complete set of Unicode® characters is not supported for code generation. Characters are restricted to 8 bits of precision in generated code. Because many mathematical operations require more than 8 bits of precision, it is recommended that you do not perform arithmetic with characters if you intend to convert your MATLAB algorithm to fixed point.

- MATLAB classes. Range computation for MATLAB classes is supported via simulation mode only.

Automated conversion supports:

- Class properties

- Constructors
- Methods
- Specializations

It does not support class inheritance or packages. For more information, see “Fixed-Point Code for MATLAB Classes” on page 7-42.

- Ability to call functions (see “Resolution of Function Calls for Code Generation” on page 14-2)
- Subset of MATLAB toolbox functions (see “Functions Supported for Code Acceleration or C Code Generation” on page 12-4).
- Subset of DSP System Toolbox™ System objects.

The DSP System Toolbox System objects supported for automated conversion are:

- `dsp.BiquadFilter`
- `dsp.FIRDecimator`
- `dsp.FIRInterpolator`
- `dsp.FIRFilter` (Direct Form and Direct Form Transposed only)
- `dsp.FIRRateConverter`
- `dsp.VariableFractionalDelay`

MATLAB Language Features Not Supported for Automated Fixed-Point Conversion

Fixed-Point Designer does not support the following features in automated fixed-point conversion:

- Anonymous functions
- Cell arrays
- String scalars
- Objects of value classes as entry-point function inputs or outputs
- Function handles
- Java®
- Nested functions
- Recursion
- Sparse matrices
- `try/catch` statements
- `varargin`, `varargout`, or generation of fewer input or output arguments than an entry-point function defines
- Dot indexing properties of fixed-point data types.

Avoid using properties of fixed-point types in the code being converted by the Fixed-Point Converter app, and in MATLAB Function blocks being converted by the Fixed-Point Tool.

- Dynamic field/property references

Generated Fixed-Point Code

| In this section... |
|---|
| “Location of Generated Fixed-Point Files” on page 7-37 |
| “Minimizing fi-casts to Improve Code Readability” on page 7-37 |
| “Avoiding Overflows in the Generated Fixed-Point Code” on page 7-38 |
| “Controlling Bit Growth” on page 7-38 |
| “Avoiding Loss of Range or Precision” on page 7-39 |
| “Handling Non-Constant mpower Exponents” on page 7-40 |

Location of Generated Fixed-Point Files

By default, the fixed-point conversion process generates files in a folder named `codegen/fcn_name/fixpt` in your local working folder. `fcn_name` is the name of the MATLAB function that you are converting to fixed point.

| File name | Description |
|---|--|
| <code>fcn_name_fixpt.m</code> | Generated fixed-point MATLAB code. To integrate this fixed-point code into a larger application, consider generating a MEX-function for the function and calling this MEX-function in place of the original MATLAB code. |
| <code>fcn_name_fixpt_exVal.mat</code> | MAT-file containing: <ul style="list-style-type: none"> • A structure for the input arguments. • The name of the fixed-point file. |
| <code>fcn_name_fixpt_report.html</code> | Link to the type proposal report that displays the generated fixed-point code and the proposed type information. |
| <code>fcn_name_report.html</code> | Link to the type proposal report that displays the original MATLAB code and the proposed type information. |
| <code>fcn_name_wrapper_fixpt.m</code> | File that converts the floating-point data values supplied by the test file to the fixed-point types determined for the inputs during the conversion step. These fixed-point values are fed into the converted fixed-point function, <code>fcn_name_fixpt</code> . |

Minimizing fi-casts to Improve Code Readability

The conversion process tries to reduce the number of `fi`-casts by analyzing the floating-point code. If an arithmetic operation is comprised of only compile-time constants, the conversion process does not cast the operands to fixed point individually. Instead, it casts the entire expression to fixed point.

For example, here is the fixed-point code generated for the constant expression $x = 1/\sqrt{2}$ when the selected word length is 14.

| Original MATLAB Code | Generated Fixed-Point Code |
|-----------------------------|--|
| <code>x = 1/sqrt(2);</code> | <code>x = fi(1/sqrt(2), 0, 14, 14, fm);</code> <code>fm is the local fimath.</code> |

Avoiding Overflows in the Generated Fixed-Point Code

The conversion process avoids overflows by:

- Using full-precision arithmetic unless you specify otherwise.
- Avoiding arithmetic operations that involve double and `fi` data types. Otherwise, if the word length of the `fi` data type is not able to represent the value in the double constant expression, overflows occur.
- Avoiding overflows when adding and subtracting non fixed-point variables and fixed-point variables.

The fixed-point conversion process casts non-`fi` expressions to the corresponding `fi` type.

For example, consider the following MATLAB algorithm.

```
% A = 5;
% B = ones(300, 1)
function y = fi_plus_non_fi(A, B)
    % '1024' is non-fi, cast it
    y = A + 1024;
    % 'size(B, 1)*length(A)' is a non-fi, cast it
    y = A + size(B, 1)*length(A);
end
```

The generated fixed-point code is:

```
##codegen
% A = 5;
% B = ones(300,1)
function y = fi_plus_non_fi_fixpt(A,B)
    % '1024' is non-fi, cast it
    fm = fimath('RoundingMethod','Floor',...
        'OverflowAction','Wrap',...
        'ProductMode','FullPrecision',...
        'MaxProductWordLength',128,...
        'SumMode','FullPrecision',...
        'MaxSumWordLength',128);

    y = fi(A + fi(1024,0,11,0, fm),0,11,0, fm);
    % 'size(B, 1)*length(A)' is a non-fi, cast it
    y(:) = A + fi(size(B,fi(1,0,1,0, fm))*length(A),0,9,0, fm);
end
```

Controlling Bit Growth

The conversion process controls bit growth by using subscripted assignments, that is, assignments that use the colon (`:`) operator, in the generated code. When you use subscripted assignments,

MATLAB overwrites the value of the left-hand side argument but retains the existing data type and array size. Using subscripted assignment keeps fixed-point variables fixed point rather than inadvertently turning them into doubles. Maintaining the fixed-point type reduces the number of type declarations in the generated code. Subscripted assignment also prevents bit growth which is useful when you want to maintain a particular data type for the output.

Avoiding Loss of Range or Precision

Avoiding Loss of Range or Precision in Unsigned Subtraction Operations

When the result of the subtraction is negative, the conversion process promotes the left operand to a signed type.

For example, consider the following MATLAB algorithm.

```
% A = 1;
% B = 5
function [y,z] = unsigned_subtraction(A,B)
    y = A - B;

    C = -20;
    z = C - B;
end
```

In the original code, both A and B are unsigned and the result of A-B can be negative. In the generated fixed-point code, A is promoted to signed. In the original code, C is signed, so does not require promotion in the generated code.

```
##codegen
% A = 1;
% B = 5
function [y,z] = unsigned_subtraction_fixpt(A,B)

fm = fimath('RoundingMethod','Floor',...
    'OverflowAction','Wrap',...
    'ProductMode','FullPrecision',...
    'MaxProductWordLength',128,...
    'SumMode','FullPrecision',...
    'MaxSumWordLength',128);
y = fi(fi_signed(A) - B,1,3,0,fm);
C = fi(-20,1,6,0,fm);
z = fi(C - B,1,6,0,fm);
end

function y = fi_signed(a)
coder.inline('always');
if isfi(a) && ~(issigned(a))
    nt = numerictype(a);
    new_nt = numerictype(1,nt.WordLength + 1,nt.FractionLength);
    y = fi(a,new_nt,fimath(a));
else
    y = a;
end
end
```

Avoiding Loss of Range When Concatenating Arrays of Fixed-Point Numbers

If you concatenate matrices using `vertcat` and `horzcat`, the conversion process uses the largest `numericType` among the expressions of a row and casts the leftmost element to that type. This type is then used for the concatenated matrix to avoid loss of range.

For example, consider the following MATLAB algorithm.

```
% A = 1, B = 100, C = 1000
function [y, z] = lb_node(A, B, C)
    %% single rows
    y = [A B C];
    %% multiple rows
    z = [A 5; A B; A C];
end
```

In the generated fixed-point code:

- For the expression `y = [A B C]`, the leftmost element, `A`, is cast to the type of `C` because `C` has the largest type in the row.
- For the expression `[A 5; A B; A C]`:
 - In the first row, `A` is cast to the type of `C` because `C` has the largest type of the whole expression.
 - In the second row, `A` is cast to the type of `B` because `B` has the larger type in the row.
 - In the third row, `A` is cast to the type of `C` because `C` has the larger type in the row.

```
##codegen
% A = 1, B = 100, C = 1000
function [y,z] = lb_node_fixpt(A,B,C)
    %% single rows
    fm = fimath('RoundingMethod','Floor',...
        'OverflowAction','Wrap',...
        'ProductMode','FullPrecision',...
        'MaxProductWordLength',128, ...
        'SumMode','FullPrecision',...
        'MaxSumWordLength',128);

    y = fi([fi(A,0,10,0, fm) B C],0,10,0, fm);

    %% multiple rows
    z = fi([fi(A,0,10,0, fm) 5; fi(A,0,7,0, fm) B; ...
        fi(A,0,10,0, fm) C],0,10,0, fm);
end
```

Handling Non-Constant mpower Exponents

If the function that you are converting has a scalar input, and the `mpower` exponent input is not constant, the conversion process sets the `fimath ProductMode` to `SpecifyPrecision` in the generated code. With this setting, the output data type can be determined at compile time.

For example, consider the following MATLAB algorithm.

```
% a = 1
% b = 3
```



```

function y = exp_operator(a, b)
    % exponent is a constant so no need to specify precision
    y = a^3;
    % exponent is not a constant, use 'SpecifyPrecision' for 'ProductMode'
    y = b^a;
end

```

In the generated fixed-point code, for the expression $y = a^3$, the exponent is a constant, so there is no need to specify precision. For the expression, $y = b^a$, the exponent is not constant, so the ProductMode is set to SpecifyPrecision.

```

%#codegen
% a = 1
% b = 3
function y = exp_operator_fixpt(a,b)
    % exponent is a constant so no need to specify precision
    fm = fimath('RoundingMethod','Floor',...
        'OverflowAction','Wrap',...
        'ProductMode','FullPrecision',...
        'MaxProductWordLength',128,...
        'SumMode','FullPrecision',...
        'MaxSumWordLength',128);

    y = fi(a^3,0,2,0,fm);
    % exponent is not a constant, use 'SpecifyPrecision' for 'ProductMode'
    y(:) = fi(b,'ProductMode','SpecifyPrecision',...
        'ProductWordLength',2,...
        'ProductFractionLength',0)^a;
end

```

Fixed-Point Code for MATLAB Classes

In this section...

“Automated Conversion Support for MATLAB Classes” on page 7-42

“Unsupported Constructs” on page 7-42

“Coding Style Best Practices” on page 7-42

Automated Conversion Support for MATLAB Classes

The automated fixed-point conversion process:

- Proposes fixed-point data types based on simulation ranges for MATLAB classes. It does not propose data types based on derived ranges for MATLAB classes.

After simulation, the Fixed-Point Converter app:

- Function list contains class constructors, methods, and specializations.
- Code window displays the objects used in each function.
- Provides code coverage for methods.

For more information, see “Viewing Information for MATLAB Classes” on page 7-18.

- Supports class methods, properties, and specializations. For each specialization of a class, `class_name`, the conversion generates a separate `class_name_fixpt.m` file. For every instantiation of a class, the generated fixed-point code contains a call to the constructor of the appropriate specialization.
- Supports classes that have `get` and `set` methods such as `get.PropertyName`, `set.PropertyName`. These methods are called when properties are read or assigned. The `set` methods can be specialized. Sometimes, in the generated fixed-point code, assignment statements are transformed to function calls.

Unsupported Constructs

The automated conversion process does not support:

- Class inheritance.
- Packages.
- Constructors that use `nargin` and `varargin`.

Coding Style Best Practices

When you write MATLAB code that uses MATLAB classes:

- Initialize properties in the class constructor.
- Replace constant properties with static methods.

For example, consider the `counter` class.

```
classdef Counter < handle
    properties
```

```

    Value = 0;
end

properties(Constant)
    MAX_VALUE = 128
end

methods
    function out = next(this)
        out = this.Count;
        if this.Value == this.MAX_VALUE
            this.Value = 0;
        else
            this.Value = this.Value + 1;
        end
    end
end
end
end

```

To use the automated fixed-point conversion process, rewrite the class to have a static class that initializes the constant property `MAX_VALUE` and a constructor that initializes the property `Value`.

```

classdef Counter < handle
    properties
        Value;
    end

    methods(Static)
        function t = MAX_VALUE()
            t = 128;
        end
    end

    methods
        function this = Counter()
            this.Value = 0;
        end
        function out = next(this)
            out = this.Value;
            if this.Value == this.MAX_VALUE
                this.Value = 0;
            else
                this.Value = this.Value + 1;
            end
        end
    end
end
end

```

Automated Fixed-Point Conversion Best Practices

In this section...

“Create a Test File” on page 7-44
 “Prepare Your Algorithm for Code Acceleration or Code Generation” on page 7-45
 “Check for Fixed-Point Support for Functions Used in Your Algorithm” on page 7-45
 “Manage Data Types and Control Bit Growth” on page 7-46
 “Convert to Fixed Point” on page 7-46
 “Use the Histogram to Fine-Tune Data Type Settings” on page 7-47
 “Optimize Your Algorithm” on page 7-47
 “Avoid Explicit Double and Single Casts” on page 7-49

Create a Test File

A best practice for structuring your code is to separate your core algorithm from other code that you use to test and verify the results. Create a test file to call your original MATLAB algorithm and fixed-point versions of the algorithm. For example, as shown in the following table, you might set up some input data to feed into your algorithm, and then, after you process that data, create some plots to verify the results. Since you need to convert only the algorithmic portion to fixed point, it is more efficient to structure your code so that you have a test file, in which you create your inputs, call your algorithm, and plot the results, and one (or more) algorithmic files, in which you do the core processing.

| Original code | Best Practice | Modified code |
|---|--|--|
| <pre>% TEST INPUT x = randn(100,1); % ALGORITHM y = zeros(size(x)); y(1) = x(1); for n=2:length(x) y(n)=y(n-1) + x(n); end % VERIFY RESULTS yExpected=cumsum(x); plot(y-yExpected) title('Error')</pre> | <p>Issue</p> <p>Generation of test input and verification of results are intermingled with the algorithm code.</p> <p>Fix</p> <p>Create a test file that is separate from your algorithm. Put the algorithm in its own function.</p> | <p>Test file</p> <pre>% TEST INPUT x = randn(100,1); % ALGORITHM y = cumulative_sum(x); % VERIFY RESULTS yExpected = cumsum(x); plot(y-yExpected) title('Error')</pre> <p>Algorithm in its own function</p> <pre>function y = cumulative_sum(x) y = zeros(size(x)); y(1) = x(1); for n=2:length(x) y(n) = y(n-1) + x(n); end end</pre> |

You can use the test file to:

- Verify that your floating-point algorithm behaves as you expect before you convert it to fixed point. The floating-point algorithm behavior is the baseline against which you compare the behavior of the fixed-point versions of your algorithm.
- Propose fixed-point data types.
- Compare the behavior of the fixed-point versions of your algorithm to the floating-point baseline.
- Help you determine initial values for static ranges.

By default, the Fixed-Point Converter app shows code coverage results. Your test files should exercise the algorithm over its full operating range so that the simulation ranges are accurate. For example, for a filter, realistic inputs are impulses, sums of sinusoids, and chirp signals. With these inputs, using linear theory, you can verify that the outputs are correct. Signals that produce maximum output are useful for verifying that your system does not overflow. The quality of the proposed fixed-point data types depends on how well the test files cover the operating range of the algorithm with the accuracy that you want. Reviewing code coverage results help you verify that your test file is exercising the algorithm adequately. Review code flagged with a red code coverage bar because this code is not executed. If the code coverage is inadequate, modify the test file or add more test files to increase coverage. See “Code Coverage” on page 7-5.

Prepare Your Algorithm for Code Acceleration or Code Generation

The automated conversion process instruments your code and provides data type proposals to help you convert your algorithm to fixed point.

MATLAB algorithms that you want to convert to fixed point automatically must comply with code generation requirements and rules. To view the subset of the MATLAB language that is supported for code generation, see “Functions and Objects Supported for C/C++ Code Generation” on page 26-2.

To help you identify unsupported functions or constructs in your MATLAB code, add the `%#codegen` pragma to the top of your MATLAB file. The MATLAB Code Analyzer flags functions and constructs that are not available in the subset of the MATLAB language supported for code generation. This advice appears in real time as you edit your code in the MATLAB editor. For more information, see “Check Code Using the MATLAB Code Analyzer” on page 12-65. The software provides a link to a report that identifies calls to functions and the use of data types that are not supported for code generation. For more information, see “Check Code Using the Code Generation Readiness Tool” on page 12-64.

Check for Fixed-Point Support for Functions Used in Your Algorithm

The app flags unsupported function calls found in your algorithm on the **Function Replacements** tab. For example, if you use the `fft` function, which is not supported for fixed point, the tool adds an entry to the table on this tab and indicates that you need to specify a replacement function to use for fixed-point operations.

| Variables | | Function Replacements |
|---|---|-----------------------|
| Enter a function to replace | | |
| Function or Operator | Replacement | |
| <ul style="list-style-type: none"> Custom Function | Function Name | |
| fft | Replacement required to use fixed-point | |

You can specify additional replacement functions. For example, functions like `sin`, `cos`, and `sqrt` might support fixed point, but for better efficiency, you might want to consider an alternative implementation like a lookup table or CORDIC-based algorithm. The app provides an option to generate lookup table approximations for continuous and stateless single-input, single-output functions in your original MATLAB code. See “Replacing Functions Using Lookup Table Approximations” on page 7-50.

Manage Data Types and Control Bit Growth

The automated fixed-point conversion process automatically manages data types and controls bit growth. It controls bit growth by using subscripted assignments, that is, assignments that use the colon (`:`) operator, in the generated code. When you use subscripted assignments, MATLAB overwrites the value of the left-hand side argument but retains the existing data type and array size. In addition to preventing bit growth, subscripted assignment reduces the number of casts in the generated fixed-point code and makes the code more readable.

Convert to Fixed Point

What Are Your Goals for Converting to Fixed Point?

Before you start the conversion, consider your goals for converting to fixed point. Are you implementing your algorithm in C or HDL? What are your target constraints? The answers to these questions determine many fixed-point properties such as the available word length, fraction length, and math modes, as well as available math libraries.

To set up these properties, use the **Advanced** settings.


| Setting | Value |
|-----------------------------------|------------------------|
| Default word length | 16 |
| Default fraction length | 4 |
| Advanced | |
| When proposing types | use all collected data |
| Propose target container types | No |
| Optimize whole numbers | Yes |
| Signedness | Automatic |
| Safety margin for sim min/max (%) | 0 |
| Search paths | |
| fimath | |
| Rounding method | Floor |
| Overflow action | Wrap |

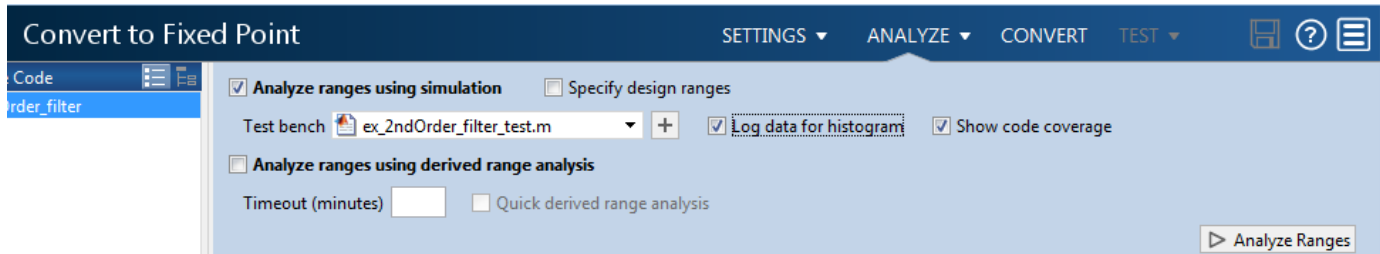
For more information, see “Specify Type Proposal Options” on page 8-2.

Run With Fixed-Point Types and Compare Results

Create a test file to validate that the floating-point algorithm works as expected before converting it to fixed point. You can use the same test file to propose fixed-point data types, and to compare fixed-point results to the floating-point baseline after the conversion. For more information, see “Running a Simulation” on page 7-8 and “Log Data for Histogram” on page 7-19 .

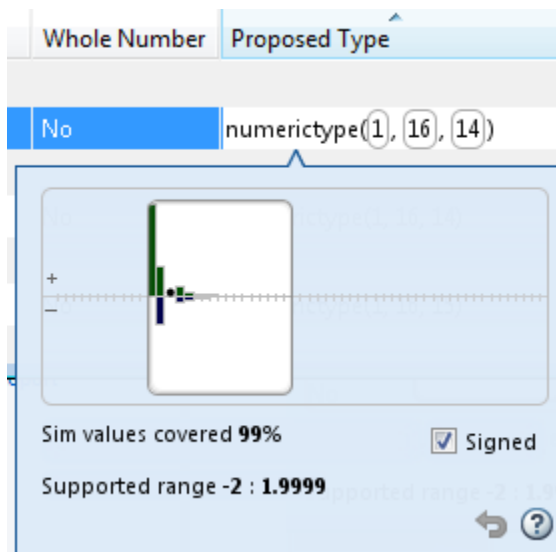
Use the Histogram to Fine-Tune Data Type Settings

To fine-tune fixed-point type settings, use the histogram. To log data for histograms, in the app, click the **Analyze** arrow  and select Log data for histogram.



After simulation and static analysis:

- To view the histogram for a variable, on the **Variables** tab, click the **Proposed Type** field for that variable.



You can view the effect of changing the proposed data types by dragging the edges of the bounding box in the histogram window to change the proposed data type and selecting or clearing the **Signed** option.

- If the values overflow and the range cannot fit the proposed type, the table shows proposed types in red.

When the tool applies data types, it generates an html report that provides overflow information and highlights overflows in red. Review the proposed data types.

Optimize Your Algorithm

Use `fimath` to Get Optimal Types for C or HDL

`fimath` properties define the rules for performing arithmetic operations on `fi` objects, including math, rounding, and overflow properties. You can use the `fimath` `ProductMode` and `SumMode`

properties to retain optimal data types for C or HDL. HDL can have arbitrary word length types in the generated HDL code whereas C requires container types (uint8, uint16, uint32). Use the **Advanced** settings, see “Specify Type Proposal Options” on page 8-2.

C

The **KeepLSB** setting for **ProductMode** and **SumMode** models the behavior of integer operations in the C language, while **KeepMSB** models the behavior of many DSP devices. Different rounding methods require different amounts of overhead code. Setting the **RoundingMethod** property to **Floor**, which is equivalent to two's complement truncation, provides the most efficient rounding implementation. Similarly, the standard method for handling overflows is to wrap using modulo arithmetic. Other overflow handling methods create costly logic. Whenever possible, set **OverflowAction** to **Wrap**.

| MATLAB Code | Best Practice | Generated C Code | | |
|---|---|--|-------|--|
| <p>Code being compiled</p> <pre>function y = adder(a,b) y = a + b; end</pre> <p>Note In the app, set Default word length to 16.</p> | <p>Issue</p> <p>With the default word length set to 16 and the default fimath settings, additional code is generated to implement saturation overflow, nearest rounding, and full-precision arithmetic.</p> | <pre>int adder(short a, short b) { int y; int i; int i1; int i2; int i3; i = a; i1 = b; if ((i & 65536) != 0) { i2 = i -65536; } else { i2 = i & 65535; } if ((i1 & 65536) != 0) { i3 = i1 -65536; } else { i3 = i1 & 65535; } i = i2 + i3; if ((i & 65536) != 0) { y = i -65536; } else { y = i & 65535; } return y; }</pre> | | |
| | <p>Fix</p> <p>To make the generated C code more efficient, choose fixed-point math settings that match your processor types.</p> <p>To customize fixed-point type proposals, use the app Settings. Select fimath and then set:</p> <table border="1" data-bbox="509 1835 1081 1873"> <tr> <td>Rounding method</td> <td>Floor</td> </tr> </table> | Rounding method | Floor | <pre>int adder(short a, short b) { return a + b; }</pre> |
| Rounding method | Floor | | | |

| MATLAB Code | Best Practice | | Generated C Code |
|-------------|---------------------|---------|------------------|
| | Overflow action | Wrap | |
| | Product mode | KeepLSB | |
| | Sum mode | KeepLSB | |
| | Product word length | 32 | |
| | Sum word length | 32 | |

HDL

For HDL code generation, set:

- ProductMode and SumMode to FullPrecision
- Overflow action to Wrap
- Rounding method to Floor

Replace Built-in Functions with More Efficient Fixed-Point Implementations

Some MATLAB built-in functions can be made more efficient for fixed-point implementation. For example, you can replace a built-in function with a Lookup table implementation, or a CORDIC implementation, which requires only iterative shift-add operations. For more information, see “Function Replacements” on page 7-21.

Reimplement Division Operations Where Possible

Often, division is not fully supported by hardware and can result in slow processing. When your algorithm requires a division, consider replacing it with one of the following options:

- Use bit shifting when the denominator is a power of two. For example, `bitsra(x,3)` instead of `x/8`.
- Multiply by the inverse when the denominator is constant. For example, `x*0.2` instead of `x/5`.
- If the divisor is not constant, use a temporary variable for the division. Doing so results in a more efficient data type proposal and, if overflows occur, makes it easier to see which expression is overflowing.

Eliminate Floating-Point Variables

For more efficient code, the automated fixed-point conversion process eliminates floating-point variables. The one exception to this is loop indices because they usually become integer types. It is good practice to inspect the fixed-point code after conversion to verify that there are no floating-point variables in the generated fixed-point code.

Avoid Explicit Double and Single Casts

For the automated workflow, do not use explicit double or single casts in your MATLAB algorithm to insulate functions that do not support fixed-point data types. The automated conversion tool does not support these casts.

Instead of using casts, supply a replacement function. For more information, see “Function Replacements” on page 7-21.

Replacing Functions Using Lookup Table Approximations

The Fixed-Point Designer software provides an option to generate lookup table approximations for continuous and stateless single-input, single-output functions in your original MATLAB code. These functions must be on the MATLAB path.

You can use this capability to handle functions that are not supported for fixed point and to replace your own custom functions. The fixed-point conversion process infers the ranges for the function and then uses an interpolated lookup table to replace the function. You can control the interpolation method and number of points in the lookup table. By adjusting these settings, you can tune the behavior of replacement function to match the behavior of the original function as closely as possible.

The fixed-point conversion process generates one lookup table approximation per call site of the function that needs replacement.

To use lookup table approximations in the Fixed-Point Converter app, see “Replace the exp Function with a Lookup Table” on page 8-38 and “Replace a Custom Function with a Lookup Table” on page 8-47.

To use lookup table approximations in the programmatic workflow, see `coder.approximation`, “Replace the exp Function with a Lookup Table” on page 9-16, and “Replace a Custom Function with a Lookup Table” on page 9-18.

Custom Plot Functions

The Fixed-Point Converter app provides a default time series based plotting function. The conversion process uses this function at the test numerics step to show the floating-point and fixed-point results and the difference between them. However, during fixed-point conversion you might want to visualize the numerical differences in a view that is more suitable for your application domain. For example, plots that show eye diagrams and bit error differences are more suitable in the communications domain and histogram difference plots are more suitable in image processing designs.

You can choose to use a custom plot function at the test numerics step. The Fixed-Point Converter app facilitates custom plotting by providing access to the raw logged input and output data before and after fixed-point conversion. You supply a custom plotting function to visualize the differences between the floating-point and fixed-point results. If you specify a custom plot function, the fixed-point conversion process calls the function for each input and output variable, passes in the name of the variable and the function that uses it, and the results of the floating-point and fixed-point simulations.

Your function should accept three inputs:

- A structure that holds the name of the variable and the function that uses it.

Use this information to:

- Customize plot headings and axes.
- Choose which variables to plot.
- Generate different error metrics for different output variables.
- A cell array to hold the logged floating-point values for the variable.

This cell array contains values observed during floating-point simulation of the algorithm during the test numerics phase. You might need to reformat this raw data.

- A cell array to hold the logged values for the variable after fixed-point conversion.

This cell array contains values observed during fixed-point simulation of the converted design.

For example, function `customComparisonPlot(varInfo, floatVarVals, fixedPtVarVals)`.

To use a custom plot function, in the Fixed-Point Converter app, select **Advanced**, and then set **Custom plot function** to the name of your plot function. See “Visualize Differences Between Floating-Point and Fixed-Point Results” on page 8-52.

In the programmatic workflow, set the `coder.FixPtConfig` configuration object `PlotFunction` property to the name of your plot function. See “Visualize Differences Between Floating-Point and Fixed-Point Results” on page 9-20.

Generate Fixed-Point MATLAB Code for Multiple Entry-Point Functions

When your end goal is to generate fixed-point C/C++ library functions, generating a single C/C++ library for more than one entry-point MATLAB function allows you to:

- Create C/C++ libraries containing multiple, compiled MATLAB files to integrate with larger C/C++ applications. Generating C/C++ code requires a MATLAB Coder license.
- Share code efficiently between library functions.
- Communicate between library functions using shared memory.

Note If any of the entry-point functions in a project share memory (for example, persistent data), an error will occur. In this case, you should rewrite your code to avoid invoking functions with persistent data from multiple entry-points.

Example 7.1. Convert Two Entry-Point Functions to Fixed-Point Using the Fixed-Point Converter App

In this example, you convert two entry-point functions, `ep1` and `ep2`, to fixed point.

- 1 In a local writable folder, create the functions `ep1.m` and `ep2.m`.

```
function y = ep1(u) %#codegen
y = u;
end

function y = ep2(u, v) %#codegen
y = u + v;
end
```

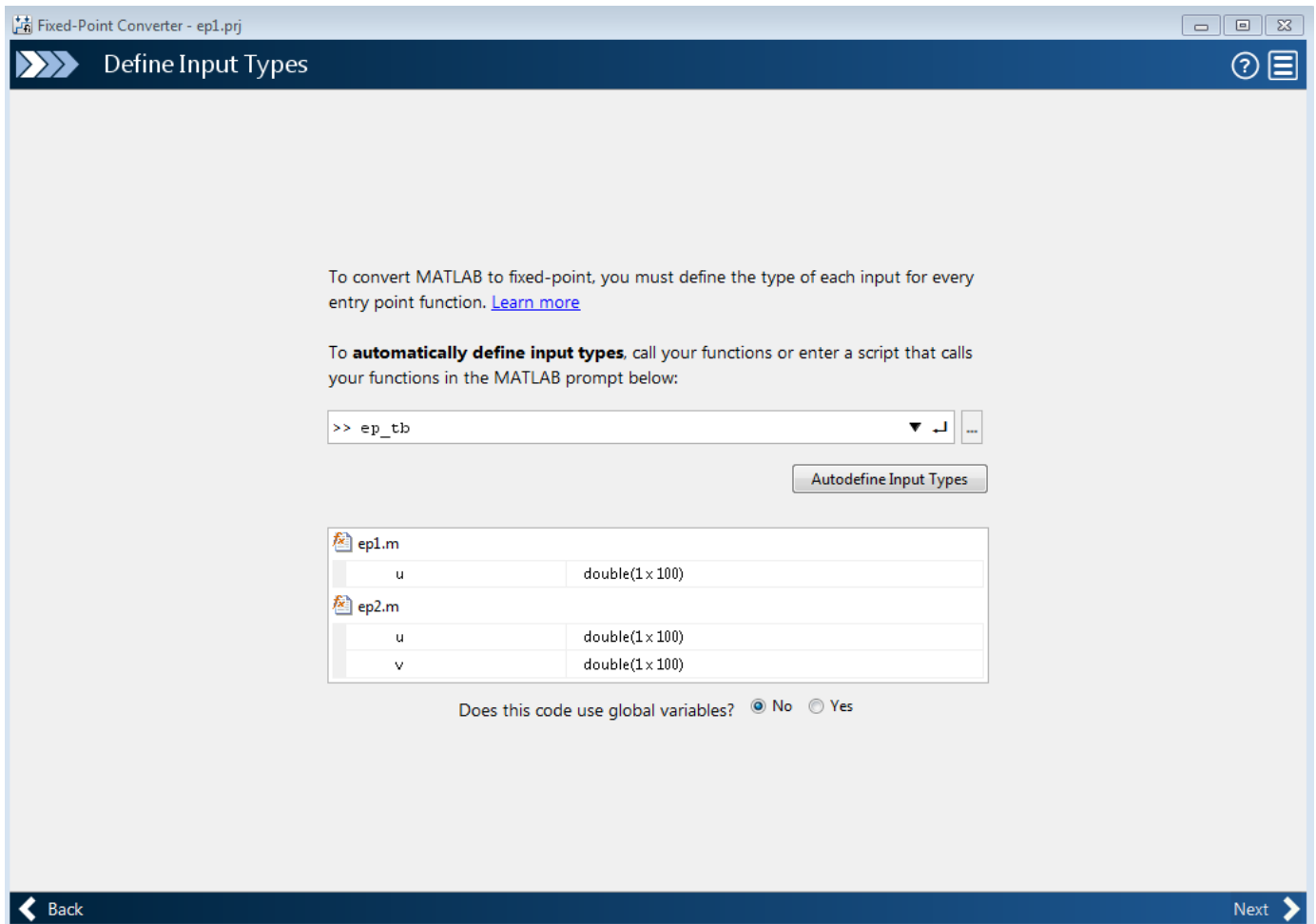
- 2 In the same folder, create a test file, `ep_tb.m`, that calls both functions.

```
% test file for ep1 and ep2
u = 1:100;
v = 5:104;
z = ep1(u);
y = ep2(v,z);
```

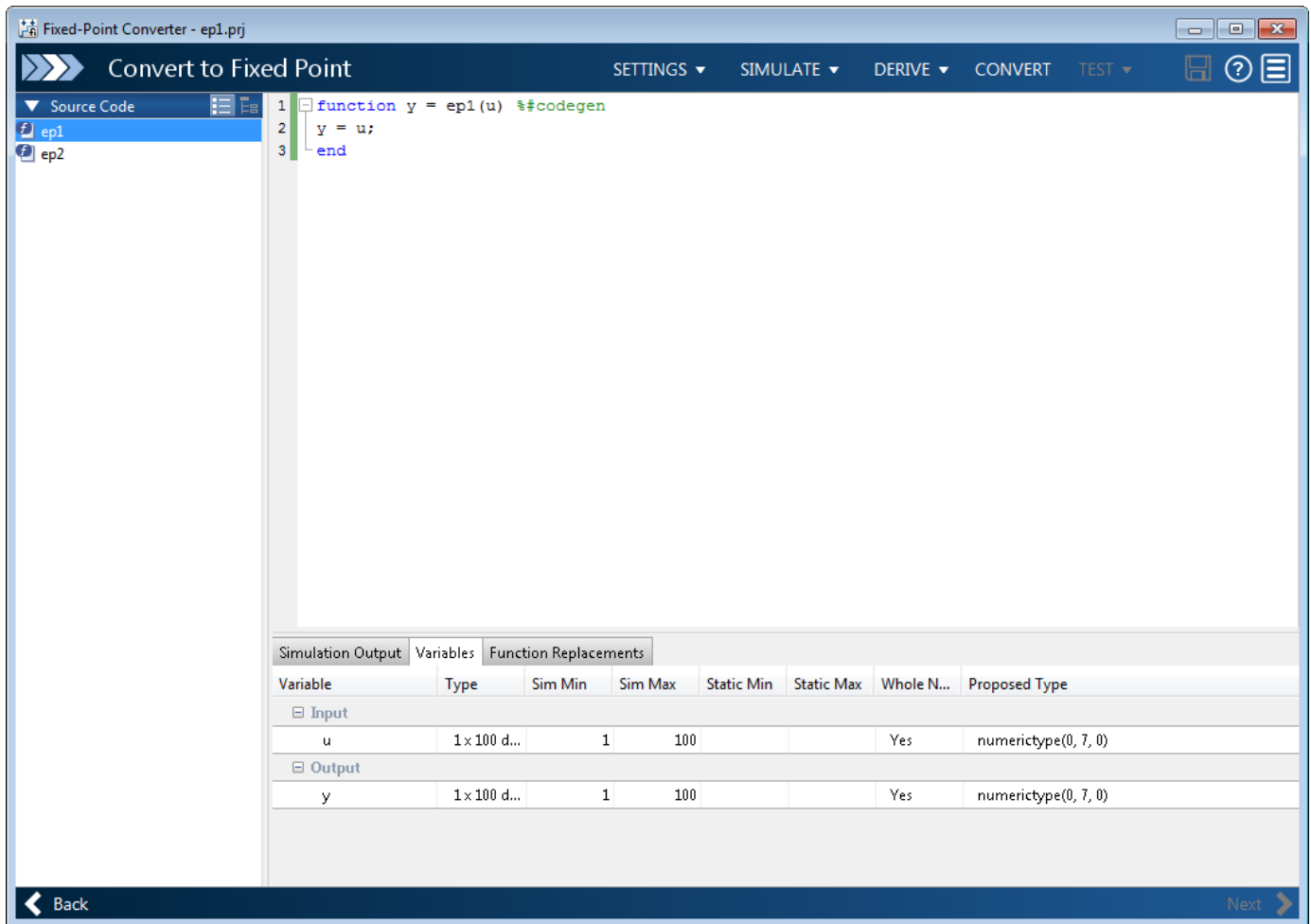
- 3 From the apps gallery, open the Fixed-Point Converter app.
- 4 To add the first entry-point function, `ep1`, to the project, on the **Select Source Files** page, browse to the `ep1` file, and click **Open**.

By default, the app uses the name of the first entry-point function as the name of the project.

- 5 Click **Add Entry-Point Function** and add the second entry-point function, `ep2`. Click **Next**.
- 6 On the **Define Input Types** page, enter a test file that exercises your two entry-point functions. Browse to select the `ep_tb` file. Click **Autodefine Input Types**.

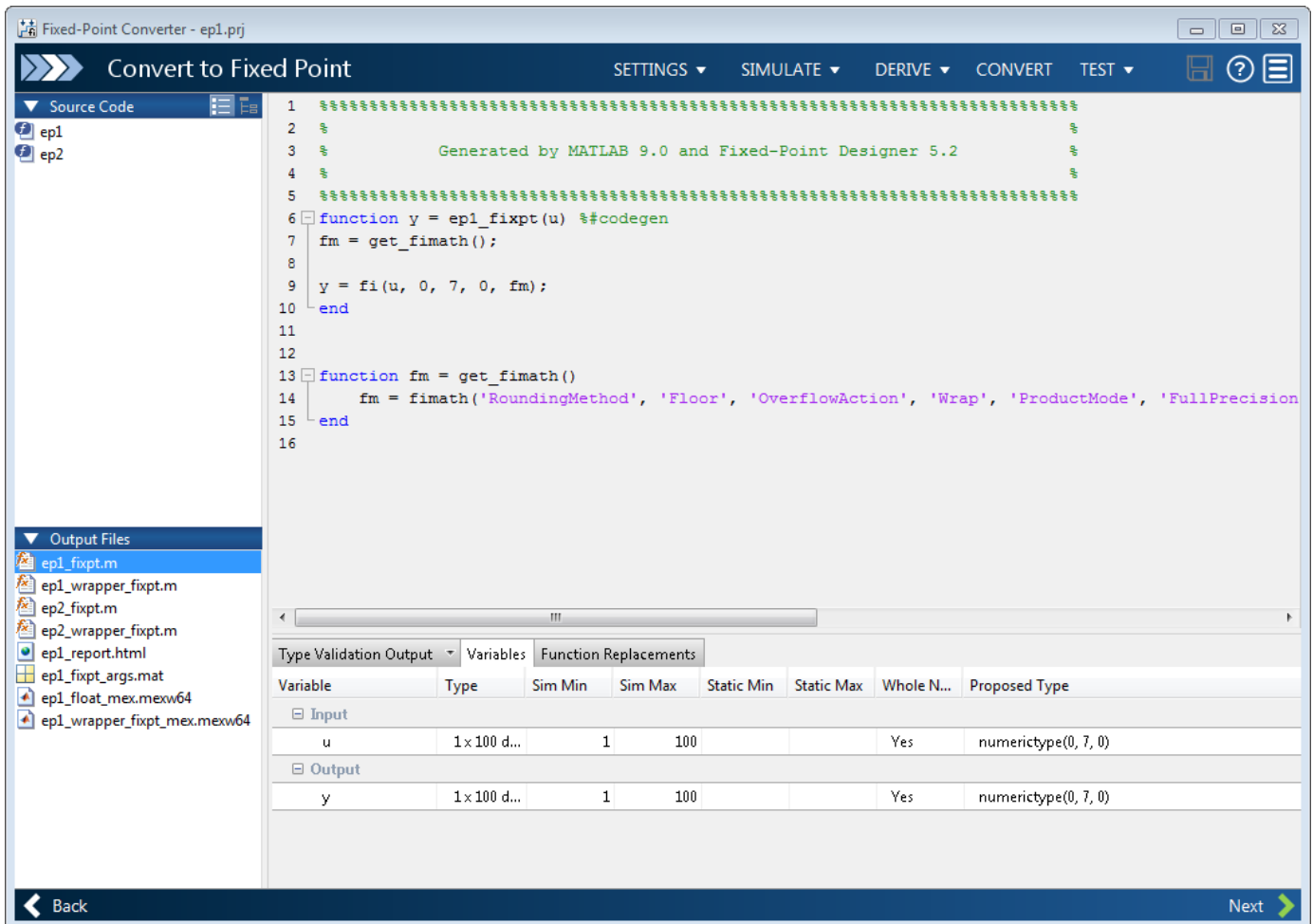


- 7 Click **Next**. The app generates an instrumented MEX function for your entry-point MATLAB function. On the **Convert to Fixed-Point** page, click **Simulate** to simulate the entry-point functions, gather range information, and get proposed data types.



8 Click **Convert**.

The entry-point functions `ep1` and `ep2` convert to fixed point. The **Output Files** pane lists the generated fixed-point and wrapper files for both entry-point functions.



- Click **Next**. The **Finish Workflow** page contains the project summary. The generated Fixed-Point Conversion Report contains the reports for both entry-point functions. The app stores the generated files in the subfolder `codegen/ep1/fixpt`.

Convert Code Containing Global Data to Fixed Point

In this section...

“Workflow” on page 7-56

“Declare Global Variables” on page 7-56

“Define Global Data” on page 7-56

“Define Constant Global Data” on page 7-57

Workflow

To convert MATLAB code that uses global data to fixed-point:

- 1 Declare the variables as global in your code.

For more information, see “Declare Global Variables” on page 7-56

- 2 Before using the global data, define and initialize it.

For more information, see “Define Global Data” on page 7-56.

- 3 Convert code to fixed-point from the Fixed-Point Converter or using `fiaccl`.

The Fixed-Point Converter always synchronizes global data between MATLAB and the generated MEX function.

Declare Global Variables

When using global data, you must first declare the global variables in your MATLAB code. This code shows the `use_globals` function, which uses two global variables, `AR` and `B`.

```
function y = use_globals(u)
%#codegen
% Declare AR and B as global variables
global AR;
global B;
AR(1) = u + B(1);
y = AR * 2;
```

Define Global Data

You can define global data in the MATLAB global workspace, in a Fixed-Point Converter project, or at the command line. If you do not initialize global data in a project or at the command line, the software looks for the variable in the MATLAB global workspace.

Define Global Data in the MATLAB Global Workspace

To convert the `use_globals` function described in “Declare Global Variables” on page 7-56, you must first define and initialize the global data.

```
global AR B;
AR = ones(4);
B=[1 2 3];
```


Define Global Data in a Fixed-Point Converter Project

- 1 On the **Define Input Types** page, after selecting and running a test file, select **Yes** next to **Does this code use global variables**.

By default, the Fixed-Point Converter names the first global variable in a project `g`.

- 2 Enter the names of the global variables used in your code. After adding a global variable, specify its type.
- 3 Click **Add global** to enter more global variables.

Note If you do not specify the type, you must create a variable with the same name in the global workspace.

Define Global Data at the Command Line

To define global data at the command line, use the `fiaccl -globals` option. For example, to convert the `use_globals` function described in “Declare Global Variables” on page 7-56 to fixed-point, specify two global inputs, `AR` and `B`, at the command line. Use the `-args` option to specify that the input `u` is a real, scalar double.

```
fiaccl -float2fixed cfg -global {'AR',ones(4),'B',[1 2 3]} use_globals -args {0}
```

Alternatively, specify the type and initial value with the `-globals` flag using the format `-globals {'g', {type, initial_value}}`.

To provide initial values for variable-size global data, specify the type and initial value with the `-globals` flag using the format `-globals {'g', {coder.typeof(0,[2 2],1),[1 1]}}`. For example, to specify a global variable `g` that has an initial value `[1 1]` and upper bound `[2 2]`, enter:

```
fiaccl -float2fixed cfg -global {'g', {coder.typeof(0,[2 2],1),[1 1]}} myfunction
```

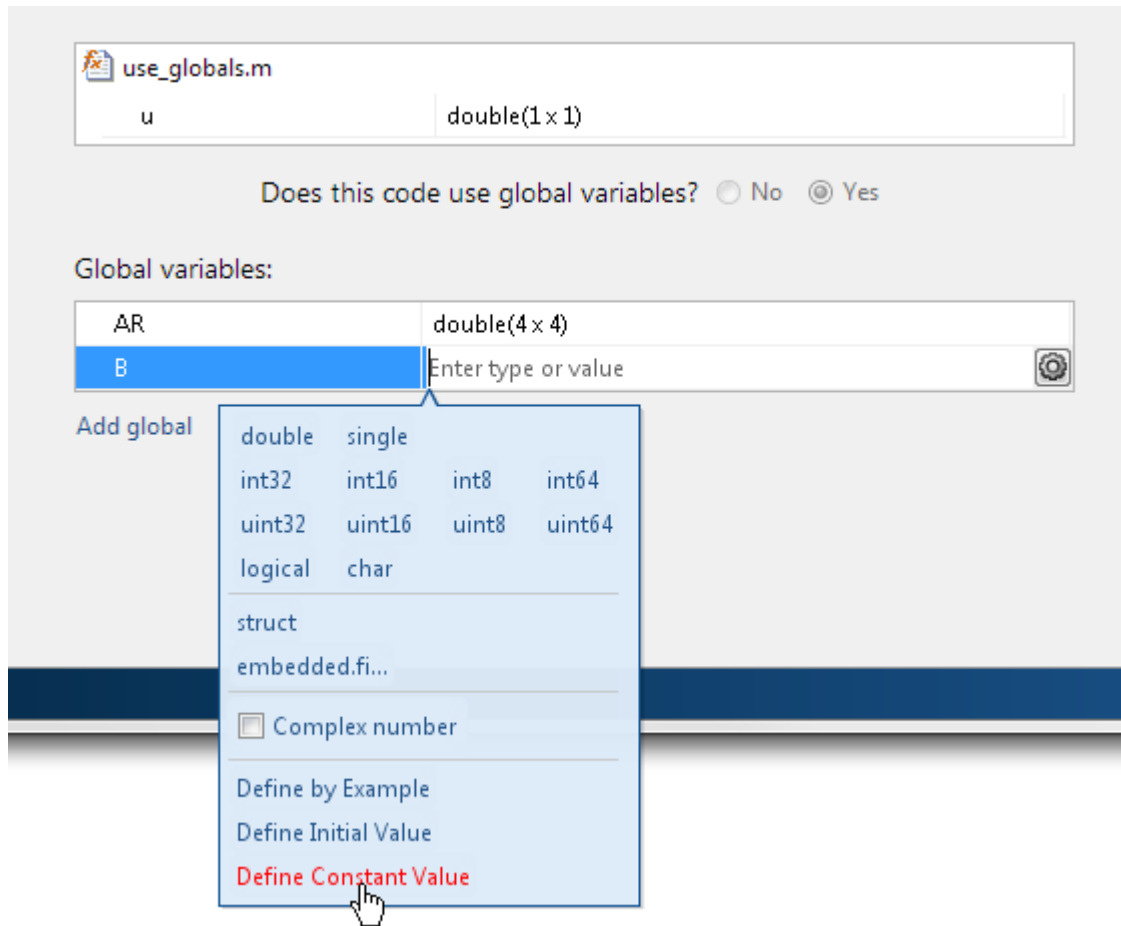
For a detailed explanation of the syntax, see `coder.typeof`.

Define Constant Global Data

If you know that the value of a global variable does not change at run time, you can reduce overhead in the fixed-point code by specifying that the global variable has a constant value. You cannot write to the constant global variable.

Define Constant Global Data in the Fixed-Point Converter

- 1 On the **Define Input Types** page, after selecting and running a test file, select **Yes** next to **Does this code use global variables**.
- 2 Enter the name of the global variables used in your code.
- 3 Click the field to the right of the global variable.
- 4 Select **Define Constant Value**.



- 5 In the field to the right of the constant global variable, enter a MATLAB expression.

Define Constant Global Data at the Command Line

To specify that a global variable is constant using the `fiaccel` command, use the `-globals` option with the `coder.Constant` class.

- 1 Define a fixed-point conversion configuration object.

```
cfg = coder.config('fixpt');
```

- 2 Use `coder.Constant` to specify that a global variable has a constant value. For example, this code specifies that the global variable `g` has an initial value 4 and that global variable `gc` has the constant value 42.

```
global_values = {'g', 4, 'gc', coder.Constant(42)};
```

- 3 Convert the code to fixed-point using the `-globals` option. For example, convert `myfunction` to fixed-point, specifying that the global variables are defined in the cell array `global_values`.

```
fiaccel -float2fixed cfg -global global_values myfunction
```

Constant Global Data in a Code Generation Report

The code generation report provides this information about a constant global variable:

- Type of Global on the **Variables** tab.
- Highlighted variable name in the **Function** pane.

See Also

Related Examples

- “Convert Code Containing Global Variables to Fixed-Point” on page 7-60

Convert Code Containing Global Variables to Fixed-Point

This example shows how to convert a MATLAB algorithm containing global variables to fixed point using the Fixed-Point Converter app.

- 1 In a local writable folder, create the function `use_globals.m`.

```
function y = use_globals(u)
    %#codegen
    % Declare AR and B as global variables
    global AR;
    global B;
    AR(1) = u + B(1);
    y = AR * 2;
```

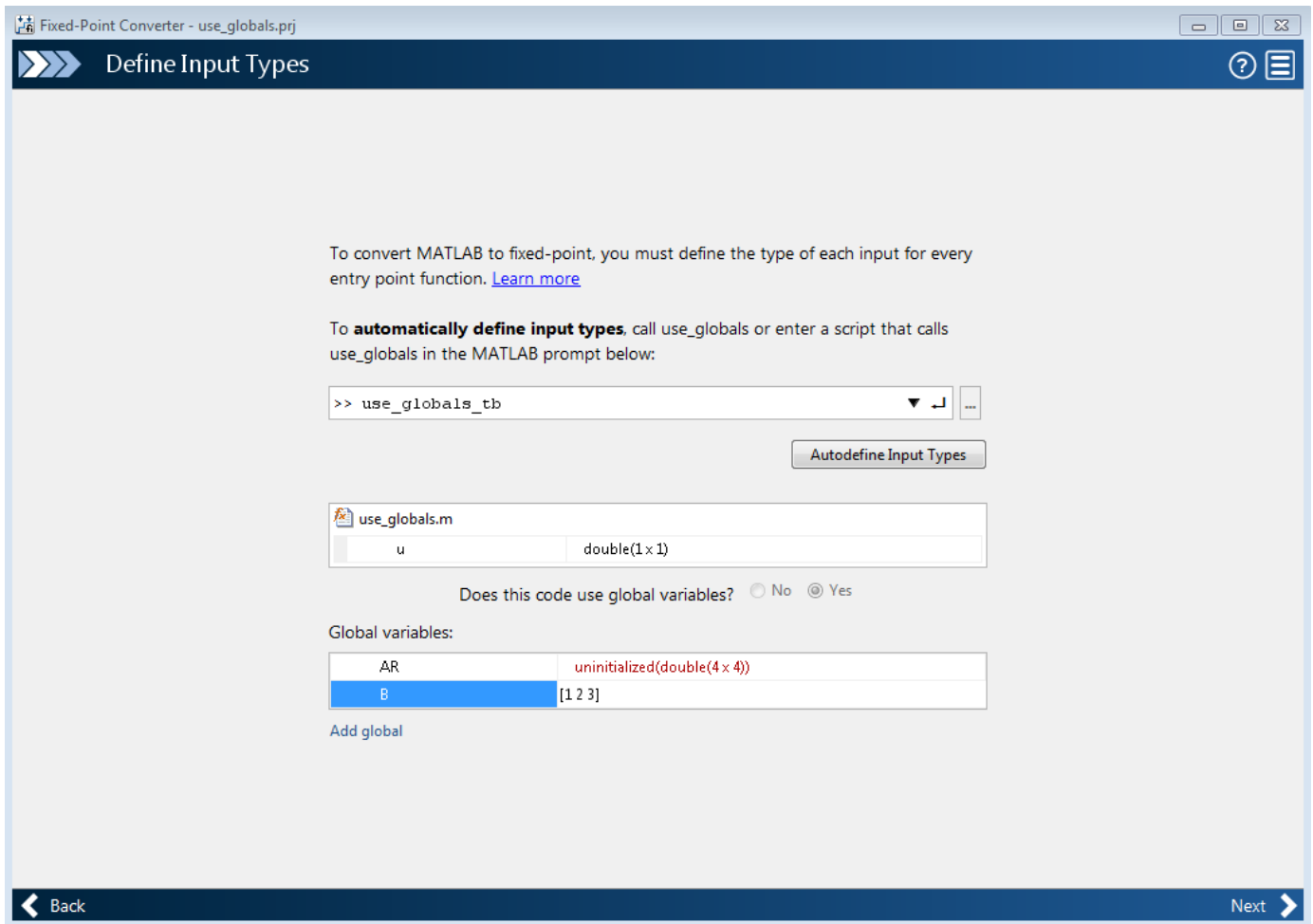
- 2 In the same folder, create a test file, `use_globals_tb.m` that calls the function.

```
u = 55;
global AR B;
AR = ones(4);
B=[1 2 3];
y = use_globals(u);
```

- 3 On the MATLAB toolstrip, in the **Apps** tab, under **Code Generation**, click the Fixed-Point Converter app icon.
- 4 To add the entry-point function, `use_globals.m` to the project, on the **Select Source Files** page, browse to the file, and click **Open**. Click **Next**.
- 5 On the **Define Input Types** page, add `use_globals_tb.m` as the test file. Click **Autodefine Input Types**.

The app determines from the test file that the input type of the input `u` is `double(1x1)`.

- 6 Select **Yes** next to **Does this code use global variables**. By default, the Fixed-Point Converter app names the first global variable in a project `g`.
- 7 Type in the names of the global variables in your code. In the field to the right of the global variable `AR`, specify its type as `double(4x4)`.
- 8 The global variable `B` is not assigned in the `use_globals` function. Define this variable as a global constant by clicking the field to the right of the constant and selecting **Define Constant Value**. Type in the value of `B` as it is defined in the test file, `[1 2 3]`. The app indicates that `B` has the value `[1 2 3]`. The app indicates that `AR` is not initialized.



- 9 Click **Next**. The app generates an instrumented MEX function for your entry-point MATLAB function. On the **Convert to Fixed-Point** page, click **Simulate** to simulate the function, gather range information, and get proposed data types.
- 10 Click **Convert** to accept the proposed data types and convert the function to fixed-point.

In the generated fixed-point code, the global variable AR is now AR_g.

The wrapper function contains three global variables: AR, AR_g, and B, where AR_g is set equal to a fi-casted AR, and AR is set equal to a double casted AR_g. The global variable B does not have a separate variable in the fixed-point code because it is a constant.

The screenshot shows the Fixed-Point Converter interface. The main window displays the source code for the 'use_globals' function. The code is as follows:

```

1 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
2 %
3 %       Generated by MATLAB 9.0 and Fixed-Point Designer 5.2
4 %
5 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
6 function y = use_globals_fixpt(u)
7 %#codegen
8 % Declare AR and B as global variables
9 fm = get_fimath();
10
11 global AR_g;
12 global B;
13 AR_g(1) = fi(u + B(1), 0, 6, 0, fm);
14 y = fi(AR_g * fi(2, 0, 2, 0, fm), 0, 7, 0, fm);
15 end
16
17
18 function fm = get_fimath()
19     fm = fimath('RoundingMethod', 'Floor', 'OverflowAction', 'Wrap', 'ProductMode', 'FullPrecision
20 end
21

```

Below the code, the Type Validation Output table is displayed:

| Variable | Type | Sim Min | Sim Max | Static Min | Static Max | Whole N... | Proposed Type |
|----------|-------------|---------|---------|------------|------------|------------|------------------|
| Input | | | | | | | |
| u | double | 55 | 55 | | | Yes | numeric(0, 6, 0) |
| Output | | | | | | | |
| y | 4 x 4 do... | 2 | 112 | | | Yes | numeric(0, 7, 0) |
| Global | | | | | | | |
| AR | 4 x 4 do... | 56 | 56 | | | Yes | numeric(0, 6, 0) |
| B | 1 x 3 do... | Unknown | Unknown | | | Yes | |

```

function y = use_globals_fixpt(u)
%#codegen
% Declare AR and B as global variables
fm = get_fimath();

global AR_g;
global B;
AR_g(1) = fi(u + B(1), 0, 6, 0, fm);
y = fi(AR_g * fi(2, 0, 2, 0, fm), 0, 7, 0, fm);
end

```

```

function fm = get_fimath()
    fm = fimath('RoundingMethod', 'Floor',...
        'OverflowAction', 'Wrap', 'ProductMode', 'FullPrecision',...
        'MaxProductWordLength', 128, 'SumMode', 'FullPrecision',...
        'MaxSumWordLength', 128);
end

```

See Also

More About

- “Convert Code Containing Global Data to Fixed Point” on page 7-56

Convert Code Containing Structures to Fixed Point

This example shows how to convert a MATLAB algorithm containing structures to fixed point using the Fixed-Point Converter app.

- 1 In a local writable folder, create the functions `struct_fcn.m`

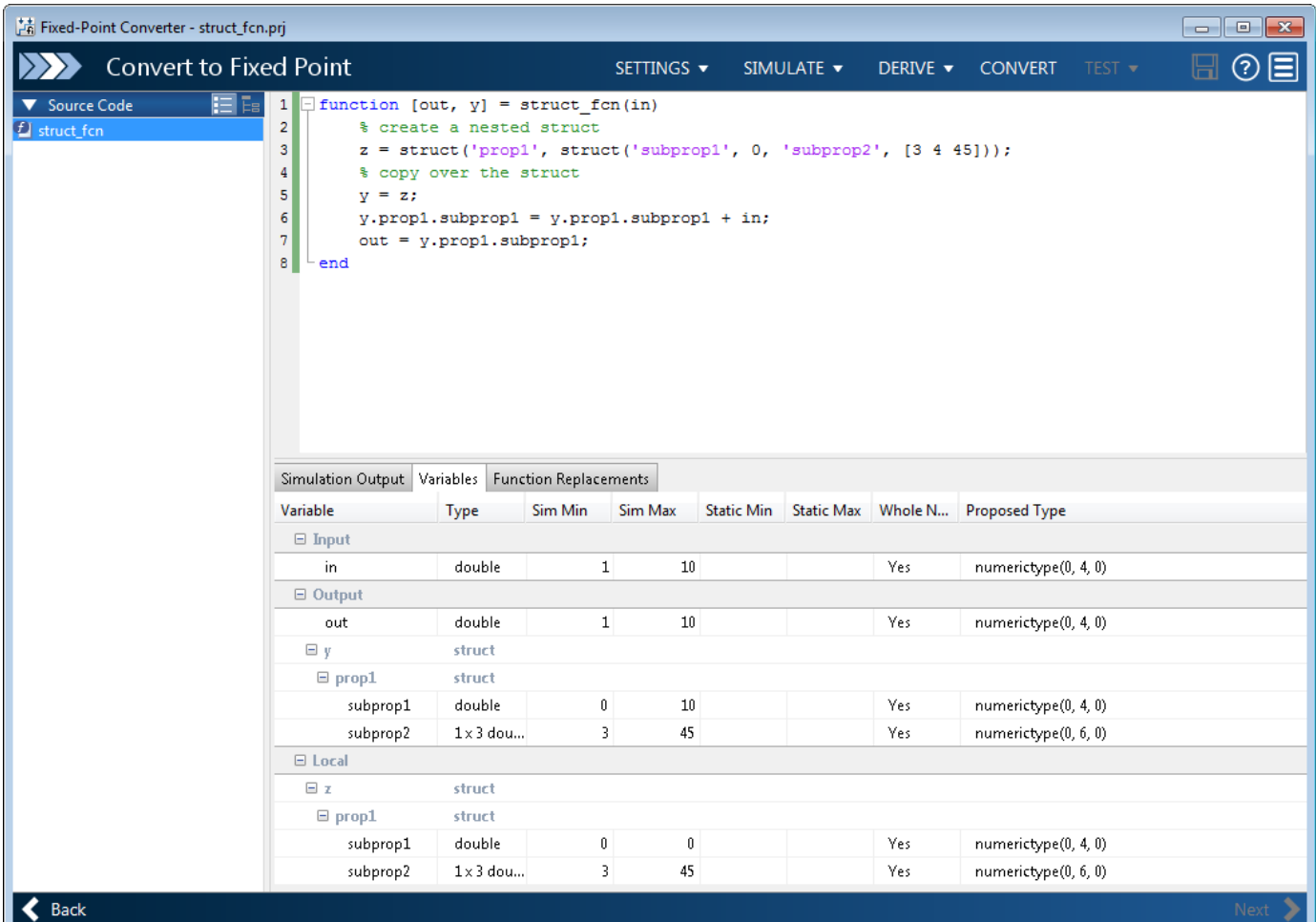
```
function [out, y] = struct_fcn(in)
    % create a nested struct
    z = struct('prop1', struct('subprop1', 0, 'subprop2', [3 4 45]));
    % copy over the struct
    y = z;
    y.prop1.subprop1 = y.prop1.subprop1 + in;
    out = y.prop1.subprop1;
end
```

- 2 In the same folder, create a test file, `struct_fcn_tb.m`, that calls the function.

```
for ii = 1:10
    struct_fcn(ii);
end
```

- 3 From the apps gallery, open the Fixed-Point Converter app.
- 4 On the **Select Source Files** page, browse to the `struct_fcn.m` file, and click **Open**.
- 5 Click **Next**. On the **Define Input Types** page, enter the test file that exercises the `struct_fcn` function. Browse to select the `struct_fcn_tb.m` file. Click **Autodefine Input Types**.
- 6 Click **Next**. The app generates an instrumented MEX function for your entry-point MATLAB function. On the **Convert to Fixed-Point** page, click **Simulate** to simulate the function, gather range information, and propose data types.

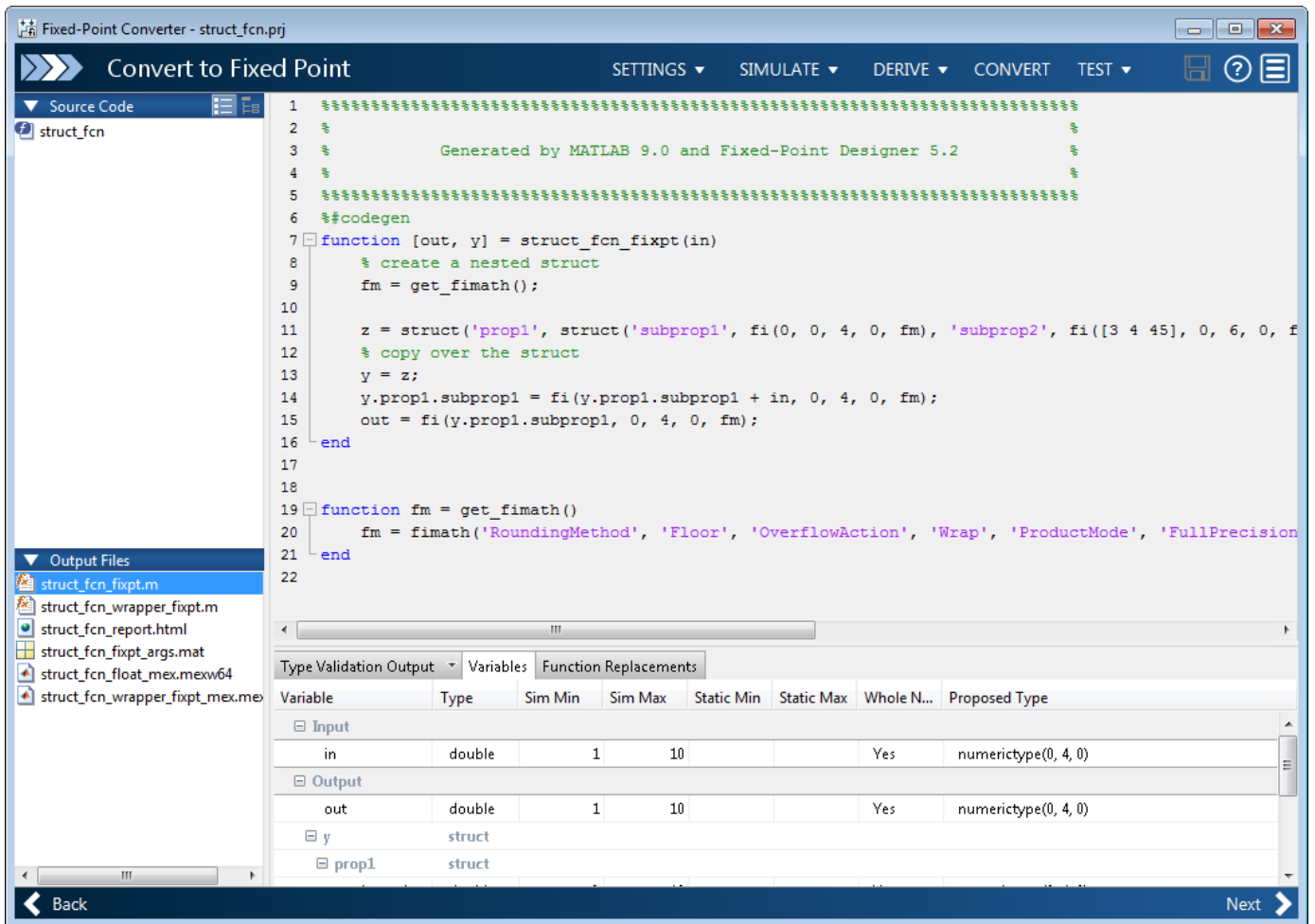
When the names, number, and types of fields of two or more structures match, the Fixed-Point Converter app proposes a unified type. In this example, the range of `z.prop1.subprop1` is $[0, 0]$, while the range of `y.prop1.subprop1` is $[0, 10]$. The app proposes a data type of `numeric(0, 4, 0)` for both `z.prop1.subprop1` and `y.prop1.subprop1` based on the union of the ranges of the two fields.



| Variable | Type | Sim Min | Sim Max | Static Min | Static Max | Whole N... | Proposed Type |
|----------|--------------|---------|---------|------------|------------|------------|----------------------|
| Input | | | | | | | |
| in | double | 1 | 10 | | | Yes | numerictype(0, 4, 0) |
| Output | | | | | | | |
| out | double | 1 | 10 | | | Yes | numerictype(0, 4, 0) |
| y | | | | | | | |
| prop1 | | | | | | | |
| subprop1 | double | 0 | 10 | | | Yes | numerictype(0, 4, 0) |
| subprop2 | 1 x 3 dou... | 3 | 45 | | | Yes | numerictype(0, 6, 0) |
| Local | | | | | | | |
| z | | | | | | | |
| prop1 | | | | | | | |
| subprop1 | double | 0 | 0 | | | Yes | numerictype(0, 4, 0) |
| subprop2 | 1 x 3 dou... | 3 | 45 | | | Yes | numerictype(0, 6, 0) |

7 Click **Convert**.

The Fixed-Point Converter converts the function containing the structures to fixed point and generates the `struct_fcn_fixpt.m` file.



See Also

More About

- “Convert Code Containing Global Data to Fixed Point” on page 7-56
- “Convert Code Containing Global Variables to Fixed-Point” on page 7-60
- “Convert Identical Functions Called with Different Data Types” on page 7-67

Convert Identical Functions Called with Different Data Types

This example shows how to convert a MATLAB algorithm containing specialized functions to fixed point using the **Fixed-Point Converter** app. The example also shows how each instance of the specialized function can have a different data type after conversion.

- 1 In a local writable folder, create the function `mtictac.m`. The input is a Tic-Tac-Toe game board. The output is the winning player. The output is 0 if there is no winner.

```
function a = mtictac(b)
% Input is a Tic-Tac-Toe board.
% The output is the winning player
% or zero if none.

% Copyright 2022 The MathWorks, Inc.

for i = 1:3
    a = winner(b(i,:));
    if a > 0
        return;
    end
    a = winner(b(:,i));
    if a > 0
        return;
    end
end

a = winner(diag(b));
if a>0
    return;
end
```

- 2 In the same folder, create the function, `winner.m`. The `winner.m` function is called by the function `mtictac.m`.

```
function a = winner(b)

% Copyright 2022 The MathWorks, Inc.

x = mean(b);

if x <= - 1
    a = 1;
elseif x >= 1
    a = 2;
else
    a = 0;
end
```

- 3 In the same folder, create a test file, `mtictac_tb.m`, that calls the `mtictac` function.

```
mtictac([-1 0 0; 0 -1 0; 0 0 0]);
mtictac([1 0 0; 0 1 0; 0 0 1]);
mtictac([0 0 -1; 0 0 -1; 0 0 -1]);
```

- 4 From the apps gallery, open the **Fixed-Point Converter** app.
- 5 Select `mtictac.m` as your entry-point function.

- 6 Click **Next**. On the **Define Input Types** page, enter the test file that exercises the `mtictac.m` function. Browse to select the `mtictac_tb.m` file. Click **Autodefine Input Types**.
- 7 Click **Next**. The app generates an instrumented MEX function for your entry-point MATLAB function.
- 8 Click **Analyze**. Observe that two instances have different data types. The data types for each individual call of the `winner` function can be customized to use different data types.

Fixed-Point Converter - mtictacprj.prj

Convert to Fixed Point

SETTINGS ▾ ANALYZE ▾ CONVERT TEST ▾

Source Code

```

1 function a = winner(b)
2 %
3
4 % Copyright 2021 The MathWorks, Inc.
5
6 x = mean(b);
7
8 if x <= - 1
9     a = 1;
10 elseif x >= 1
11     a = 2;
12 else
13     a = 0;
14 end
15

```

Output Files

- mtictac_fixpt.m
- mtictac_wrapper_fixpt.m
- mtictac_report.html
- mtictac_fixpt_args.mat
- mtictac_wrapper_fixpt_mex.mexw
- report.mldatx

| Variable | Type | Sim Min | Sim Max | Whole Nu... | Proposed Type | Log Data | Max Diff |
|----------|--------------|---------|---------|-------------|--------------------|----------|----------|
| Input | | | | | | | |
| b | 1 x 3 double | -1 | 1 | Yes | numeric(1, 2, 0) | ✓ | |
| Output | | | | | | | |
| a | double | 0 | 0 | Yes | numeric(0, 1, 0) | ✓ | |
| Local | | | | | | | |
| x | double | -0.33 | 0.33 | No | numeric(1, 16, 16) | | |

Back Next

The screenshot shows the Fixed-Point Converter interface. The source code for the `winner` function is as follows:

```

1 function a = winner(b)
2 %
3
4 % Copyright 2021 The MathWorks, Inc.
5
6 x = mean(b);
7
8 if x <= - 1
9     a = 1;
10 elseif x >= 1
11     a = 2;
12 else
13     a = 0;
14 end
15

```

The 'Output' tab of the conversion table is active, showing the following data:

| Variable | Type | Sim Min | Sim Max | Whole Nu... | Proposed Type | Log Data | Max Diff |
|----------|--------------|---------|---------|-------------|------------------------|----------|----------|
| Input | | | | | | | |
| b | 3 x 1 double | -1 | 1 | Yes | numerictype(1, 2, 0) | ✓ | |
| Output | | | | | | | |
| a | double | 0 | 2 | Yes | numerictype(0, 2, 0) | ✓ | |
| Local | | | | | | | |
| x | double | -1 | 1 | No | numerictype(1, 16, 14) | | |

9 Click **Convert**.

The **Fixed-Point Converter** app converts the function containing the structures to fixed point and generates the `mtictac_fixpt.m` file.

The screenshot shows the Fixed-Point Converter interface for a project named 'mtictacprj'. The main window displays the source code for a function named 'mtictac_fixpt'. The code is as follows:

```

1 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
2 %
3 %       Generated by MATLAB 9.14 and Fixed-Point Designer 7.5
4 %
5 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
6 %#codegen
7 function a = mtictac_fixpt(b)
8 % Input is a tic-tac-toe board. Output is the winning player or zero if
9 % none.
10 %
11 % Trivial example used to test F2F conversion.
12 %
13 % Copyright 2021 The MathWorks, Inc.
14
15 fm = get_fimath();
16
17 for i = 1:3
18     a = fi(winner_s1(b(i,:), 0, 2, 0, fm);
19     if a > fi(0, 0, 1, 0, fm)
20         return;
21     end
22     a(:) = winner_s2(b(:,i));
23     if a > fi(0, 0, 1, 0, fm)
24         return;
25     end
26 end

```

Below the code editor, the 'Output' tab is selected, showing a table of variable declarations and their fixed-point properties:

| Variable | Type | Size | Signed | Word Length | Fraction Length |
|----------|-----------------|-------|--------|-------------|-----------------|
| Input | | | | | |
| b | embedded.fi | 3 x 3 | Yes | 2 | 0 |
| Output | | | | | |
| a | embedded.fi | 1 x 1 | No | 2 | 0 |
| Local | | | | | |
| fm | embedded.fimath | 1 x 1 | No | | |
| i | double | 1 x 1 | No | | |

See Also

More About


- “Convert Code Containing Global Data to Fixed Point” on page 7-56
- “Convert Code Containing Global Variables to Fixed-Point” on page 7-60
- “Convert Code Containing Structures to Fixed Point” on page 7-64

Data Type Issues in Generated Code

Within the fixed-point conversion report, you have the option to highlight MATLAB code that results in double, single, or expensive fixed-point operations. Consider enabling these checks when trying to achieve a strict single, or fixed-point design.

These checks are disabled by default.

Enable the Highlight Option in the Fixed-Point Converter App

- 1 On the **Convert to Fixed Point** page, to open the **Settings** dialog box, click the **Settings** arrow .
- 2 Under **Plotting and Reporting**, set **Highlight potential data type issues** to Yes.

When conversion is complete, open the fixed-point conversion report to view the highlighting. Click **View report** in the **Type Validation Output** tab.

Enable the Highlight Option at the Command Line

- 1 Create a fixed-point code configuration object:

```
cfg = coder.config('fixpt');
```

- 2 Set the `HighlightPotentialDataTypeIssues` property of the configuration object to `true`.

```
cfg.HighlightPotentialDataTypeIssues = true;
```

Stowaway Doubles

When trying to achieve a strict-single or fixed-point design, manual inspection of code can be time-consuming and error prone. This check highlights all expressions that result in a double operation.

Stowaway Singles

This check highlights all expressions that result in a single operation.

Expensive Fixed-Point Operations

The expensive fixed-point operations check identifies optimization opportunities for fixed-point code. It highlights expressions in the MATLAB code that require cumbersome multiplication or division, expensive rounding, expensive comparison, or multiword operations. For more information on optimizing generated fixed-point code, see “Tips for Making Generated Code More Efficient” on page 49-9.

Cumbersome Operations

Cumbersome operations most often occur due to insufficient range of output. Avoid inputs to a multiply or divide operation that has word lengths larger than the base integer type of your processor. Operations with larger word lengths can be handled in software, but this approach requires much more code and is much slower.

Expensive Rounding

Traditional handwritten code, especially for control applications, almost always uses "no effort" rounding. For example, for unsigned integers and two's complement signed integers, shifting right and dropping the bits is equivalent to rounding to floor. To get results comparable to, or better than, what you expect from traditional handwritten code, use the `floor` rounding method. This check identifies expensive rounding operations in multiplication and division.

Expensive Comparison Operations

Comparison operations generate extra code when a casting operation is required to do the comparison. For example, when comparing an unsigned integer to a signed integer, one of the inputs must first be cast to the signedness of the other before the comparison operation can be performed. Consider optimizing the data types of the input arguments so that a cast is not required in the generated code.

Multiword Operations

Multiword operations can be inefficient on hardware. When an operation has an input or output data type larger than the largest word size of your processor, the generated code contains multiword operations. You can avoid multiword operations in the generated code by specifying local `fimath` properties for variables. You can also manually specify input and output word lengths of operations that generate multiword code.

System Objects Supported by Fixed-Point Converter App

Use the Fixed-Point Converter app to automatically propose and apply data types for commonly used system objects. The proposed data types are based on simulation data from the System object™.

Automated conversion is available for these DSP System Toolbox System Objects:

- `dsp.BiquadFilter`
- `dsp.FIRDecimator`
- `dsp.FIRInterpolator`
- `dsp.FIRFilter` (Direct Form and Direct Form Transposed only)
- `dsp.FIRRateConverter`
- `dsp.VariableFractionalDelay`

The Fixed-Point Converter app can display simulation minimum and maximum values, whole number information, and histogram data.

- You cannot propose data types for these System objects based on static range data.
- You must configure the System object to use 'Custom' fixed-point settings.
- The app applies the proposed data types only if the input signal is floating point, not fixed-point.

The app treats scaled doubles as fixed-point. The scaled doubles workflow for System objects is the same as that for regular variables.

- The app ignores the **Default word length** setting in the **Settings** menu. The app also ignores specified rounding and overflow modes. Data-type proposals are based on the settings of the System object.

See Also

Related Examples

- “Convert `dsp.FIRFilter` Object to Fixed-Point Using the Fixed-Point Converter App” on page 7-74

Convert dsp.FIRFilter Object to Fixed-Point Using the Fixed-Point Converter App

Convert a `dsp.FIRFilter` System object, which filters a high-frequency sinusoid signal, to fixed-point using the Fixed-Point Converter app. This example requires Fixed-Point Designer and DSP System Toolbox licenses.

Create DSP Filter Function and Test Bench

Create a `myFIRFilter` function from a `dsp.FIRFilter` System object.

By default, System objects are configured to use full-precision fixed-point arithmetic. To gather range data and get data type proposals from the Fixed-Point Converter app, configure the System object to use 'Custom' settings.

Save the function to a local writable folder.

```
function output = myFIRFilter(input, num)

    persistent lowpassFIR;
    if isempty(lowpassFIR)
        lowpassFIR = dsp.FIRFilter('NumeratorSource', 'Input port', ...
            'FullPrecisionOverride', false, ...
            'ProductDataType', 'Full precision', ... % default
            'AccumulatorDataType', 'Custom', ...
            'CustomAccumulatorDataType', numerictype(1,16,4), ...
            'OutputDataType', 'Custom', ...
            'CustomOutputDataType', numerictype(1,8,2));
    end
    output = lowpassFIR(input, num);

end
```

Create a test bench, `myFIRFilter_tb`, for the filter. The test bench generates a signal that gathers range information for conversion. Save the test bench.

```
% Test bench for myFIRFilter
% Remove high-frequency sinusoid using an FIR filter.

% Initialize
f1 = 1000;
f2 = 3000;
Fs = 8000;
Fcutoff = 2000;

% Generate input
SR = dsp.SineWave('Frequency',[f1,f2],'SampleRate',Fs,...
    'SamplesPerFrame',1024);

% Filter coefficients
num = fir1(130,Fcutoff/(Fs/2));

% Visualize input and output spectra
plot = spectrumAnalyzer('SampleRate',Fs,...
    'PlotAsTwoSidedSpectrum',false,...
```

```

        'ShowLegend',true,'YLimits',[-120 30],...
        'Title','Input Signal (Channel 1) Output Signal (Channel 2)');

% Stream
for k = 1:100
    input = sum(SR(),2); % Add the two sinusoids together
    filteredOutput = myFIRFilter(input, num); % Filter
    plot([input,filteredOutput]); % Visualize
end

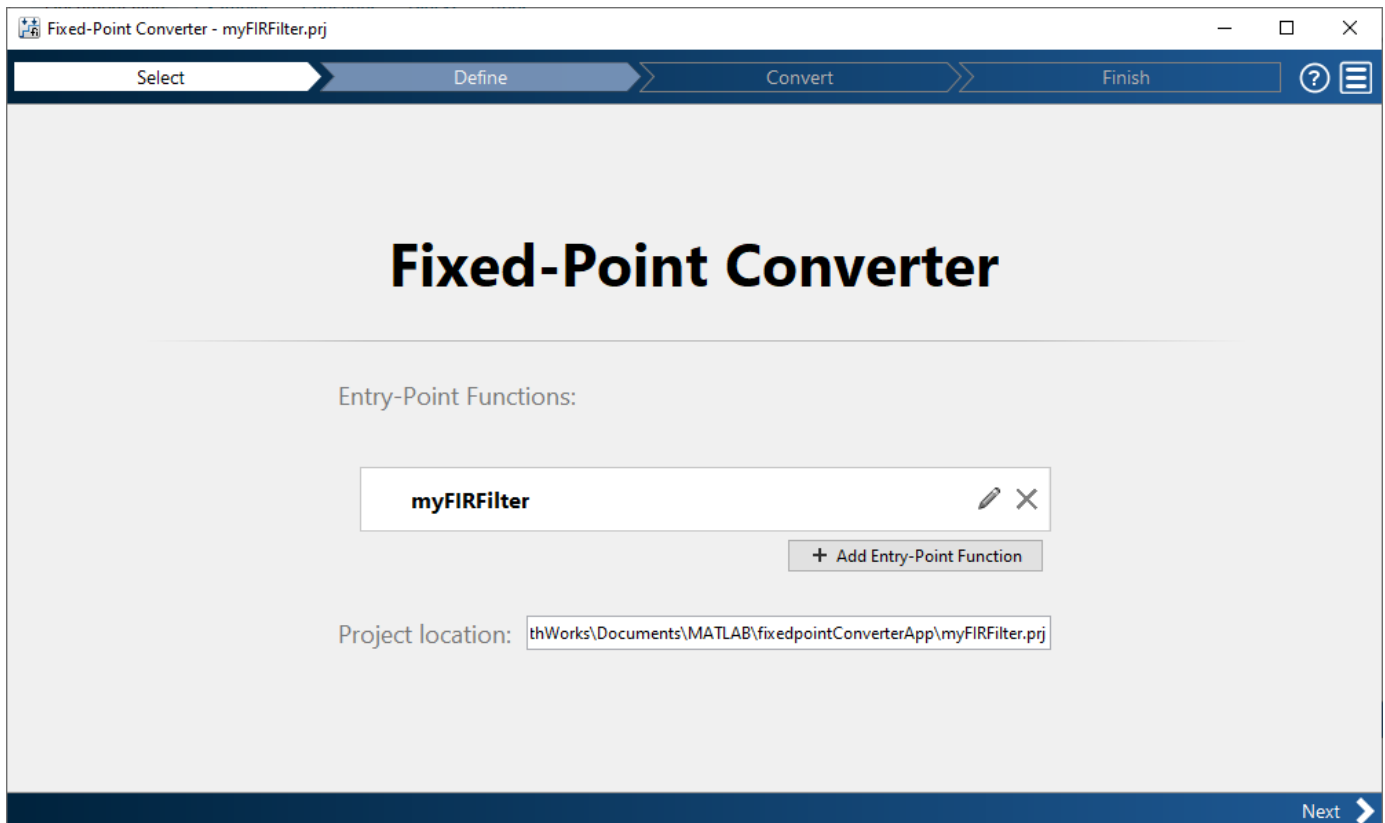
```

Convert the Function to Fixed-Point

- 1 Open the Fixed-Point Converter app.
 - MATLAB Toolstrip: On the **Apps** tab, under **Code Generation**, click the app icon.
 - MATLAB command prompt: Enter

```
fixedPointConverter
```
- 2 To add the entry-point function `myFIRFilter` to the project, browse to the file `myFIRFilter.m`, and then click **Open**.

By default, the app saves information and settings for this project in the current folder in a file named `myFirFilter.prj`.

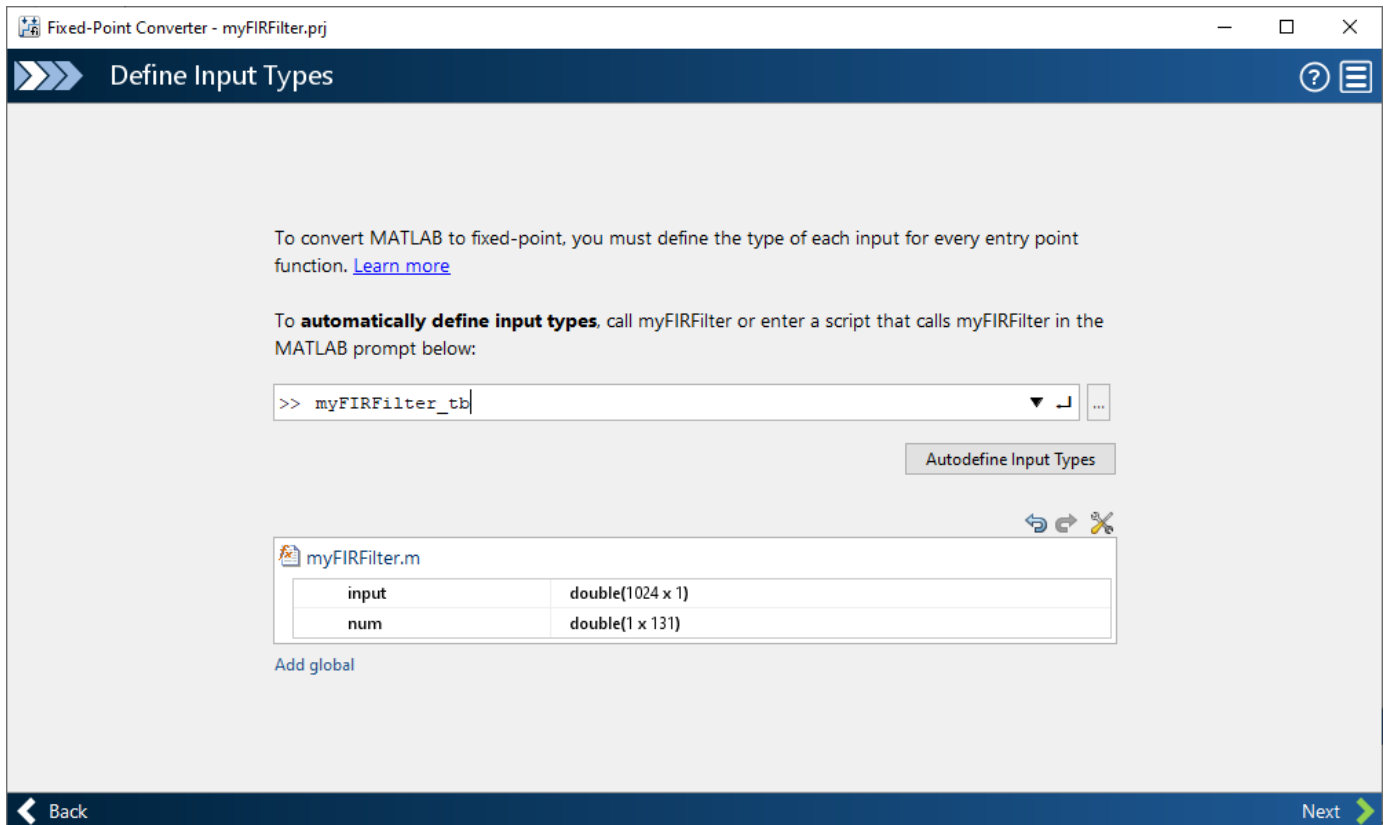


- 3 Click **Next** to go to the **Define Input Types** step.

The app screens `myFIRFilter.m` for code violations and fixed-point conversion readiness issues. The app does not find issues in `myFIRFilter.m`.

- 4 On the **Define Input Types** page, to add myFIRFilter_tb as a test file, browse to myFIRFilter_tb.m, and then click **Autodefine Input Types**.

The app determines from the test file that the type of input is double(1024 x 1) and the type of num is double(1 x 131).



- 5 Click **Next** to go to the **Convert to Fixed Point** step.
- 6 On the **Convert to Fixed Point** page, click **Analyze** to collect range information.

The screenshot shows the 'Fixed-Point Converter - myFIRFilter.prj' application window. The title bar includes standard window controls and the text 'Fixed-Point Converter - myFIRFilter.prj'. The main window has a dark blue header with the text 'Convert to Fixed Point' and several buttons: 'SETTINGS', 'ANALYZE' (highlighted with a red box), 'CONVERT', and 'TEST'. Below the header is a 'Source Code' tab showing the MATLAB code for the 'myFIRFilter' function. The code is as follows:

```

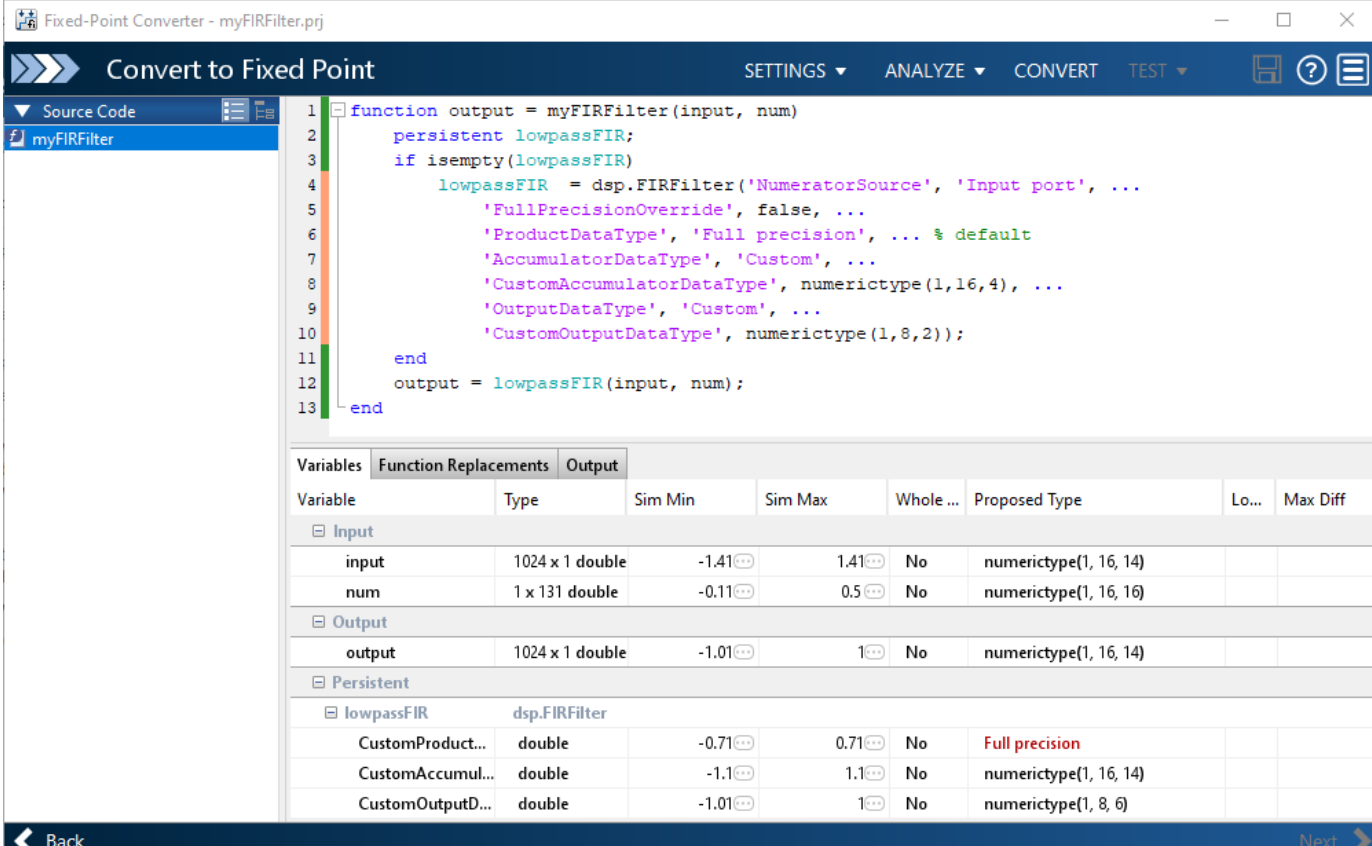
1 function output = myFIRFilter(input, num)
2     persistent lowpassFIR;
3     if isempty(lowpassFIR)
4         lowpassFIR = dsp.FIRFilter('NumeratorSource', 'Input port', ...
5             'FullPrecisionOverride', false, ...
6             'ProductDataType', 'Full precision', ... % default
7             'AccumulatorDataType', 'Custom', ...
8             'CustomAccumulatorDataType', numericity(1,16,4), ...
9             'OutputDataType', 'Custom', ...
10            'CustomOutputDataType', numericity(1,8,2));
11     end
12     output = lowpassFIR(input, num);
13 end
    
```

Below the code editor is a 'Variables' tab with a table showing the analysis results. The table has columns for 'Variable', 'Type', 'Sim Min', 'Sim Max', 'Whole ...', 'Proposed Type', 'Lo...', and 'Max Diff'. The data is as follows:

| Variable | Type | Sim Min | Sim Max | Whole ... | Proposed Type | Lo... | Max Diff |
|--------------------------|-----------------|---------|---------|-----------|---------------|-------|----------|
| Input | | | | | | | |
| input | 1024 x 1 double | | | No | | | |
| num | 1 x 131 double | | | No | | | |
| Output | | | | | | | |
| output | 1024 x 1 double | | | No | | | |
| Persistent | | | | | | | |
| lowpassFIR dsp.FIRFilter | | | | | | | |
| CustomProduct... | double | | | No | | | |
| CustomAccumul... | double | | | No | | | |
| CustomOutputD... | double | | | No | | | |

At the bottom of the window, there are 'Back' and 'Next' navigation buttons.

The **Variables** tab displays the collected range information and type proposals. Manually edit the data type proposals as needed.



Fixed-Point Converter - myFIRFilter.pj

Convert to Fixed Point

SETTINGS ANALYZE CONVERT TEST

```

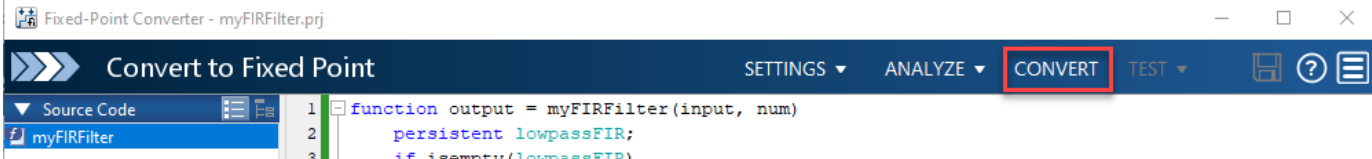
1 function output = myFIRFilter(input, num)
2     persistent lowpassFIR;
3     if isempty(lowpassFIR)
4         lowpassFIR = dsp.FIRFilter('NumeratorSource', 'Input port', ...
5             'FullPrecisionOverride', false, ...
6             'ProductDataType', 'Full precision', ... % default
7             'AccumulatorDataType', 'Custom', ...
8             'CustomAccumulatorDataType', numerictype(1,16,4), ...
9             'OutputDataType', 'Custom', ...
10            'CustomOutputDataType', numerictype(1,8,2));
11     end
12     output = lowpassFIR(input, num);
13 end

```

| Variable | Type | Sim Min | Sim Max | Whole ... | Proposed Type | Lo... | Max Diff |
|--------------------------|-----------------|---------|---------|-----------|------------------------|-------|----------|
| Input | | | | | | | |
| input | 1024 x 1 double | -1.41 | 1.41 | No | numerictype(1, 16, 14) | | |
| num | 1 x 131 double | -0.11 | 0.5 | No | numerictype(1, 16, 16) | | |
| Output | | | | | | | |
| output | 1024 x 1 double | -1.01 | 1 | No | numerictype(1, 16, 14) | | |
| Persistent | | | | | | | |
| lowpassFIR dsp.FIRFilter | | | | | | | |
| CustomProduct... | double | -0.71 | 0.71 | No | Full precision | | |
| CustomAccumul... | double | -1.1 | 1.1 | No | numerictype(1, 16, 14) | | |
| CustomOutputD... | double | -1.01 | 1 | No | numerictype(1, 8, 6) | | |

Back Next

7 Click **Convert** to apply the proposed data types to the function.



Fixed-Point Converter - myFIRFilter.pj

Convert to Fixed Point

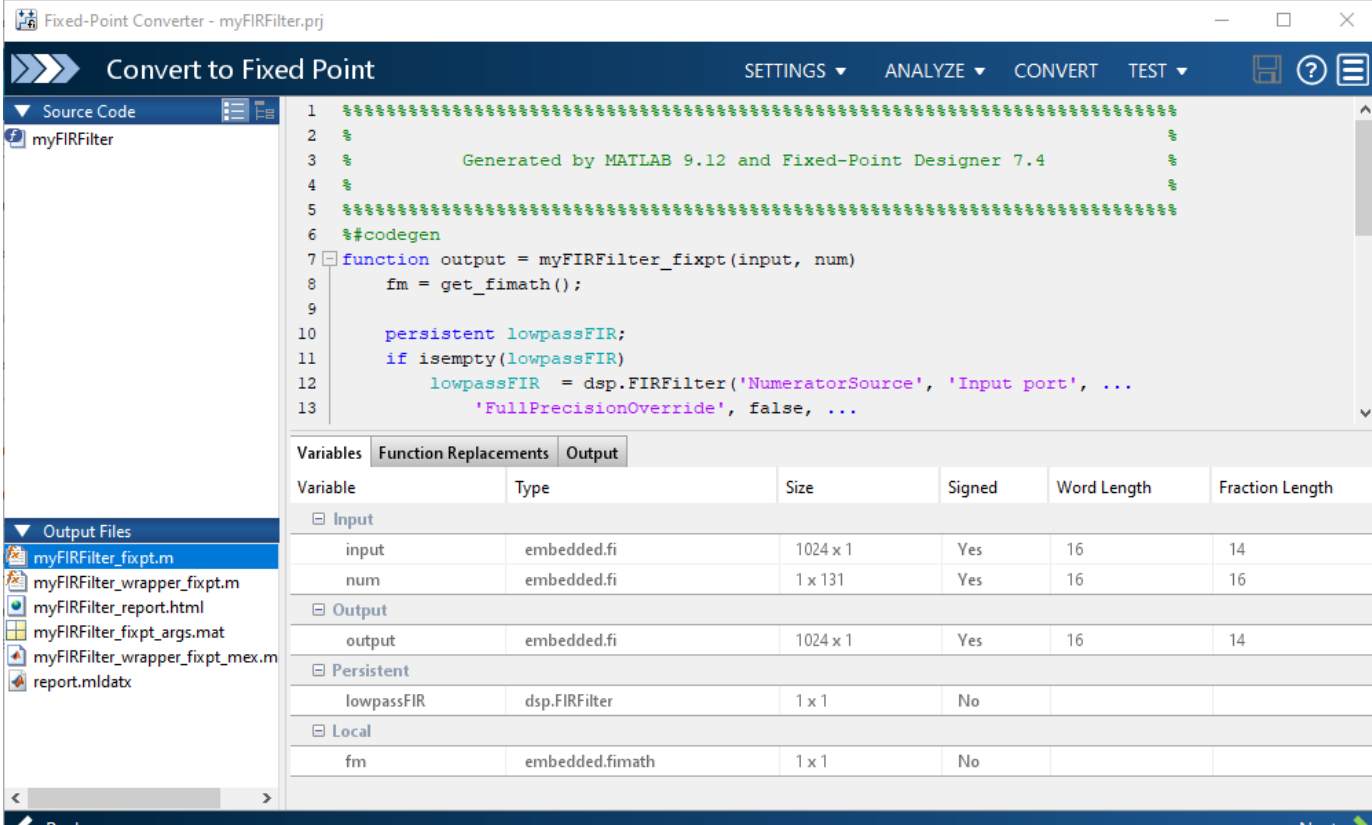
SETTINGS ANALYZE **CONVERT** TEST

```

1 function output = myFIRFilter(input, num)
2     persistent lowpassFIR;
3     if isempty(lowpassFIR)

```

The Fixed-Point Converter app applies the proposed data types and generates a fixed-point function, `myFIRFilter_fixpt`.



The screenshot shows the Fixed-Point Converter app interface. The main window displays the generated code for the function `myFIRFilter_fixpt`. The code is as follows:

```

1  %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
2  %
3  %       Generated by MATLAB 9.12 and Fixed-Point Designer 7.4
4  %
5  %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
6  %#codegen
7  function output = myFIRFilter_fixpt(input, num)
8      fm = get_fimath();
9
10     persistent lowpassFIR;
11     if isempty(lowpassFIR)
12         lowpassFIR = dsp.FIRFilter('NumeratorSource', 'Input port', ...
13             'FullPrecisionOverride', false, ...

```

Below the code, a table lists the variables and their properties:

| Variable | Type | Size | Signed | Word Length | Fraction Length |
|------------|-----------------|----------|--------|-------------|-----------------|
| Input | | | | | |
| input | embedded.fi | 1024 x 1 | Yes | 16 | 14 |
| num | embedded.fi | 1 x 131 | Yes | 16 | 16 |
| Output | | | | | |
| output | embedded.fi | 1024 x 1 | Yes | 16 | 14 |
| Persistent | | | | | |
| lowpassFIR | dsp.FIRFilter | 1 x 1 | No | | |
| Local | | | | | |
| fm | embedded.fimath | 1 x 1 | No | | |

```

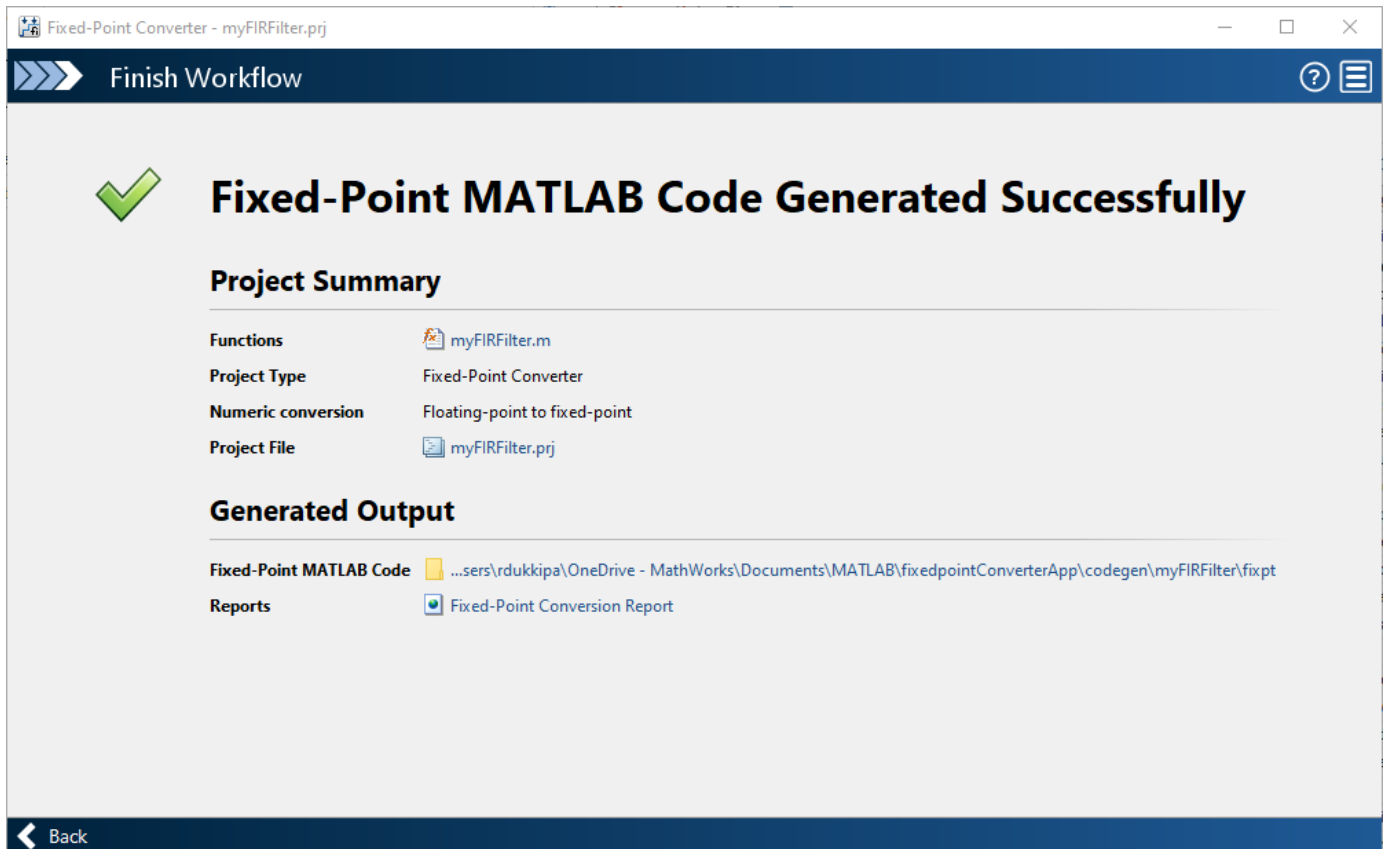
%#codegen
function output = myFIRFilter_fixpt(input, num)
    fm = get_fimath();

    persistent lowpassFIR;
    if isempty(lowpassFIR)
        lowpassFIR = dsp.FIRFilter('NumeratorSource', 'Input port', ...
            'FullPrecisionOverride', false, ...
            'ProductDataType', 'Full precision', ... % default
            'AccumulatorDataType', 'Custom', ...
            'CustomAccumulatorDataType', numerictype(1, 16, 14), ...
            'OutputDataType', 'Custom', ...
            'CustomOutputDataType', numerictype(1, 8, 6));
    end
    output = fi(lowpassFIR(input, num), 1, 16, 14, fm);
end

function fm = get_fimath()
    fm = fimath('RoundingMethod', 'Floor', ...
        'OverflowAction', 'Wrap', ...
        'ProductMode', 'FullPrecision', ...
        'MaxProductWordLength', 128, ...
        'SumMode', 'FullPrecision', ...
        'MaxSumWordLength', 128);
end

```

- 8 Click **Next** to see the project summary details and links to the fixed-point MATLAB code and conversion report.



See Also

More About

- “System Objects Supported by Fixed-Point Converter App” on page 7-73

Automated Conversion Using Fixed-Point Converter App

- “Specify Type Proposal Options” on page 8-2
- “Detect Overflows” on page 8-5
- “Propose Data Types Based on Simulation Ranges” on page 8-13
- “Propose Data Types Based on Derived Ranges” on page 8-24
- “View and Modify Variable Information” on page 8-35
- “Replace the exp Function with a Lookup Table” on page 8-38
- “Convert Fixed-Point Conversion Project to MATLAB Scripts” on page 8-45
- “Replace a Custom Function with a Lookup Table” on page 8-47
- “Visualize Differences Between Floating-Point and Fixed-Point Results” on page 8-52
- “Enable Plotting Using the Simulation Data Inspector” on page 8-62
- “Add Global Variables by Using the App” on page 8-63
- “Automatically Define Input Types by Using the App” on page 8-64
- “Define Constant Input Parameters Using the App” on page 8-65
- “Define or Edit Input Parameter Type by Using the App” on page 8-66
- “Define Input Parameter by Example by Using the App” on page 8-73
- “Specify Global Variable Type and Initial Value Using the App” on page 8-80
- “Specify Properties of Entry-Point Function Inputs Using the App” on page 8-83
- “Detect Unexecuted and Constant-Folded Code” on page 8-84

Specify Type Proposal Options

To view type proposal options, in the Fixed-Point Converter app, on the **Convert to Fixed Point** page, click the **Settings** arrow .

The following options are available.

| Basic Type Proposal Settings | Values | Description |
|--------------------------------|--|---|
| Fixed-point type proposal mode | Propose fraction lengths for specified word length | Use the specified word length for data type proposals and propose the minimum fraction lengths to avoid overflows. |
| | Propose word lengths for specified fraction length (default) | Use the specified fraction length for data type proposals and propose the minimum word lengths to avoid overflows. |
| Default word length | 16 (default) | Default word length to use when Fixed-point type proposal mode is set to Propose fraction lengths for specified word lengths |
| Default fraction length | 4 (default) | Default fraction length to use when Fixed-point type proposal mode is set to Propose word lengths for specified fraction lengths |

| Advanced Type Proposal Settings | Values | Description |
|---------------------------------|----------------------------------|--|
| When proposing types | ignore simulation ranges | Propose data types based on derived ranges. |
| | ignore derived ranges | Propose data types based on simulation ranges. |
| | use all collected data (default) | Propose data types based on both simulation and derived ranges. |
| Propose target container types | Yes | Propose data type with the smallest word length that can represent the range and is suitable for C code generation (8,16,32, 64 ...). For example, for a variable with range [0. . 7], propose a word length of 8 rather than 3. |
| | No (default) | Propose data types with the minimum word length needed to represent the value. |

| Advanced Type Proposal Settings | Values | Description |
|-----------------------------------|---------------------|---|
| Optimize whole numbers | No | Do not use integer scaling for variables that were whole numbers during simulation. |
| | Yes (default) | Use integer scaling for variables that were whole numbers during simulation. |
| Signedness | Automatic (default) | Proposes signed and unsigned data types depending on the range information for each variable. |
| | Signed | Propose signed data types. |
| | Unsigned | Propose unsigned data types. |
| Safety margin for sim min/max (%) | 0 (default) | Specify safety factor for simulation minimum and maximum values. The simulation minimum and maximum values are adjusted by the percentage designated by this parameter, allowing you to specify a range different from that obtained from the simulation run. For example, a value of 55 specifies that you want a range at least 55 percent larger. A value of -15 specifies that a range up to 15 percent smaller is acceptable. |
| Search paths | ' ' (default) | Add paths to the list of paths to search for MATLAB files. Separate list items with a semicolon. |

| fimath Settings | Values | Description |
|-----------------|-------------------------|--|
| Rounding method | Ceiling | Specify the fimath properties for the generated fixed-point data types. The default fixed-point math properties use the Floor rounding and Wrap overflow. These settings generate the most efficient code but might cause problems with overflow. |
| | Convergent | |
| | Floor (default) | |
| | Nearest | |
| | Round | |
| | Zero | |
| Overflow action | Saturate | After code generation, if required, modify these settings to optimize the generated code, or example, avoid overflow or eliminate bias, and then rerun the verification. |
| | Wrap (default) | |
| Product mode | FullPrecision (default) | |
| | KeepLSB | |
| | KeepMSB | |
| | SpecifyPrecision | |
| Sum mode | FullPrecision (default) | |
| | KeepLSB | |

| fimath Settings | Values | Description |
|------------------------|------------------|---|
| | KeepMSB | For more information on <code>fimath</code> properties, see “ <code>fimath</code> Object Properties” on page 3-4. |
| | SpecifyPrecision | |

| Generated File Settings | Value | Description |
|--|------------------|--|
| Generated fixed-point file name suffix | _fixpt (default) | Specify the suffix to add to the generated fixed-point file names. |

| Plotting and Reporting Settings | Values | Description |
|--|---------------|---|
| Custom plot function | ' ' (default) | Specify the name of a custom plot function to use for comparison plots. |
| Plot with Simulation Data Inspector | No (default) | Specify whether to use the Simulation Data Inspector for comparison plots. |
| | Yes | |
| Highlight potential data type issues | No (default) | Specify whether to highlight potential data types in the generated html report. If this option is turned on, the report highlights single-precision, double-precision, and expensive fixed-point operation usage in your MATLAB code. |
| | Yes | |

Detect Overflows

This example shows how to detect overflows using the Fixed-Point Converter app. At the numerical testing stage in the conversion process, you choose to simulate the fixed-point code using scaled doubles. The app then reports which expressions in the generated code produce values that overflow the fixed-point data type.

Prerequisites

This example requires the following products:

- MATLAB
- Fixed-Point Designer
- C compiler

See https://www.mathworks.com/support/compilers/current_release/.

You can use `mex -setup` to change the default compiler. See “Change Default Compiler”.

Create a New Folder and Copy Relevant Files

- 1 In a local, writable folder, create a function `overflow.m`.

```
function y = overflow(b,x,reset)
    if nargin<3, reset = true; end
    persistent z p
    if isempty(z) || reset
        p = 0;
        z = zeros(size(b));
    end
    [y,z,p] = fir_filter(b,x,z,p);
end
function [y,z,p] = fir_filter(b,x,z,p)
    y = zeros(size(x));
    nx = length(x);
    nb = length(b);
    for n = 1:nx
        p=p+1; if p>nb, p=1; end
        z(p) = x(n);
        acc = 0;
        k = p;
        for j=1:nb
            acc = acc + b(j)*z(k);
            k=k-1; if k<1, k=nb; end
        end
        y(n) = acc;
    end
end
```

- 2 Create a test file, `overflow_test.m`, to exercise the overflow algorithm. You use this test file to define input types for `b`, `x`, and `reset`, and, later, to verify the fixed-point version of the algorithm.

```
function overflow_test
    % The filter coefficients were computed
    % using the FIR1 function from
    % Signal Processing Toolbox.
```

```

% b = fir1(11,0.25);
b = [-0.004465461051254
     -0.004324228005260
      +0.012676739550326
      +0.074351188907780
      +0.172173206073645
      +0.249588554524763
      +0.249588554524763
      +0.172173206073645
      +0.074351188907780
      +0.012676739550326
      -0.004324228005260
      -0.004465461051254]';

% Input signal
nx = 256;
t = linspace(0,10*pi,nx)';

% Impulse
x_impulse = zeros(nx,1); x_impulse(1) = 1;

% Max Gain
% The maximum gain of a filter will occur when the
% inputs line up with the signs of the filter's
% impulse response.
x_max_gain = sign(b)';
x_max_gain = repmat(x_max_gain,ceil(nx/length(b)),1);
x_max_gain = x_max_gain(1:nx);

% Sums of sines
f0=0.1; f1=2;
x_sines = sin(2*pi*t*f0) + 0.1*sin(2*pi*t*f1);

% Chirp
f_chirp = 1/16; % Target frequency
x_chirp = sin(pi*f_chirp*t.^2); % Linear chirp

x = [x_impulse,x_max_gain,x_sines,x_chirp];
titles = {'Impulse','Max gain','Sum of sines','Chirp'};
y = zeros(size(x));

for i=1:size(x,2)
    reset = true;
    y(:,i) = overflow(b,x(:,i),reset);
end

test_plot(1,titles,t,x,y)

end
function test_plot(fig,titles,t,x,y1)
figure(fig)
clf
sub_plot = 1;
font_size = 10;
for i=1:size(x,2)
    subplot(4,1,sub_plot)
    sub_plot = sub_plot+1;
    plot(t,x(:,i),'c',t,y1(:,i),'k')
end

```

```

        axis('tight')
        xlabel('t','FontSize',font_size);
        title(titles{i},'FontSize',font_size);
        ax = gca;
        ax.FontSize = 10;
    end
    figure(gcf)
end

```

It is a best practice is to create a separate test script to do pre- and post-processing, such as:

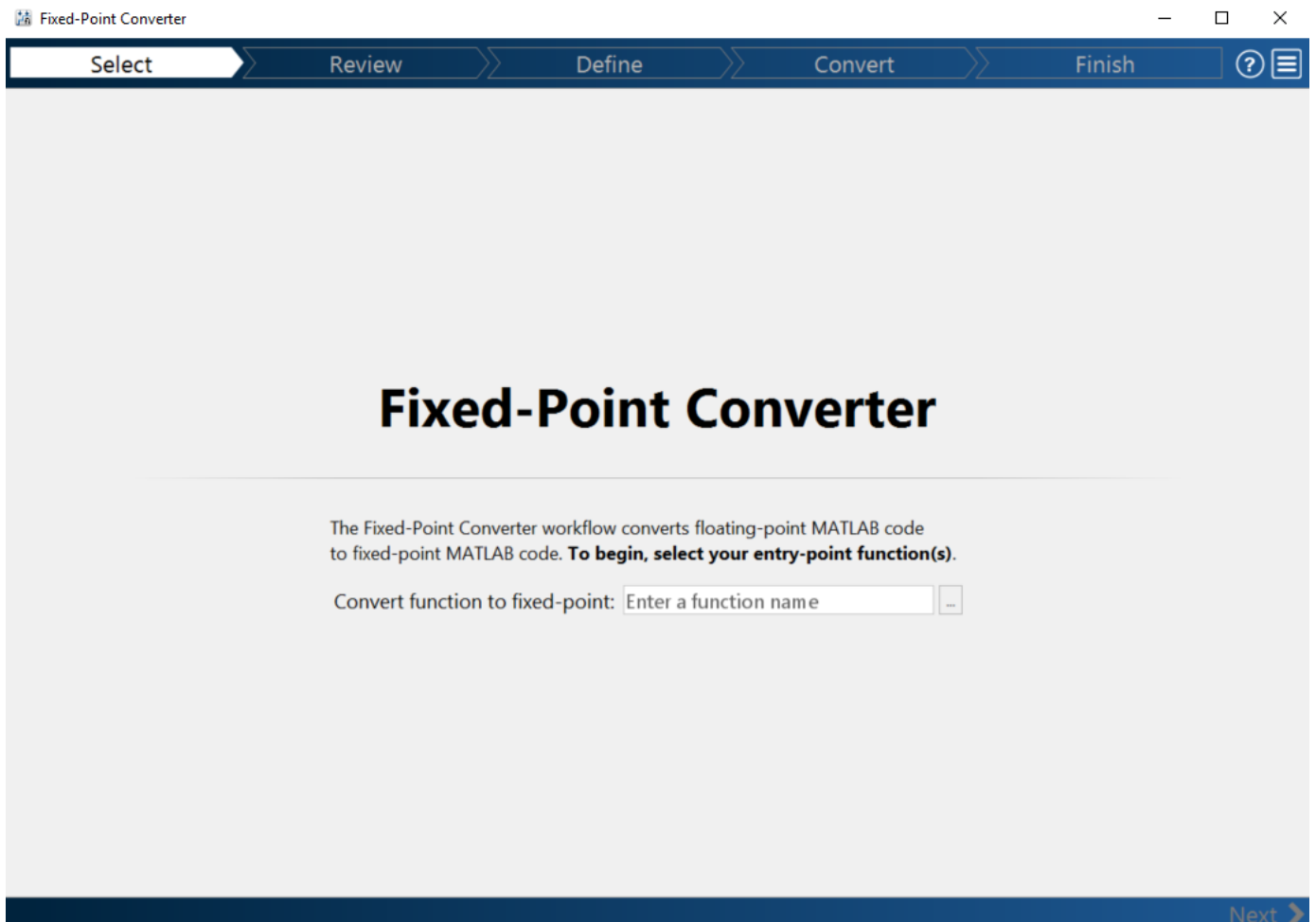
- Loading inputs.
- Setting up input values.
- Outputting test results.

For more information, see “Create a Test File” on page 11-3.

| Type | Name | Description |
|---------------|-----------------|-------------------------------------|
| Function code | overflow.m | Entry-point MATLAB function |
| Test file | overflow_test.m | MATLAB script that tests overflow.m |

Open the Fixed-Point Converter App

- 1 Navigate to the work folder that contains the file for this example.
- 2 On the MATLAB Toolstrip **Apps** tab, under **Code Generation**, click the app icon.



Select Source Files

- 1 To add the entry-point function `overflow` to the project, browse to the file `overflow.m`, and then click **Open**. By default, the app saves information and settings for this project in the current folder in a file named `overflow.prj`.
- 2 Click **Next** to go to the **Define Input Types** step.

The app screens `overflow.m` for code violations and fixed-point conversion readiness issues. The app does not find issues in `overflow.m`.

Define Input Types

- 1 On the **Define Input Types** page, to add `overflow_test` as a test file, browse to `overflow_test.m`, and then click **Open**.
- 2 Click **Autodefine Input Types**.

The test file runs. The app determines from the test file that the input type of `b` is `double(1x12)`, `x` is `double(256x1)`, and `reset` is `logical(1x1)`.

To **automatically define input types**, call `overflow` or enter a script that calls `overflow` in the MATLAB prompt below:

The screenshot shows the MATLAB interface for the 'overflow' app. At the top, a text box contains the command `>> overflow_test`. Below this is a button labeled 'Autodefine Input Types'. Underneath is a table titled 'overflow.m' with the following content:

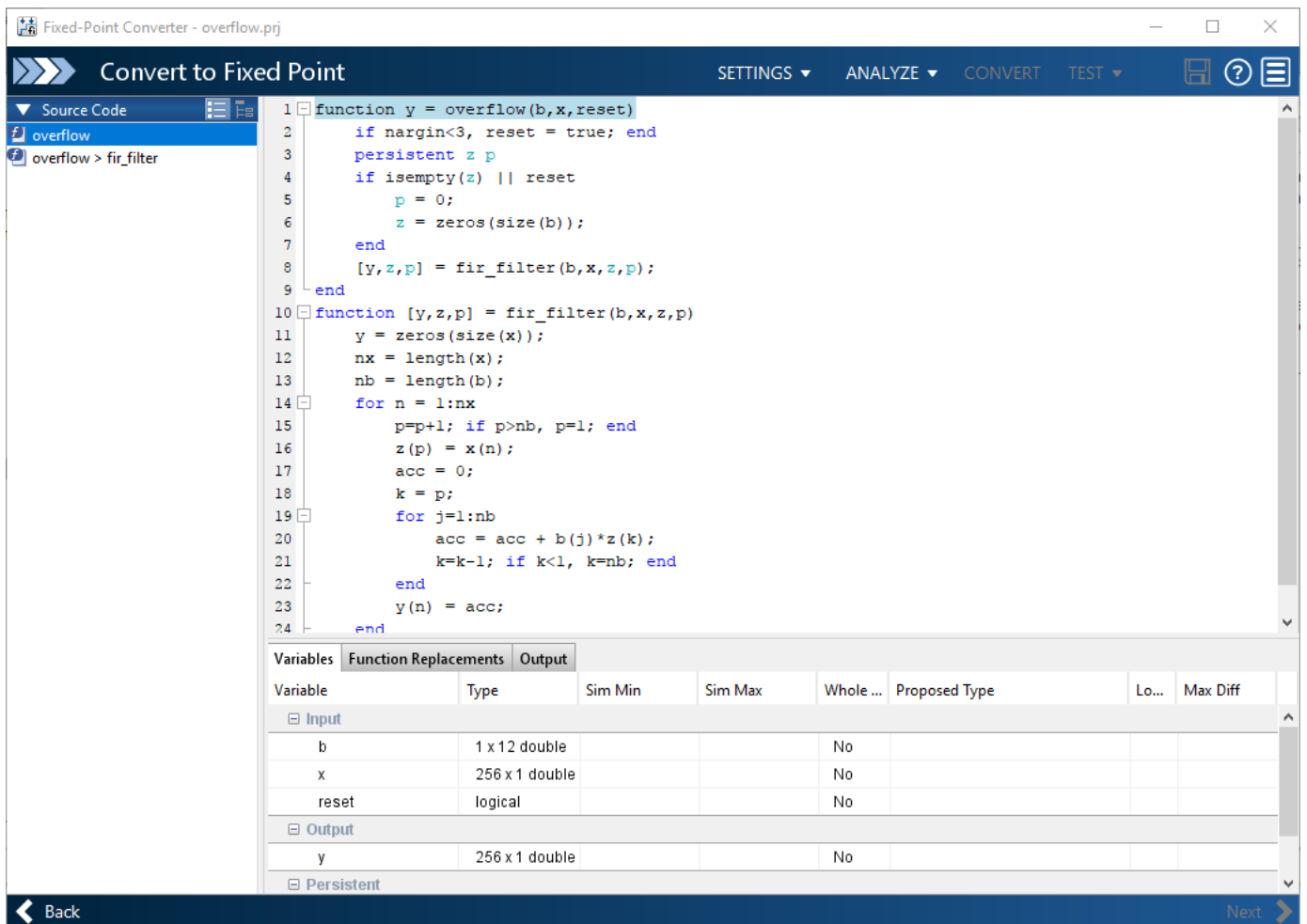
| | |
|-------|-----------------|
| b | double(1 x 12) |
| x | double(256 x 1) |
| reset | logical(1 x 1) |

Below the table is a link labeled 'Add global'.


- 3 Click **Next** to go to the **Convert to Fixed Point** step.

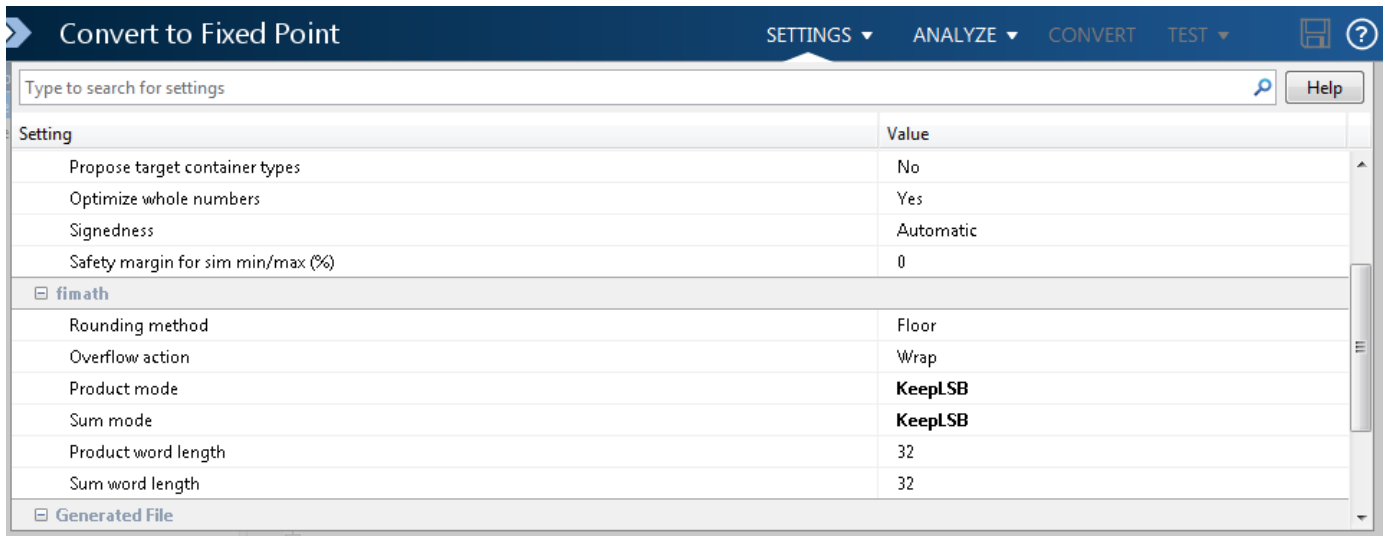
Convert to Fixed Point

- 1 The app generates an instrumented MEX function for your entry-point MATLAB function. The app displays compiled information — type, size, and complexity — for variables in your code. For more information, see “View and Modify Variable Information” on page 8-35.



On the **Function Replacements** tab the app displays functions that are not supported for fixed-point conversion. See “Running a Simulation” on page 7-8.

- To view the fimath settings, click the **Settings** arrow . Set the fimath **Product mode** and **Sum mode** to KeepLSB. These settings model the behavior of integer operations in the C language.



3 Click **Analyze**.

The test file, `overflow_test`, runs. The app displays simulation minimum and maximum ranges on the **Variables** tab. Using the simulation range data, the software proposes fixed-point types for each variable based on the default type proposal settings, and displays them in the **Proposed Type** column.

| Variable | Type | Sim Min | Sim Max | Whole Number | Proposed Type |
|-------------------|----------------|---------|---------|--------------|------------------------|
| Input | | | | | |
| b | 1 × 12 double | 0 | 0.25 | No | numerictype(1, 16, 17) |
| x | 256 × 1 double | -1.1 | 1.09 | No | numerictype(1, 16, 14) |
| reset | logical | 1 | 1 | Yes | numerictype(0, 1, 0) |
| Output | | | | | |
| y | 256 × 1 double | -1 | 1.04 | No | numerictype(1, 16, 14) |
| Persistent | | | | | |
| z | 1 × 12 double | -1 | 1 | No | numerictype(1, 16, 14) |
| p | double | 0 | 4 | Yes | numerictype(0, 3, 0) |

4 To convert the floating-point algorithm to fixed point, click **Convert**.

The software validates the proposed types and generates a fixed-point version of the entry-point function.

If errors and warnings occur during validation, the app displays them on the **Output** tab. See “Validating Types” on page 7-21.

Test Numerics and Check for Overflows

1 Click the **Test** arrow . Verify that the test file is `overflow_test.m`. Select **Use scaled doubles to detect overflows**, and then click **Test**.

The app runs the test file that you used to define input types to test the fixed-point MATLAB code. Because you selected to detect overflows, it also runs the simulation using scaled double versions

of the proposed fixed-point types. Scaled doubles store their data in double-precision floating-point, so they carry out arithmetic in full range. Because they retain their fixed-point settings, they can report when a computation goes out of the range of the fixed-point type.

The simulation runs. The app detects an overflow. The app reports the overflow on the **Overflow** tab. To highlight the expression that overflowed, click the overflow.

The screenshot shows the 'Convert to Fixed Point' application window. The main editor displays a MATLAB script named 'overflow_test.m'. The script contains several lines of code, with line 39 highlighted in red: `acc(:) = acc + b(j)*z(k);`. Below the code editor, there is a table with tabs for 'Variables', 'Function Replacements', 'Output', 'Errors', 'Verification Output', and 'Overflows'. The 'Overflows' tab is active, showing a table with one entry:

| Function | Line | Description |
|------------|------|---|
| fir_filter | 39 | Overflow error in expression 'acc + b(j)*z(k)'. Percentage of Current Range = 104%. |

- Determine whether it was the sum or the multiplication that overflowed.

In the **fimath** settings, set **Product mode** to `FullPrecision`, and then repeat the conversion and test the fixed-point code again.

The overflow still occurs, indicating that it is the addition in the expression that is overflowing.

Propose Data Types Based on Simulation Ranges

This example shows how to propose fixed-point data types based on simulation range data using the Fixed-Point Converter app.

Prerequisites

This example requires the following products:

- MATLAB
- Fixed-Point Designer
- C compiler

See https://www.mathworks.com/support/compilers/current_release/.

You can use `mex -setup` to change the default compiler. See “Change Default Compiler”.

Create a New Folder and Copy Relevant Files

- 1 In a local, writable folder, create a function `ex_2ndOrder_filter.m`.

```
function y = ex_2ndOrder_filter(x) %#codegen
    persistent z
    if isempty(z)
        z = zeros(2,1);
    end
    % [b,a] = butter(2, 0.25)
    b = [0.0976310729378175, 0.195262145875635, 0.0976310729378175];
    a = [1, -0.942809041582063, 0.333333333333333];

    y = zeros(size(x));
    for i = 1:length(x)
        y(i) = b(1)*x(i) + z(1);
        z(1) = b(2)*x(i) + z(2) - a(2) * y(i);
        z(2) = b(3)*x(i)          - a(3) * y(i);
    end
end
```

- 2 Create a test file, `ex_2ndOrder_filter_test.m`, to exercise the `ex_2ndOrder_filter` algorithm.

It is a best practice is to create a separate test script to do pre- and post-processing, such as:

- Loading inputs.
- Setting up input values.
- Outputting test results.

See “Create a Test File” on page 11-3.

To cover the full intended operating range of the system, the test script runs the `ex_2ndOrder_filter` function with three input signals: chirp, step, and impulse. The script then plots the outputs.

```
% ex_2ndOrder_filter_test
%
```

```

% Define representative inputs
N = 256; % Number of points
t = linspace(0,1,N); % Time vector from 0 to 1 second
f1 = N/2; % Target frequency of chirp set to Nyquist
x_chirp = sin(pi*f1*t.^2); % Linear chirp from 0 to Fs/2 Hz in 1 second
x_step = ones(1,N); % Step
x_impulse = zeros(1,N); % Impulse
x_impulse(1) = 1;

% Run the function under test
x = [x_chirp;x_step;x_impulse];
y = zeros(size(x));
for i = 1:size(x,1)
    y(i,:) = ex_2ndOrder_filter(x(i,:));
end

% Plot the results
titles = {'Chirp', 'Step', 'Impulse'}
clf
for i = 1:size(x,1)
    subplot(size(x,1),1,i)
    plot(t,x(i,:),t,y(i,:))
    title(titles{i})
    legend('Input','Output')
end
xlabel('Time (s)')
figure(gcf)

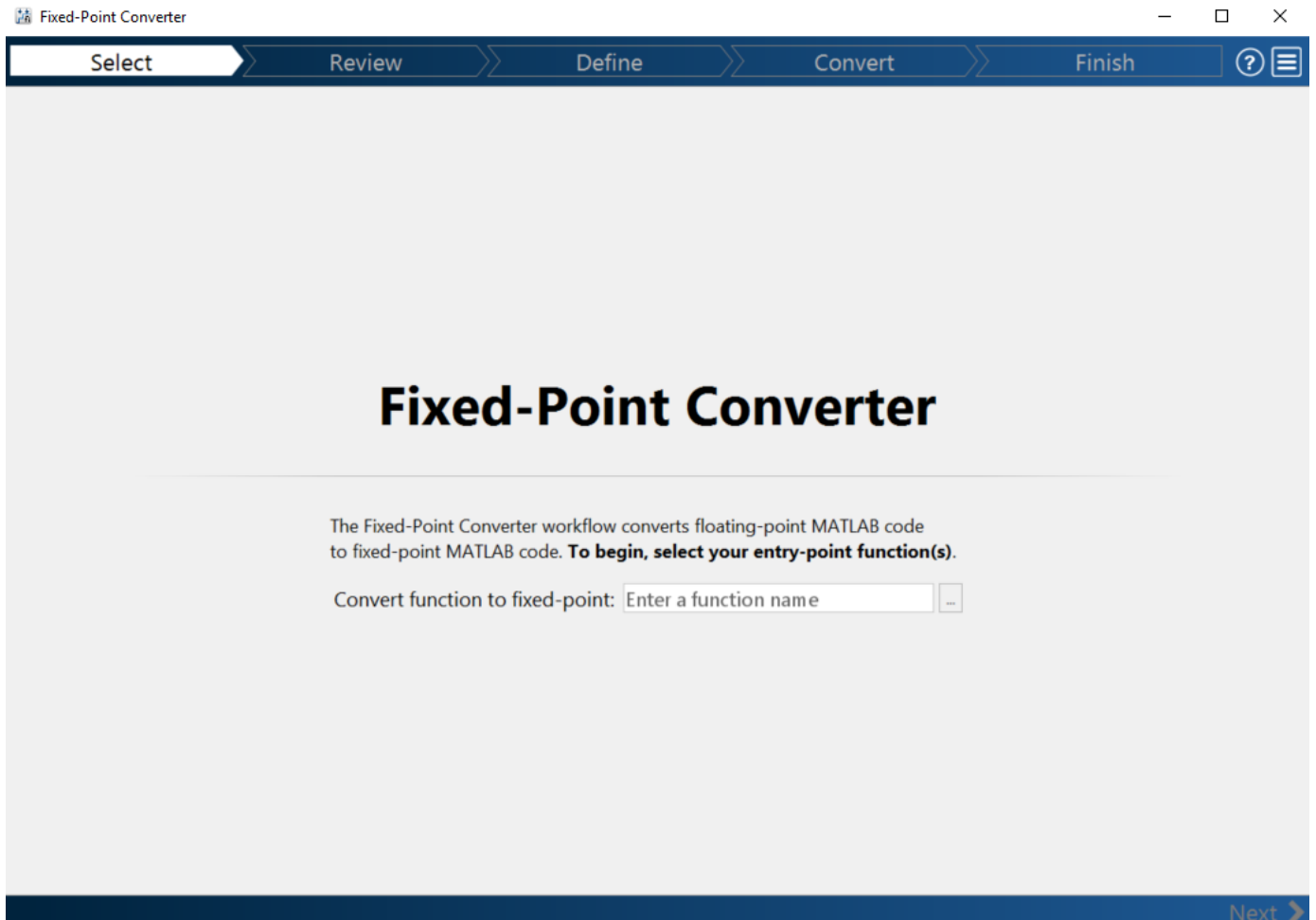
disp('Test complete.')

```

| Type | Name | Description |
|---------------|---------------------------|---|
| Function code | ex_2ndOrder_filter.m | Entry-point MATLAB function |
| Test file | ex_2ndOrder_filter_test.m | MATLAB script that tests ex_2ndOrder_filter.m |

Open the Fixed-Point Converter App

- 1 Navigate to the work folder that contains the file for this example.
- 2 On the MATLAB Toolstrip **Apps** tab, under **Code Generation**, click the app icon.



Select Source Files

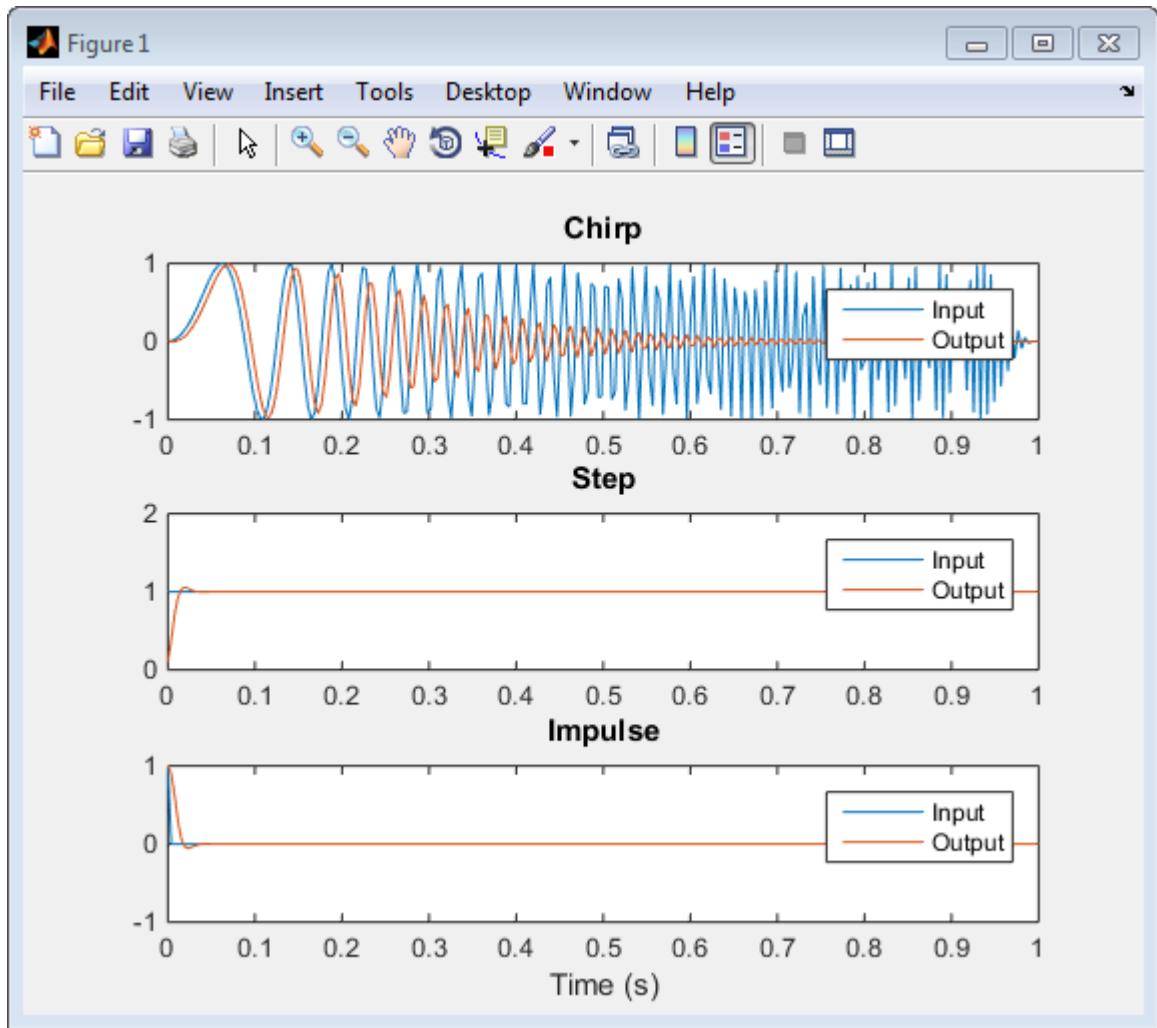
- 1 To add the entry-point function `ex_2ndOrder_filter` to the project, browse to the file `ex_2ndOrder_filter.m`, and then click **Open**. By default, the app saves information and settings for this project in the current folder in a file named `ex_2ndOrder_filter.prj`.
- 2 Click **Next** to go to the **Define Input Types** step.

The app screens `ex_2ndOrder_filter.m` for code violations and fixed-point conversion readiness issues. The app does not find issues in `ex_2ndOrder_filter.m`.

Define Input Types

- 1 On the **Define Input Types** page, to add `ex_2ndOrder_filter_test` as a test file, browse to `ex_2ndOrder_filter_test`, and then click **Open**.
- 2 Click **Autodefine Input Types**.

The test file runs and displays the outputs of the filter for each of the input signals.



The app determines from the test file that the input type of x is `double(1x256)`.

To **automatically define input types**, call `ex_2ndOrder_filter` or enter a script that calls `ex_2ndOrder_filter` in the MATLAB prompt below:

```
>> ex_2ndOrder_filter_test
```

ex_2ndOrder_filter.m

| | |
|---|-----------------|
| x | double(1 x 256) |
|---|-----------------|

[Add global](#)

- 3 Click **Next** to go to the **Convert to Fixed Point** step.

Convert to Fixed Point

- 1 The app generates an instrumented MEX function for your entry-point MATLAB function. The app displays compiled information—type, size, and complexity—for variables in your code. See “View and Modify Variable Information” on page 8-35.

The screenshot shows the Fixed-Point Converter app interface. The top window title is "Fixed-Point Converter - ex_2ndOrder_filter.prj". The main window has a title bar "Convert to Fixed Point" and a menu bar with "SETTINGS", "ANALYZE", "CONVERT", and "TEST". The left sidebar shows "Source Code" and "ex_2ndOrder_filter". The main editor displays the following MATLAB code:

```


1 function y = ex_2ndOrder_filter(x) %#codegen
2     persistent z
3     if isempty(z)
4         z = zeros(2,1);
5     end
6     % [b,a] = butter(2, 0.25)
7     b = [0.0976310729378175, 0.195262145875635, 0.0976310729378175];
8     a = [
9         1, -0.942809041582063, 0.333333333333333];
10
11    y = zeros(size(x));
12    for i=1:length(x)
13        y(i) = b(1)*x(i) + z(1);
14        z(1) = b(2)*x(i) + z(2) - a(2) * y(i);
15        z(2) = b(3)*x(i) - a(3) * y(i);
16    end
17 end
  
```

Below the code editor is a table with tabs for "Variables", "Function Replacements", and "Output". The "Variables" tab is active, showing the following information:

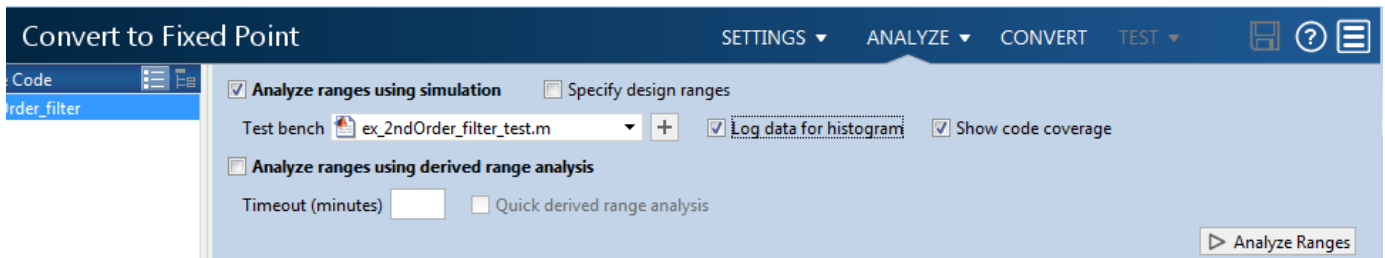
| Variable | Type | Sim Min | Sim Max | Whole Number | Proposed Type |
|------------|----------------|---------|---------|--------------|---------------|
| Input | | | | | |
| x | 1 x 256 double | | | No | |
| Output | | | | | |
| y | 1 x 256 double | | | No | |
| Persistent | | | | | |
| z | 2 x 1 double | | | No | |
| Local | | | | | |

At the bottom of the window, there are "Back" and "Next" navigation buttons.

On the **Function Replacements** tab, the app displays functions that are not supported for fixed-point conversion. See “Running a Simulation” on page 7-8.

- 2 Click the **Analyze** arrow . Verify that **Analyze ranges using simulation** is selected and that the test bench file is `ex_2ndOrder_filter_test`. You can add test files and select to run more than one test file during the simulation. If you run multiple test files, the app merges the simulation results.
- 3 Select **Log data for histogram**.

By default, the **Show code coverage** option is selected. This option provides code coverage information that helps you verify that your test file is testing your algorithm over the intended operating range.



4 Click **Analyze**.

The simulation runs and the app displays a color-coded code coverage bar to the left of the MATLAB code. Review this information to verify that the test file is testing the algorithm adequately. The dark green line to the left of the code indicates that the code runs every time the algorithm executes. The orange bar indicates that the code next to it executes only once. This behavior is expected for this example because the code initializes a persistent variable. If your test file does not cover all of your code, update the test or add more test files.

```

1 function y = ex_2ndOrder_filter(x) %#codegen
2     persistent z
3     if isempty(z)
4         z = zeros(2,1);
5     end
6     % [b,a] = butter(2, 0.25)
7     b = [0.0976310729378175, 0.195262145875635, 0.0976310729378175];
8     a = [1, -0.942809041582063, 0.333333333333333];
9
10
11    y = zeros(size(x));
12    for i=1:length(x)
13        y(i) = b(1)*x(i) + z(1);
14        z(1) = b(2)*x(i) + z(2) - a(2) * y(i);
15        z(2) = b(3)*x(i) - a(3) * y(i);
16    end
17 end
  
```

| Variable | Type | Sim Min | Sim Max | Static ... | Static ... | Whole ... | Proposed Type | Lo... | Max Diff |
|-------------------|------------|---------|---------|------------|------------|-----------|--------------------|-------|----------|
| Input | | | | | | | | | |
| x | 1 x 25... | -1 | 1 | | | No | numeric(1, 16, 14) | | |
| Output | | | | | | | | | |
| y | 1 x 25... | -0.97 | 1.06 | | | No | numeric(1, 16, 14) | | |
| Persistent | | | | | | | | | |
| z | 2 x 1 d... | -0.89 | 0.96 | | | No | numeric(1, 16, 15) | | |
| Local | | | | | | | | | |

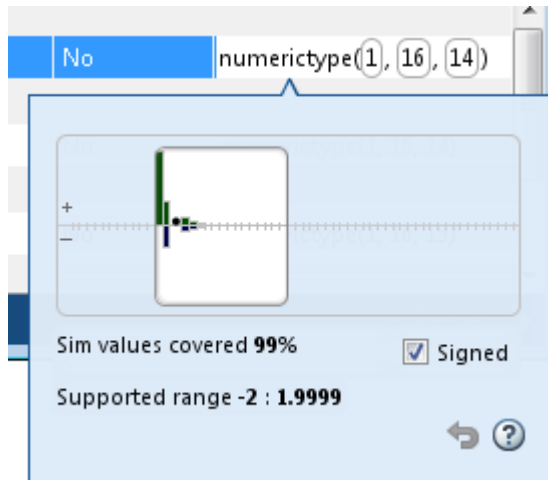
If a value has . . . next to it, the value is rounded. Pause over the . . . to view the actual value.

The app displays simulation minimum and maximum ranges on the **Variables** tab. Using the simulation range data, the software proposes fixed-point types for each variable based on the

default type proposal settings, and displays them in the **Proposed Type** column. The app enables the **Convert** option.

Note You can manually enter static ranges. These manually entered ranges take precedence over simulation ranges. The app uses the manually entered ranges to propose data types. You can also modify and lock the proposed type.

- 5 Examine the proposed types and verify that they cover the full simulation range. To view logged histogram data for a variable, click its **Proposed Type** field.



To modify the proposed data types, either enter the required type into the **Proposed Type** field or use the histogram controls. For more information about the histogram, see “Log Data for Histogram” on page 7-19.

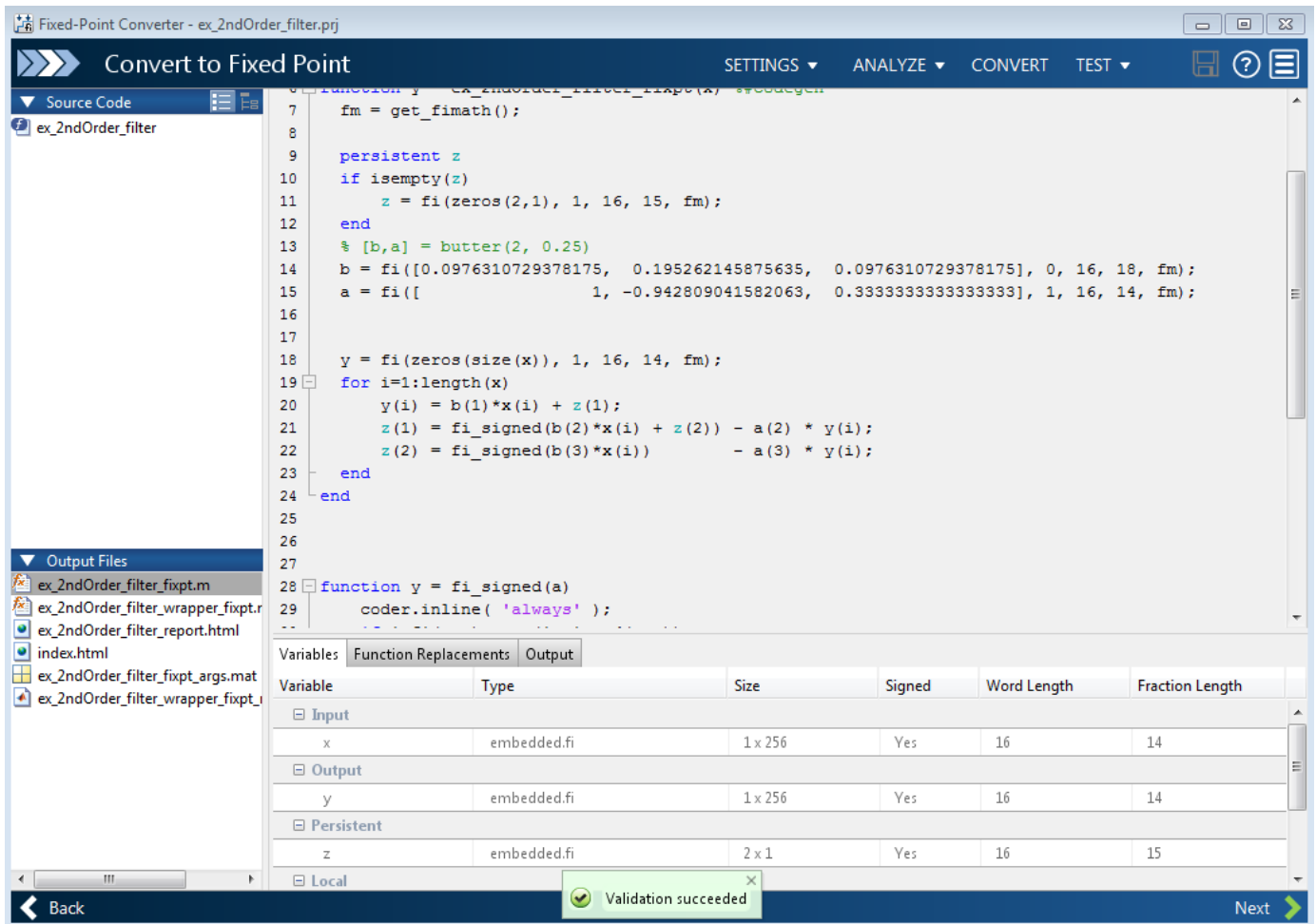
- 6 To convert the floating-point algorithm to fixed point, click **Convert**.


During the fixed-point conversion process, the software validates the proposed types and generates the following files in the `codegen\ex_2ndOrder_filter\fixpt` folder in your local working folder.

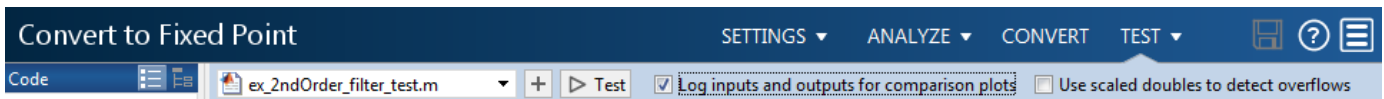
- `ex_2ndOrder_filter_fixpt.m` — the fixed-point version of `ex_2ndOrder_filter.m`.
- `ex_2ndOrder_filter_wrapper_fixpt.m` — this file converts the floating-point data values supplied by the test file to the fixed-point types determined for the inputs during conversion. These fixed-point values are fed into the converted fixed-point design, `ex_2ndOrder_filter_fixpt.m`.
- `ex_2ndOrder_filter_fixpt_report.html` — this report shows the generated fixed-point code and the fixed-point instrumentation results.
- `ex_2ndOrder_filter_report.html` — this report shows the original algorithm and the fixed-point instrumentation results.
- `ex_2ndOrder_filter_fixpt_args.mat` — MAT-file containing a structure for the input arguments, a structure for the output arguments and the name of the fixed-point file.

If errors or warnings occur during validation, you see them on the **Output** tab. See “Validating Types” on page 7-21.

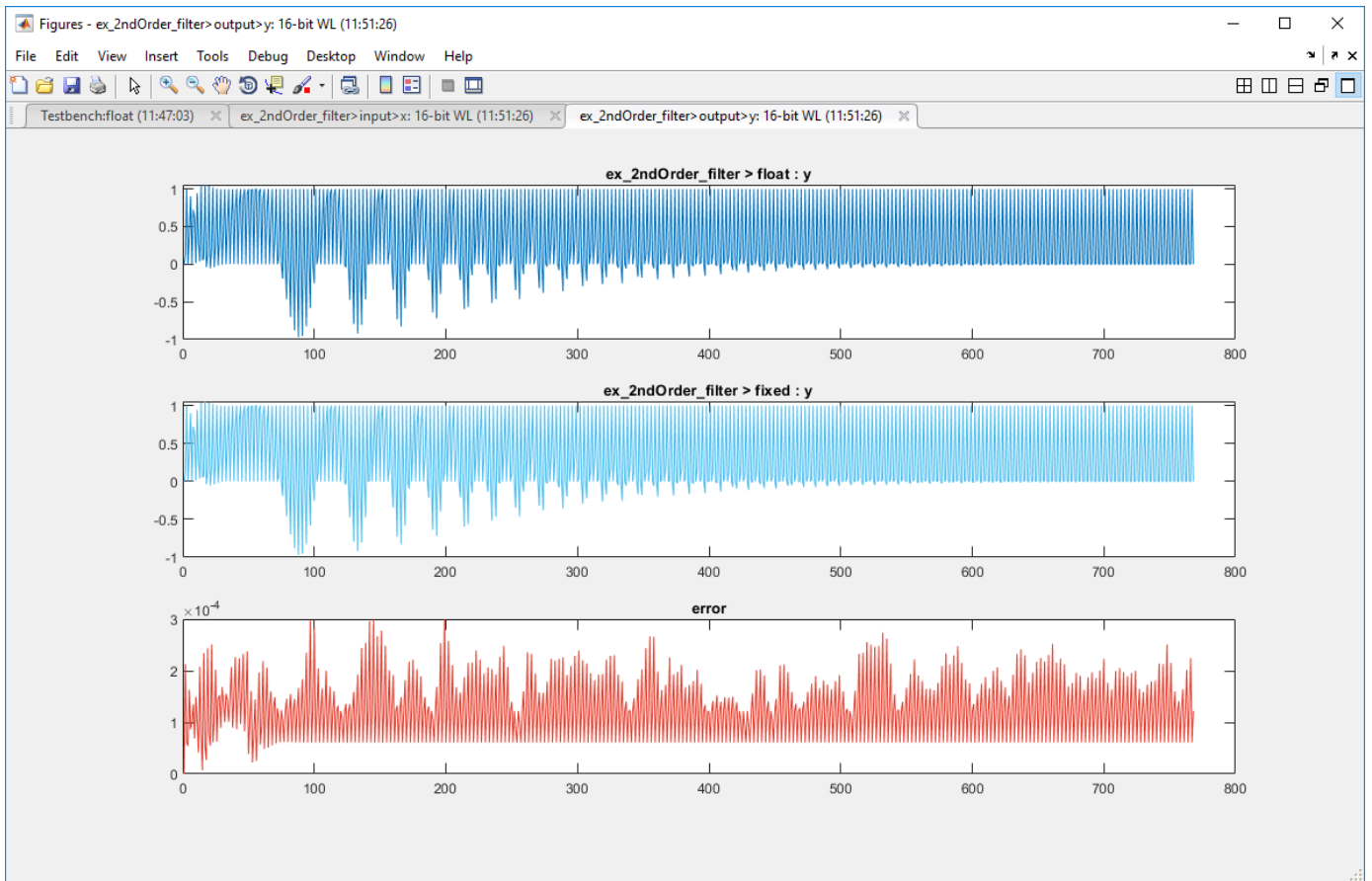
- 7 In the **Output Files** list, select `ex_2ndOrder_filter_fixpt.m`. The app displays the generated fixed-point code.



- 8 Click the **Test** arrow . Select **Log inputs and outputs for comparison plots**, and then click **Test**.



To test the fixed-point MATLAB code, the app runs the test file that you used to define input types. Optionally, you can add test files and select to run more than one test file to test numerics. The software runs both a floating-point and a fixed-point simulation and then calculates the errors for the output variable *y*. Because you selected to log inputs and outputs for comparison plots, the app generates a plot for each input and output. The app docks these plots in a single figure window.



The app also reports error information on the **Verification Output** tab. The maximum error is less than 0.03%. For this example, this margin of error is acceptable.

If the difference is not acceptable, modify the fixed-point data types or your original algorithm. For more information, see “Testing Numerics” on page 7-22.

- 9 On the **Verification Output** tab, the app provides a link to a report that shows the generated fixed-point code and the proposed type information.

Fixed-Point Report *ex_2ndOrder_filter_fixpt*

```
function y = ex_2ndOrder_filter_fixpt(x) %#codegen
    fm = get_fimath();

    persistent z
    if isempty(z)
        z = fi(zeros(2,1), 1, 16, 15, fm);
    end
    % [b,a] = butter(2, 0.25)
    b = fi([0.0976310729378175, 0.195262145875635, 0.0976310729378175], 0, 16, 18, fm);
    a = fi([
        1, -0.942809041582063, 0.333333333333333], 1, 16, 14, fm);

    y = fi(zeros(size(x)), 1, 16, 14, fm);
    for i=1:length(x)
        y(i) = b(1)*x(i) + z(1);
        z(1) = fi_signed(b(2)*x(i) + z(2)) - a(2) * y(i);
        z(2) = fi_signed(b(3)*x(i)) - a(3) * y(i);
    end
end
```

| Variable Name | Type | Sim Min | Sim Max |
|---------------|--------------------------------|---------------------|---------------------|
| a | numerictype(1, 16, 14) 1 x 3 | -0.94281005859375 | 1 |
| b | numerictype(0, 16, 18) 1 x 3 | 0.09762954711914063 | 0.19525909423828125 |
| i | double | 1 | 256 |
| x | numerictype(1, 16, 14) 1 x 256 | -1 | 1 |
| y | numerictype(1, 16, 14) 1 x 256 | -0.9698486328125 | 1.0552978515625 |
| z | numerictype(1, 16, 15) 2 x 1 | -0.890869140625 | 0.957672119140625 |

10 Click **Next** to go to the **Finish Workflow** page.

On the **Finish Workflow** page, the app displays a project summary and links to generated output files.

Integrate Fixed-Point Code

To integrate the fixed-point version of the code into system-level simulations, generate a MEX function to accelerate the fixed-point algorithm. Call this MEX function instead of the original MATLAB algorithm.

- 1 Copy `ex_2ndOrder_filter_fixpt.m` to your local working folder.
- 2 Generate a MEX function for `ex_2ndOrder_filter_fixpt.m`. Look at the `get_fimath` function in the `ex_2ndOrder_filter_fixpt.m` file to get the `fimath`, and use the type proposal report to get fixed-point data type for input `x`.

```
fm = fimath('RoundingMethod','Floor',...
    'OverflowAction','Wrap',...
    'ProductMode','FullPrecision',...
    'MaxProductWordLength',128,...
    'SumMode','FullPrecision',...
    'MaxSumWordLength',128);
fiaccel ex_2ndOrder_filter_fixpt -args {fi(0,1,16,14,fm)}
```

`fiaccel` generates a MEX function, `ex_2ndOrder_filter_fixpt_mex`, in the current folder.

- 3** You can now call this MEX function in place of the original MATLAB algorithm.

Propose Data Types Based on Derived Ranges

This example shows how to propose fixed-point data types based on static ranges using the Fixed-Point Converter app. When you propose data types based on derived ranges you, do not have to provide test files that exercise your algorithm over its full operating range. Running such test files often takes a long time. You can save time by deriving ranges instead.

Note Derived range analysis is not supported for non-scalar variables.

Prerequisites

This example requires the following products:

- MATLAB
- Fixed-Point Designer
- C compiler

See https://www.mathworks.com/support/compilers/current_release/.

You can use `mex -setup` to change the default compiler. See “Change Default Compiler”.

Create a New Folder and Copy Relevant Files

- 1 In a local, writable folder, create a function `dti.m`.

The `dti` function implements a Discrete Time Integrator in MATLAB.

```
function [y, clip_status] = dti(u_in) %#codegen
% Discrete Time Integrator in MATLAB
%
% Forward Euler method, also known as Forward
% Rectangular, or left-hand approximation.
% The resulting expression for the output of
% the block at step 'n' is
% y(n) = y(n-1) + K * u(n-1)
%
init_val = 1;
gain_val = 1;
limit_upper = 500;
limit_lower = -500;

% Variable to hold state between
% consecutive calls to this block
persistent u_state;
if isempty(u_state)
    u_state = init_val+1;
end

% Compute Output
if (u_state > limit_upper)
    y = limit_upper;
    clip_status = -2;
elseif (u_state >= limit_upper)
```



```

        y = limit_upper;
        clip_status = -1;
elseif (u_state < limit_lower)
    y = limit_lower;
    clip_status = 2;
elseif (u_state <= limit_lower)
    y = limit_lower;
    clip_status = 1;
else
    y = u_state;
    clip_status = 0;
end

```

```

% Update State
tprod = gain_val * u_in;
u_state = y + tprod;

```

- 2 Create a test file, `dti_test.m`, to exercise the `dti` algorithm.

The test script runs the `dti` function with a sine wave input. The script then plots the input and output signals.

```

% dti_test
% cleanup
clear dti

% input signal
x_in = sin(2.*pi.*(0:0.001:2)).';

pause(10);

len = length(x_in);
y_out = zeros(1,len);
is_clipped_out = zeros(1,len);

for ii=1:len
    data = x_in(ii);
    % call to the dti function
    init_val = 0;
    gain_val = 1;
    upper_limit = 500;
    lower_limit = -500;

    % call to the design that does DTI
    [y_out(ii), is_clipped_out(ii)] = dti(data);
end

figure('Name', [mfilename, '_plot']);
subplot(2,1,1)
plot(1:len,x_in)
xlabel('Time')
ylabel('Amplitude')
title('Input Signal (Sin)')

subplot(2,1,2)
plot(1:len,y_out)
xlabel('Time')

```

```
ylabel('Amplitude')
title('Output Signal (DTI)')

disp('Test complete.');
```

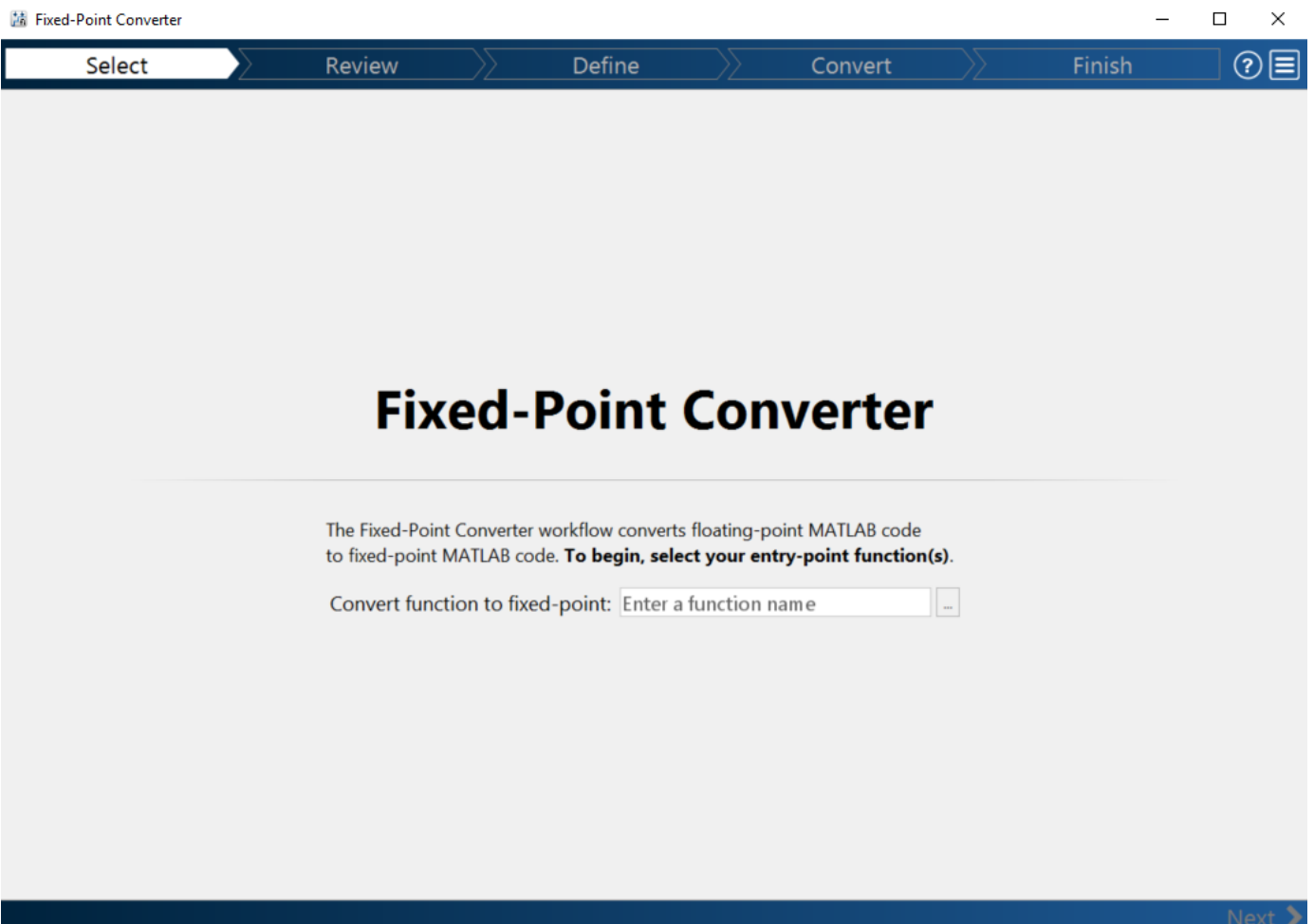
It is a best practice is to create a separate test script to do pre- and post-processing, such as:

- Loading inputs.
- Setting up input values.
- Outputting test results.

| Type | Name | Description |
|---------------|------------|--------------------------------|
| Function code | dti.m | Entry-point MATLAB function |
| Test file | dti_test.m | MATLAB script that tests dti.m |

Open the Fixed-Point Converter App

- 1 Navigate to the work folder that contains the file for this example.
- 2 On the MATLAB Toolstrip **Apps** tab, under **Code Generation**, click the app icon.



Select Source Files

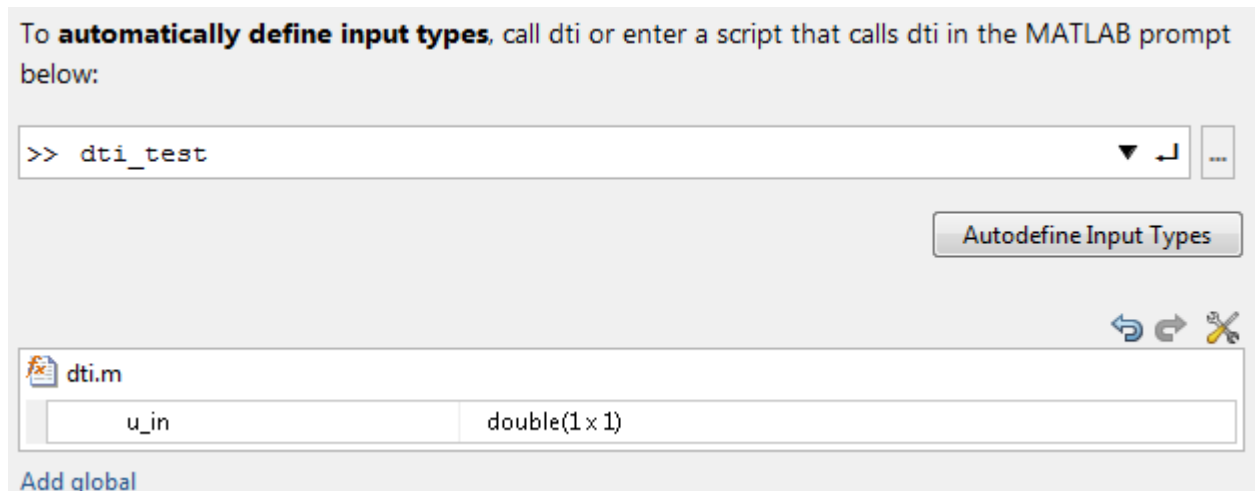
- 1 To add the entry-point function `dti` to the project, browse to the file `dti.m`, and then click **Open**. By default, the app saves information and settings for this project in the current folder in a file named `dti.prj`.
- 2 Click **Next** to go to the **Define Input Types** step.

The app screens `dti.m` for code violations and fixed-point conversion readiness issues. The app does not find issues in `dti.m`.

Define Input Types

- 1 On the **Define Input Types** page, to add `dti_test` as a test file, browse to `dti_test.m`, and then click **Open**.
- 2 Click **Autodefine Input Types**.

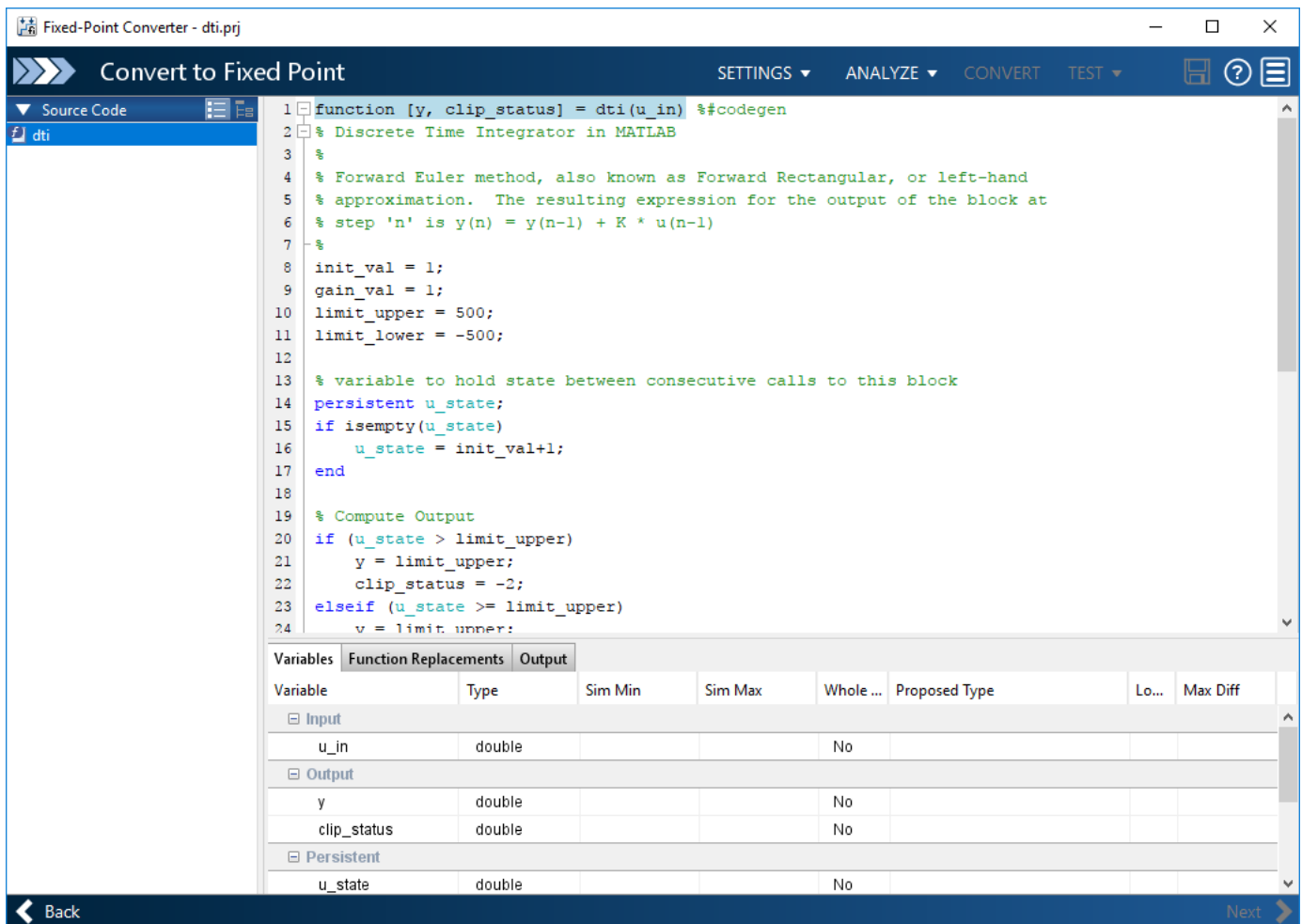
The test file runs. The app determines from the test file that the input type of `u_in` is `double(1x1)`.



- 3 Click **Next** to go to the **Convert to Fixed Point** step.

Convert to Fixed Point

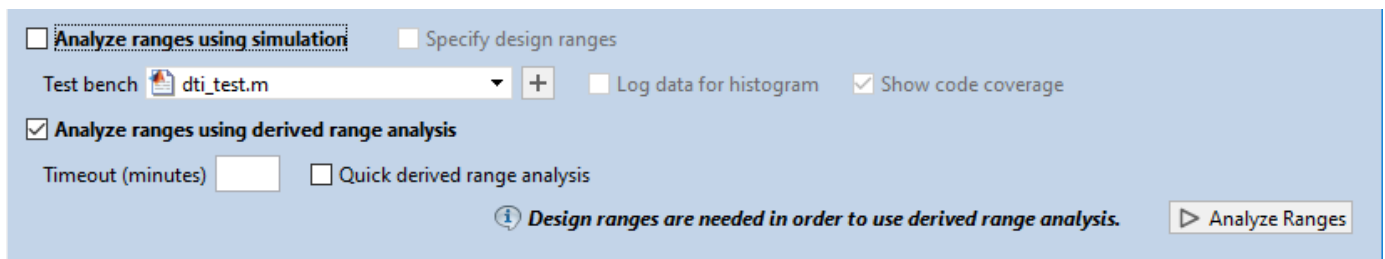
- 1 The app generates an instrumented MEX function for your entry-point MATLAB function. The app displays compiled information—type, size, and complexity—for variables in your code. For more information, see “View and Modify Variable Information” on page 8-35.



If functions are not supported for fixed-point conversion, the app displays them on the **Function Replacements** tab.

- 2 Click the **Analyze** arrow .
 - a Select **Analyze ranges using derived range analysis**.
 - b Clear the **Analyze ranges using simulation** check box.

Design ranges are required to use derived range analysis.



- 3 On the **Convert to Fixed Point** page, on the **Variables** tab, for input `u_in`, select **Static Min** and set it to -1. Set **Static Max** to 1.

To compute derived range information, at a minimum you must specify static minimum and maximum values or proposed data types for all input variables.

Note If you manually enter static ranges, these manually entered ranges take precedence over simulation ranges. The app uses the manually entered ranges to propose data types. You can also modify and lock the proposed type.

4 Click **Analyze**.

Range analysis computes the derived ranges and displays them in the **Variables** tab. Using these derived ranges, the analysis proposes fixed-point types for each variable based on the default type proposal settings. The app displays them in the **Proposed Type** column.

In the `dti` function, the `clip_status` output has a minimum value of -2 and a maximum of 2.

```
% Compute Output
if (u_state > limit_upper)
    y = limit_upper;
    clip_status = -2;
elseif (u_state >= limit_upper)
    y = limit_upper;
    clip_status = -1;
elseif (u_state < limit_lower)
    y = limit_lower;
    clip_status = 2;
elseif (u_state <= limit_lower)
    y = limit_lower;
    clip_status = 1;
else
    y = u_state;
    clip_status = 0;
end
```

When you derive ranges, the app analyzes the function and computes these minimum and maximum values for `clip_status`.

```

1 function [y, clip_status] = dti(u_in) %#codegen
2 % Discrete Time Integrator in MATLAB
3 %
4 % Forward Euler method, also known as Forward Rectangular, or left-hand
5 % approximation. The resulting expression for the output of the block at
6 % step 'n' is  $y(n) = y(n-1) + K * u(n-1)$ 
7 %
8 init_val = 1;
9 gain_val = 1;
10 limit_upper = 500;
11 limit_lower = -500;
12
13 % variable to hold state between consecutive calls to this block
14 persistent u_state;
15 if isempty(u_state)
16     u_state = init_val+1;

```

| Variables | | Function Replacements | | Output | | | |
|----------------|--------|-----------------------|---------|------------|------------|------------|------------------------|
| Variable | Type | Sim Min | Sim Max | Static Min | Static Max | Whole N... | Proposed Type |
| [-] Input | | | | | | | |
| u_in | double | | | -1 | 1 | No | numerictype(1, 16, 14) |
| [-] Output | | | | | | | |
| y | double | | | -500 | 500 | No | numerictype(1, 16, 6) |
| clip_status | double | | | -2 | 2 | No | numerictype(1, 16, 13) |
| [-] Persistent | | | | | | | |
| u_state | double | | | -501 | 501 | No | numerictype(1, 16, 6) |
| [-] Local | | | | | | | |
| init_val | double | | | 1 | 1 | Yes | numerictype(0, 1, 0) |
| gain_val | double | | | 1 | 1 | Yes | numerictype(0, 1, 0) |
| limit_upper | double | | | 500 | 500 | Yes | numerictype(0, 9, 0) |
| limit_lower | double | | | -500 | -500 | Yes | numerictype(1, 10, 0) |
| tprod | double | | | -1 | 1 | No | numerictype(1, 16, 14) |

The app provides a **Quick derived range analysis** option and the option to specify a timeout in case the analysis takes a long time. See “Computing Derived Ranges” on page 7-9.


- To convert the floating-point algorithm to fixed point, click **Convert**.

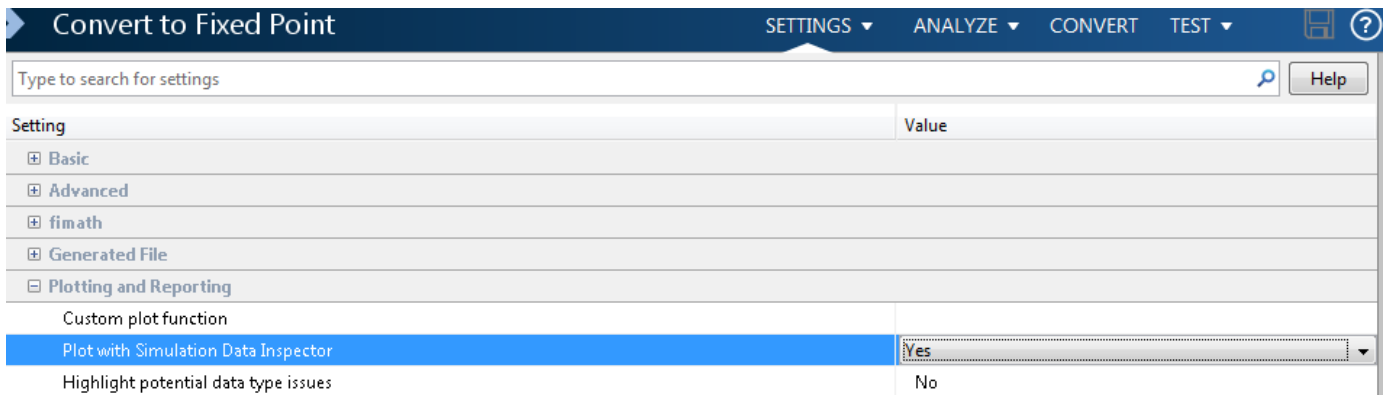
During the fixed-point conversion process, the software validates the proposed types and generates the following files in the `codegen\dti\fixpt` folder in your local working folder:

- `dti_fixpt.m` — the fixed-point version of `dti.m`.
- `dti_wrapper_fixpt.m` — this file converts the floating-point data values supplied by the test file to the fixed-point types determined for the inputs during conversion. The app feeds these fixed-point values into the converted fixed-point design, `dti_fixpt.m`.
- `dti_fixpt_report.html` — this report shows the generated fixed-point code and the fixed-point instrumentation results.
- `dti_report.html` — this report shows the original algorithm and the fixed-point instrumentation results.

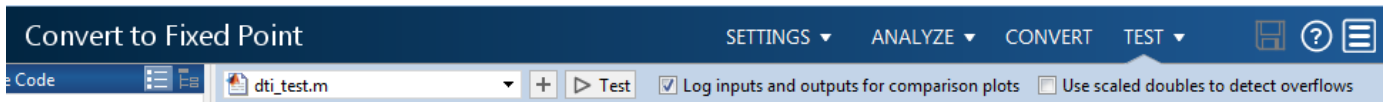
- `dti_fixpt_args.mat` — MAT-file containing a structure for the input arguments, a structure for the output arguments and the name of the fixed-point file.

If errors or warnings occur during validation, they show on the **Output** tab. See “Validating Types” on page 7-21.

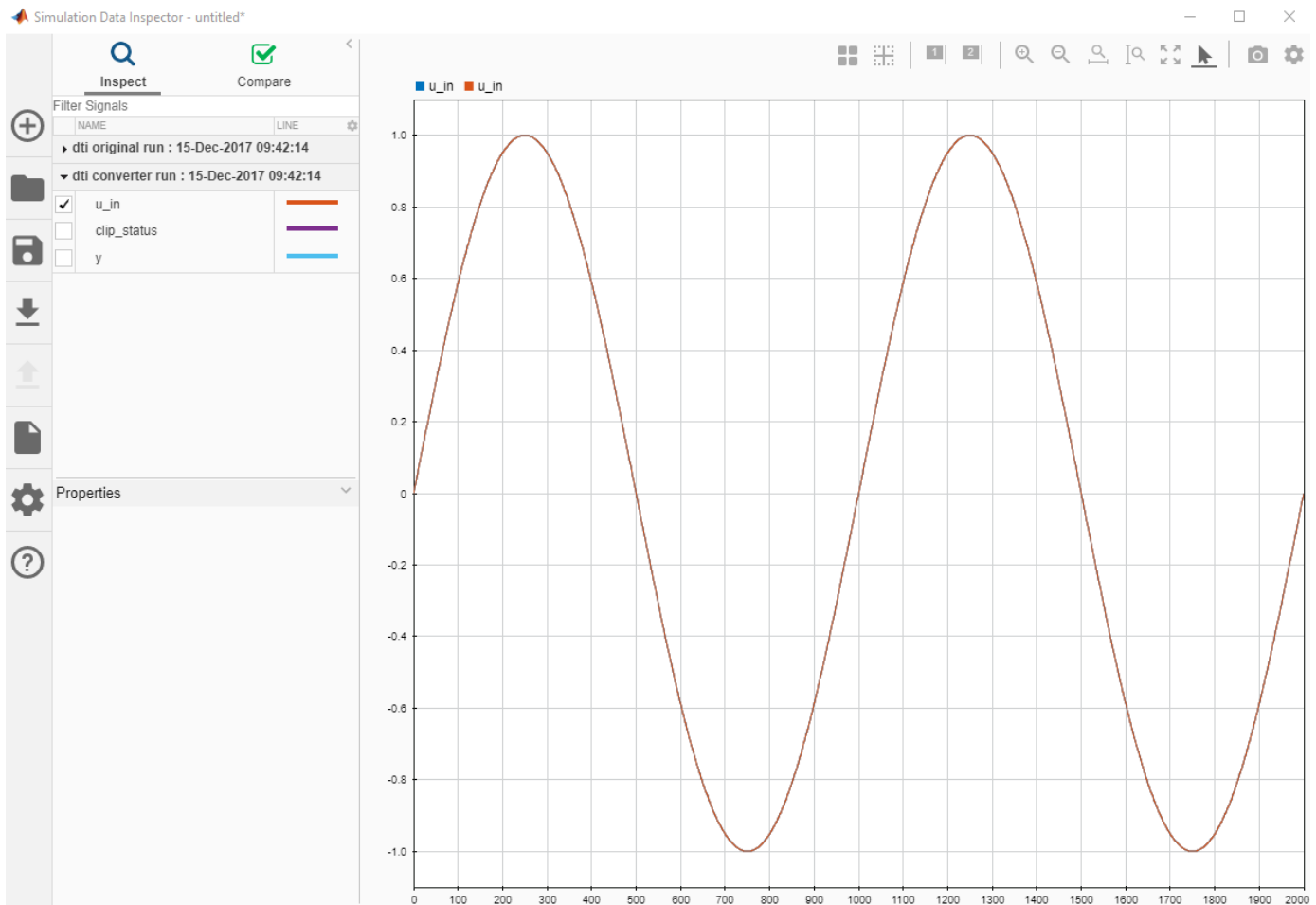
- 6 In the **Output Files** list, select `dti_fixpt.m`. The app displays the generated fixed-point code.
- 7 Use the Simulation Data Inspector to plot the floating-point and fixed-point results.
 - a Click the **Settings** arrow .
 - b Expand the **Plotting and Reporting** settings and set **Plot with Simulation Data Inspector** to Yes.



- c Click the **Test** arrow . Select **Log inputs and outputs for comparison plots**. Click **Test**.

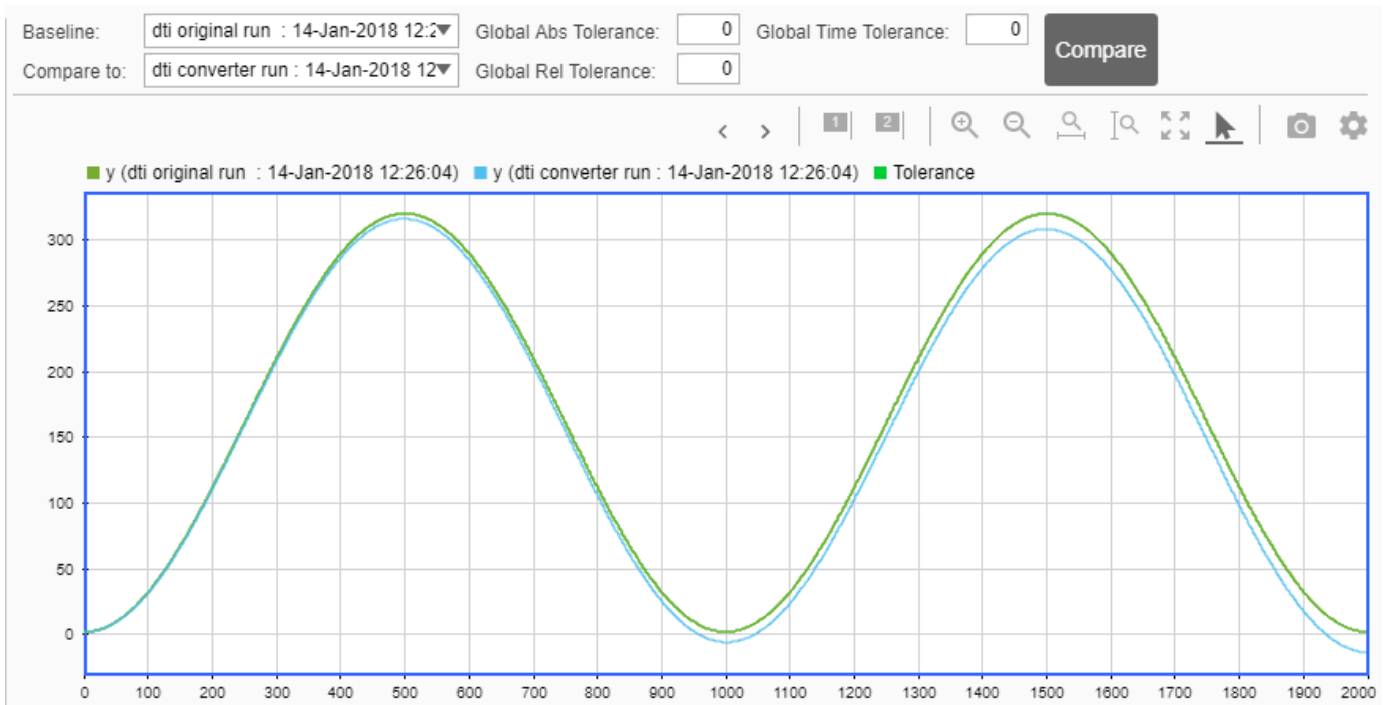


The app runs the test file that you used to define input types to test the fixed-point MATLAB code. Optionally, you can add test files and select to run more than one test file to test numerics. The software runs both a floating-point and a fixed-point simulation and then calculates the errors for the output variable `y`. Because you selected to log inputs and outputs for comparison plots and to use the Simulation Data Inspector for these plots, the Simulation Data Inspector opens.



- d You can use the Simulation Data Inspector to view floating-point and fixed-point run information and compare results. For example, to compare the floating-point and fixed-point values for the output *y*, select *y*. Click **Compare**. Set **Baseline** to the original run and **Compare to** to the converter run. Click **Compare**.

The Simulation Data Inspector displays a plot of the baseline floating-point run against the fixed-point run and the difference between them.



- 8 On the **Verification Output** tab, the app provides a link to the Fixed_Point Report.

```

Variables | Function Replacements | Output | Verification Output
-----
----- Output variable : clip_status -----
Generating comparison plot...

----- Output variable : y -----
Generating comparison plot...

### Generating Fixed-point Types Report for 'dti_fixpt' dti\_fixpt\_report.html
### Elapsed Time:          25.0139 sec(s)

```

To open the report, click the **dti_fixpt_report.html** link.

- 9 Click **Next** to go to the **Finish Workflow** page.

On the **Finish Workflow** page, the app displays a project summary and links to generated output files.

Integrate Fixed-Point Code

To integrate the fixed-point version of the code into system-level simulations, generate a MEX function to accelerate the fixed-point algorithm. Call this MEX function instead of the original MATLAB algorithm.

- 1 Copy `dti_fixpt.m` to your local working folder.
- 2 To get the `fimath` properties for the input argument, look at the `get_fimath` function in `dti_fixpt.m`.

```
function fm = get_fimath()
    fm = fimath('RoundingMethod','Floor',...
        'OverflowAction','Wrap',...
        'ProductMode','FullPrecision',...
        'MaxProductWordLength',128,...
        'SumMode','FullPrecision',...
        'MaxSumWordLength',128);
end
```

- 3 To get the fixed-point data type for input `u_in`, look at the type proposal report.
- 4 Generate a MEX function for `dti_fixpt.m`.

```
fm = fimath('RoundingMethod','Floor',...
    'OverflowAction','Wrap',...
    'ProductMode','FullPrecision',...
    'MaxProductWordLength',128,...
    'SumMode','FullPrecision',...
    'MaxSumWordLength',128);
fiaccel dti_fixpt -args {fi(0,1,16,14,fm)}
```

`fiaccel` generates a MEX function, `dti_fixpt_mex`, in the current folder.

- 5 You can now call this MEX function in place of the original MATLAB algorithm.

View and Modify Variable Information

View Variable Information

On the **Convert to Fixed Point** page of the Fixed-Point Converter app, you can view information about the variables in the MATLAB functions. To view information about the variables for the function that you selected in the **Source Code** pane, use the **Variables** tab or pause over a variable in the code window. For more information, see “Viewing Variables” on page 7-17.

You can view the variable information:

- **Variable**

Variable name. Variables are classified and sorted as inputs, outputs, persistent, or local variables.

- **Type**

The original size, type, and complexity of each variable.

- **Sim Min**

The minimum value assigned to the variable during simulation.

- **Sim Max**

The maximum value assigned to the variable during simulation.

To search for a variable in the MATLAB code window and on the **Variables** tab, use `Ctrl+F`.

Modify Variable Information

If you modify variable information, the app highlights the modified values using bold text. You can modify the following fields:

- **Static Min**

You can enter a value for **Static Min** into the field or promote **Sim Min** information. See “Promote Sim Min and Sim Max Values” on page 8-36.

Editing this field does not trigger static range analysis, but the app uses the edited values in subsequent analyses.

- **Static Max**

You can enter a value for **Static Max** into the field or promote **Sim Max** information. See “Promote Sim Min and Sim Max Values” on page 8-36.

Editing this field does not trigger static range analysis, but the app uses the edited values in subsequent analyses.

- **Whole Number**

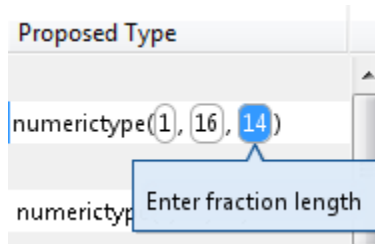
The app uses simulation data to determine whether the values assigned to a variable during simulation were always integers. You can manually override this field.

Editing this field does not trigger static range analysis, but the app uses the edited value in subsequent analyses.

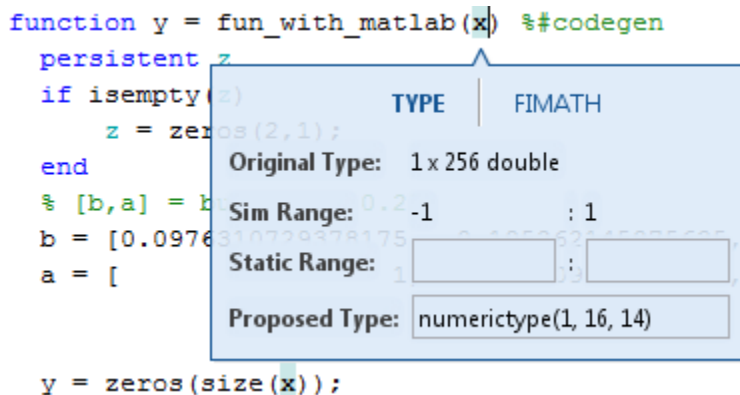
- **Proposed Type**

You can modify the signedness, word length, and fraction length settings individually:

- On the **Variables** tab, modify the value in the **ProposedType** field.



- In the code window, select a variable, and then modify the **Proposed Type** field.



If you selected to log data for a histogram, the histogram dynamically updates to reflect the modifications to the proposed type. You can also modify the proposed type in the histogram, see “Log Data for Histogram” on page 7-19.

Revert Changes

- To clear results and revert edited values, right-click the **Variables** tab and select **Reset entire table**.
- To revert the type of a selected variable to the type computed by the app, right-click the field and select **Undo changes**.
- To revert changes to variables, right-click the field and select **Undo changes for all variables**.
- To clear a static range value, right-click an edited field and select **Clear this static range**.
- To clear manually entered static range values, right-click anywhere on the **Variables** tab and select **Clear all manually entered static ranges**.

Promote Sim Min and Sim Max Values

With the Fixed-Point Converter app, you can promote simulation minimum and maximum values to static minimum and maximum values. This capability is useful if you have not specified static ranges and you have simulated the model with inputs that cover the full intended operating range.

| Simulation Output | Variables | Function Replacements | | | | | |
|-------------------|----------------|-----------------------|---------|------------|------------|--------------|-----------------|
| Variable | Type | Sim Min | Sim Max | Static Min | Static Max | Whole Number | Proposed Type |
| Input | | | | | | | |
| x | 1 x 256 double | -1 | 1 | | | | type(1, 16, 14) |
| Output | | | | | | | |
| y | 1 x 256 double | -0.97 | 1.06 | | | | type(1, 16, 14) |
| Persistent | | | | | | | |

To copy:

- A simulation range for a selected variable, select a variable, right-click, and then select Copy sim range.
- Simulation ranges for top-level inputs, right-click the Static Min or Static Max column, and then select Copy sim ranges for all top-level inputs.
- Simulation ranges for persistent variables, right-click the Static Min or Static Max column, and then select Copy sim ranges for all persistent variables.

Replace the exp Function with a Lookup Table

This example shows how to replace the `exp` function with a lookup table approximation in fixed-point code generated using the Fixed-Point Converter app.

Prerequisites

To complete this example, you must install the following products:

- MATLAB
- Fixed-Point Designer
- C compiler

See https://www.mathworks.com/support/compilers/current_release/.

You can use `mex -setup` to change the default compiler. See “Change Default Compiler”.

Create Algorithm and Test Files

- 1 Create a MATLAB function, `my_fcn.m`, that calls the `exp` function.

```
function y = my_fcn(x)
    y = exp(x);
end
```

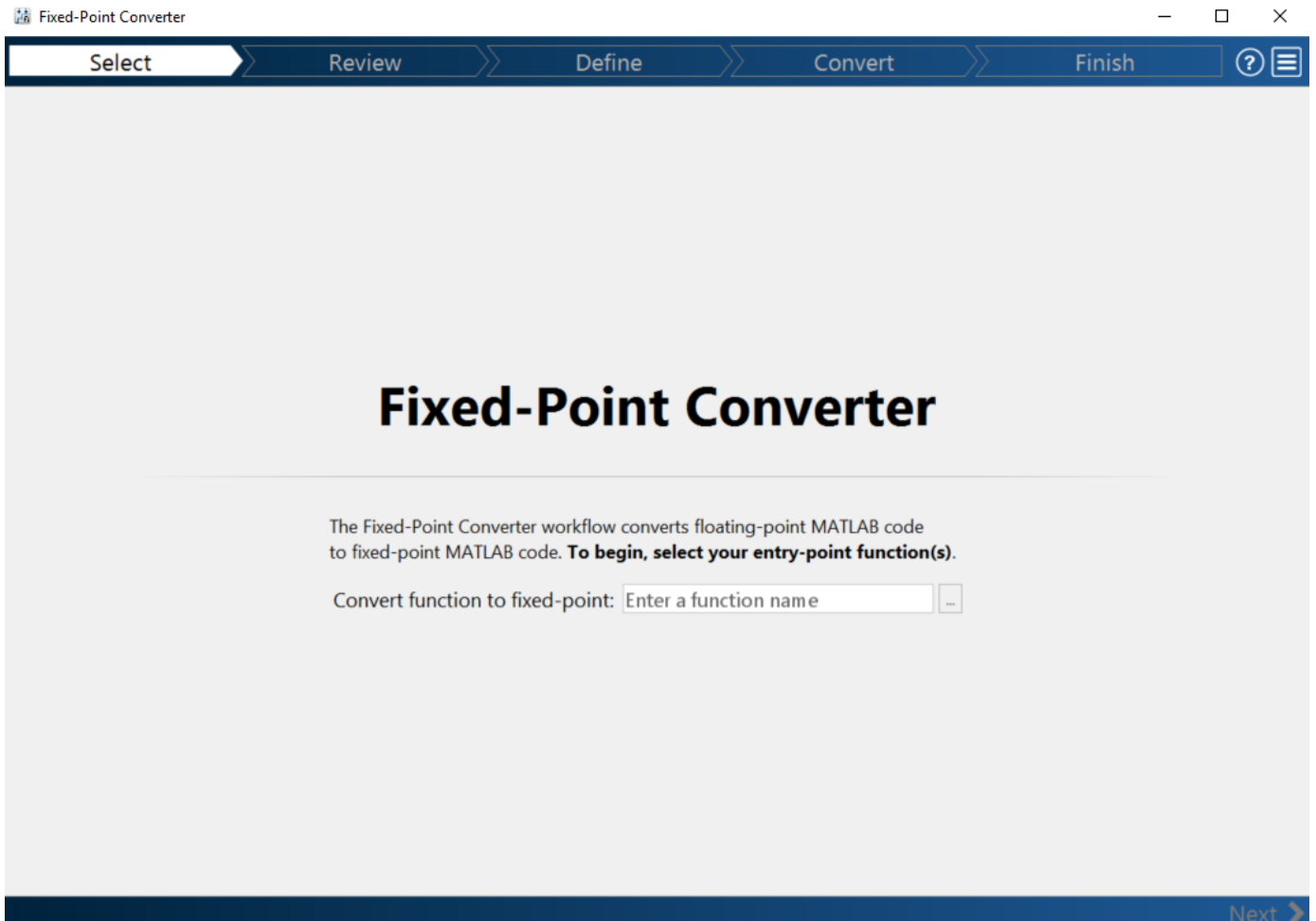
- 2 Create a test file, `my_fcn_test.m`, that uses `my_fcn.m`.

```
close all

x = linspace(-10,10,1e3);
for itr = 1e3:-1:1
    y(itr) = my_fcn( x(itr) );
end
plot( x, y );
```

Open the Fixed-Point Converter App

- 1 Navigate to the work folder that contains the file for this example.
- 2 On the MATLAB Toolstrip **Apps** tab, under **Code Generation**, click the app icon.






Select Source Files

- 1 To add the entry-point function `my_fcn` to the project, browse to the file `my_fcn.m`, and then click **Open**. By default, the app saves information and settings for this project in the current folder in a file named `my_fcn.prj`.
- 2 Click **Next** to go to the **Define Input Types** step.

The app screens `my_fcn.m` for code violations and fixed-point conversion readiness issues. The app opens the **Review Code Generation Readiness** page.


Review Code Generation Readiness

- 1 Click **Review Issues**. The app indicates that the `exp` function is not supported for fixed-point conversion. In a later step, you specify a lookup table replacement for this function.

Review Code Generation Readiness REVIEW ISSUES   

```
Code
1 function y = my_fcn(x)
2     y = exp(x);
3 end
4
5
```

Issues

| | Function | Line | Description |
|--|----------|------|---|
|  | my_fcn.m | 2 | exp is not supported for fixed-point conversion |

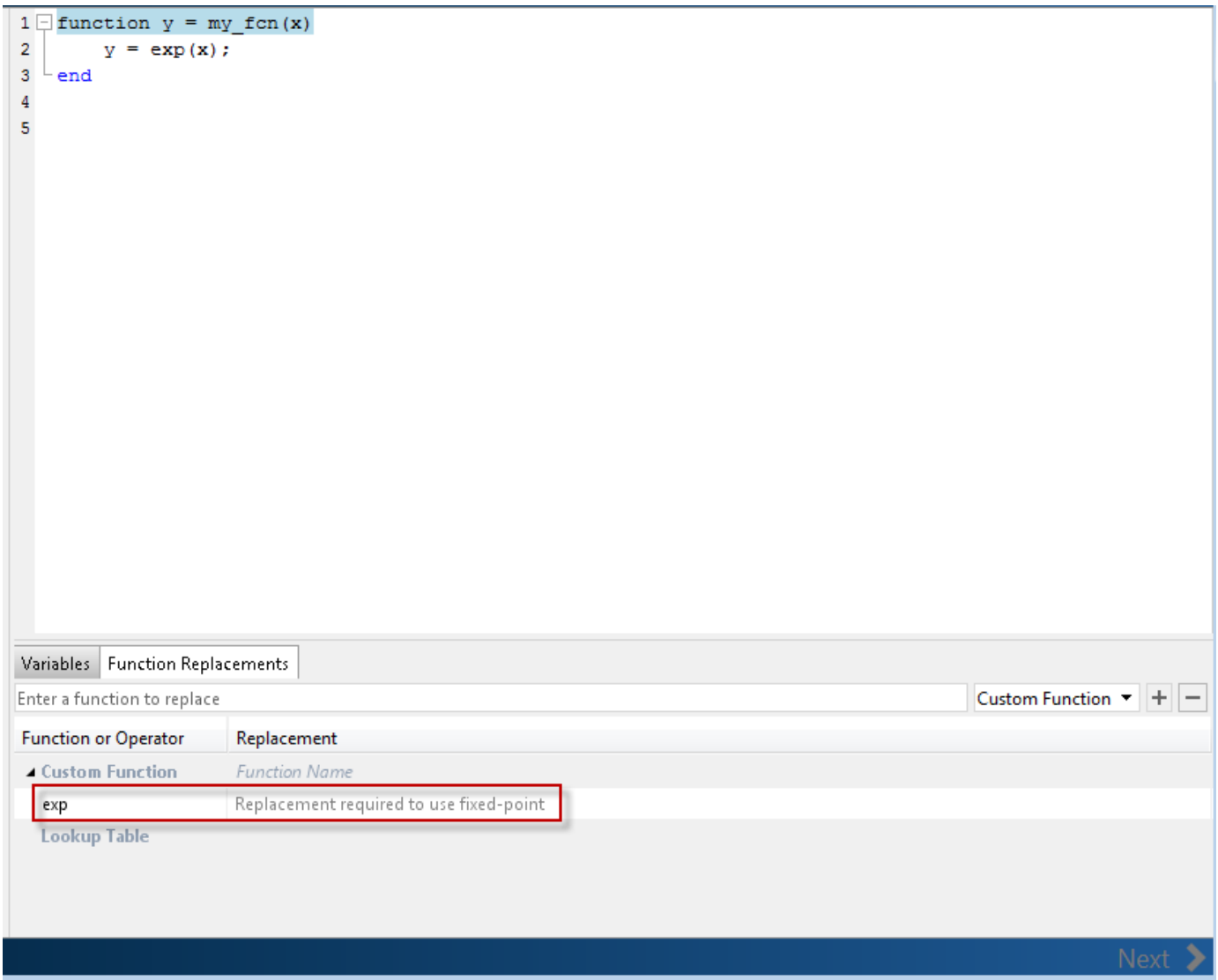
- 2 Click **Next** to go to the **Define Input Types** step.

Define Input Types

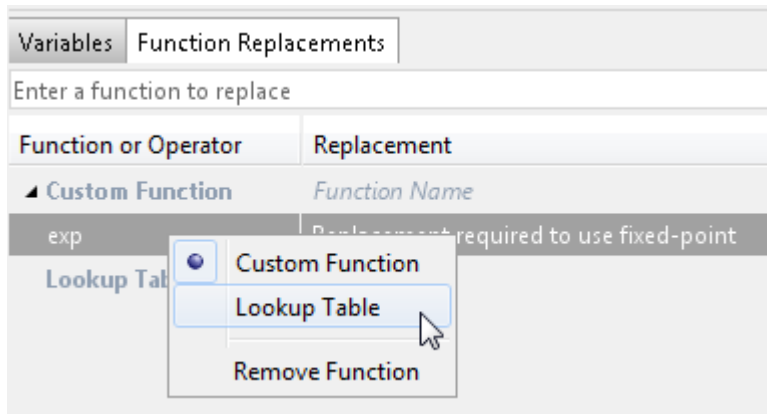
- 1 Add my_fcn_test as a test file and then click **Autodefine Input Types**.
The test file runs. The app determines from the test file that x is a scalar double.
- 2 Click **Next** to go to the **Convert to Fixed Point** step.

Replace exp Function with Lookup Table

- 1 Select the **Function Replacements** tab.
The app indicates that you must replace the exp function.




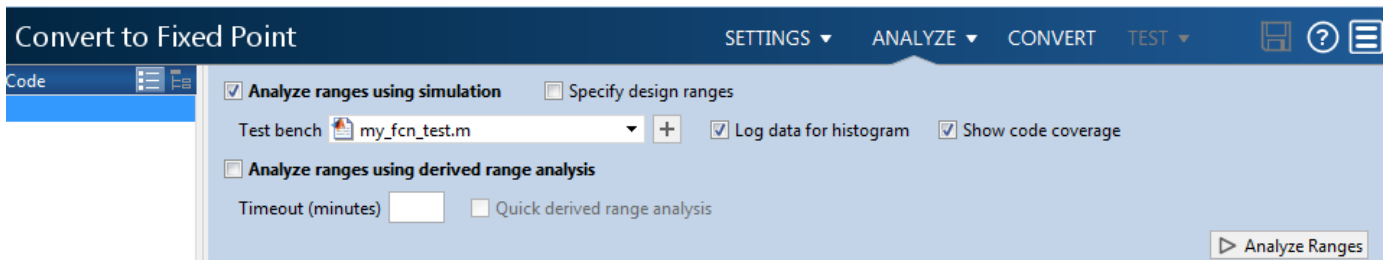
- 2 On the **Function Replacements** tab, right-click the exp function and select Lookup Table.



The app moves the `exp` function to the list of functions that it will replace with a Lookup Table. By default, the lookup table uses linear interpolation and 1000 points. **Design Min** and **Design Max** are set to `Auto` which means that the app uses the design minimum and maximum values that it detects by either running a simulation or computing derived ranges.

| Variables | | Function Replacements | | Output | | |
|---------------------------------------|----------------------|-----------------------|------------|------------------|-----------------|-----|
| Enter a function to replace | | | | | Custom Function | + - |
| Function or Operator | Replacement | | | | | |
| Custom Function | | | | | | |
| <input type="checkbox"/> Lookup Table | Interpolation Method | Design Min | Design Max | Number of Points | | |
| exp | Linear | Auto | Auto | 1000 | | |

- Click the **Analyze** arrow , select **Log data for histogram**, and verify that the test file is `my_fcn_test`.



- Click **Analyze**.

The simulation runs. On the **Variables** tab, the app displays simulation minimum and maximum ranges. Using the simulation range data, the software proposes fixed-point types for each variable based on the default type proposal settings, and displays them in the **Proposed Type** column. The app enables the **Convert** option.

- Examine the proposed types and verify that they cover the full simulation range. To view logged histogram data for a variable, click its **Proposed Type** field. The histogram provides range information and the percentage of simulation range covered by the proposed data type.

| Variable | Type | Sim Min | Sim Max | Whole Number | Proposed Type |
|----------|--------|---------|---------|--------------|------------------------|
| Input | | | | | |
| x | double | -10 | 10 | No | numerictype(1, 16, 11) |
| Output | | | | | |
| y | double | 0 | 22026.4 | | |

Sim values covered **100%** Signed

Supported range **-16 : 15.9995**

Convert to Fixed Point

- 1 Click **Convert**.

The app validates the proposed types, and generates a fixed-point version of the entry-point function, `my_fcn_fixpt.m`.

- 2 In the Output Files list, select `my_fcn_fixpt.m`.

The conversion process generates a lookup table approximation, `replacement_exp`, for the `exp` function.

The screenshot shows the 'Convert to Fixed Point' application interface. The main window displays MATLAB code for a function `my_fcn_fixpt` and its sub-function `replacement_exp`. The code uses `fi` for fixed-point conversion and `LUT` for a lookup table. The bottom panel shows a table of variables and their properties.

```

7 function y = my_fcn_fixpt(x)
8     fm = get_fimath();
9
10    y = fi(replacement_exp(x), 0, 16, 1, fm);
11 end
12
13
14 %
15 % Copyright 2017 The MathWorks, Inc.
16
17 % calculate replacement_exp via lookup table between extents x = fi([-10,10]),
18 % interpolation degree = 1, number of points = 1000
19 function y = replacement_exp( x )
20     persistent LUT
21     if ( isempty(LUT) )
22         T = numerictype( false, 16, 1);
23         LUT = fi([4.53999297624848e-05, 4.63179964587419e-05, 4.72546280836935e-05, ...
24                 4.82102000529591e-05, 4.91850953737242e-05, 5.01797047982559e-05, ...
25                 5.11944269805214e-05, 5.22296686359748e-05, 5.32858447045743e-05, ...
26                 5.43633785170962e-05, 5.54627019648123e-05, 5.6584255672598e-05, ...
27                 5.77284891755413e-05, 5.88958610991229e-05, 6.00868393430401e-05, ...
28                 6.13019012687477e-05, 6.25415338907916e-05, 6.38062340720112e-05, ...
29                 6.50965087226892e-05, 6.6412875003729e-05, 6.77558605339399e-05, ...
30                 6.91260036015147e-05, 7.05238533797832e-05, 7.19499701473294e-05, ...

```

| Variable | Type | Size | Signed | Word Length | Fraction Length |
|----------|-------------|-------|--------|-------------|-----------------|
| Input | | | | | |
| x | embedded.fi | 1 x 1 | Yes | 16 | 11 |
| Output | | | | | |
| y | embedded.fi | 1 x 1 | No | 16 | 1 |

The generated fixed-point function, `my_fcn_fixpt.m`, calls this approximation instead of calling `exp`. The fixed-point conversion process infers the ranges for the function and then uses an interpolated lookup table to replace the function. By default, the lookup table uses linear interpolation, 1000 points, and the minimum and maximum values detected by running the test file.

```

function y = my_fcn_fixpt(x)
    fm = get_fimath();

    y = fi(replacement_exp(x), 0, 16, 1, fm);
end

```

You can now test the generated fixed-point code and compare the results against the original MATLAB function. If the behavior of the generated fixed-point code does not match the behavior of the original code closely enough, modify the interpolation method or number of points used in the lookup table. Then, regenerate the code.

See Also

More About

- “Replacing Functions Using Lookup Table Approximations” on page 7-50

Convert Fixed-Point Conversion Project to MATLAB Scripts

This example shows how to convert a Fixed-Point Converter app project to a MATLAB script. You can use the `-tocode` option of the `fixedPointConverter` command to create a script for fixed-point conversion. You can use the script to repeat the project workflow in a command-line workflow. Before you convert the project to a script, you must complete the **Test** step of the fixed-point conversion process.

Prerequisites

This example uses the following files:

- Project file `ex_2ndOrder_filter.prj`
- Entry-point file `ex_2ndOrder_filter.m`
- Test bench file `ex_2ndOrder_filter_test.m`
- Generated fixed-point MATLAB file `ex_2ndOrder_filter_fixpt.m`

To obtain these files, complete the example “Propose Data Types Based on Simulation Ranges” on page 8-13, including the **Test** step.

Generate the Scripts

- 1 Change to the folder that contains the project file `ex_2ndOrder_filter.prj`.
- 2 Use the `-tocode` option of the `fixedPointConverter` command to convert the project to a script. Use the `-script` option to specify the file name for the script.

```
fixedPointConverter -tocode ex_2ndOrder_filter -script ex_2ndOrder_filter_script.m
```

The `fixedPointConverter` command generates a script in the current folder. `ex_2ndOrder_filter_script.m` contains the MATLAB commands to:

- Create a floating-point to fixed-point conversion configuration object that has the same fixed-point conversion settings as the project.
- Run the `fiaccel` command to convert the MATLAB function `ex_2ndOrder_filter` to the fixed-point MATLAB function `ex_2ndOrder_filter_fixpt`.

The `fiaccel` command overwrites existing files that have the same name as the generated script. If you omit the `-script` option, the `fiaccel` command returns the script in the Command Window.

Run Script That Generates Fixed-Point MATLAB Code

If you want to regenerate the fixed-point function, use the generated script.

- 1 Make sure that the current folder contains the entry-point function `ex_2ndOrder_filter.m` and the test bench file `ex_2ndOrder_filter_test.m`.
- 2 Run the script.

```
ex_2ndOrder_filter_script
```

The script generates `ex_2ndOrder_filter_fixpt.m` in the folder `codegen\ex_2ndOrder_filter\fixpt`. The variables `cfg` and `ARGS` appear in the base workspace.

See Also

`coder.FixPtConfig` | `fiaccel`

Related Examples

- “Propose Data Types Based on Simulation Ranges” on page 8-13

Replace a Custom Function with a Lookup Table

This example shows how to replace a custom function with a lookup table approximation function using the Fixed-Point Converter app.

Prerequisites

This example requires the following products:

- MATLAB
- Fixed-Point Designer
- C compiler

See https://www.mathworks.com/support/compilers/current_release/.

You can use `mex -setup` to change the default compiler. See “Change Default Compiler”.

Create Algorithm and Test Files

In a local, writable folder:

- 1 Create a MATLAB function, `custom_fcn.m` which is the function that you want to replace.

```
function y = custom_fcn(x)
    y = 1./(1+exp(-x));
end
```

- 2 Create a wrapper function, `call_custom_fcn.m`, that calls `custom_fcn.m`.

```
function y = call_custom_fcn(x)
    y = custom_fcn(x);
end
```

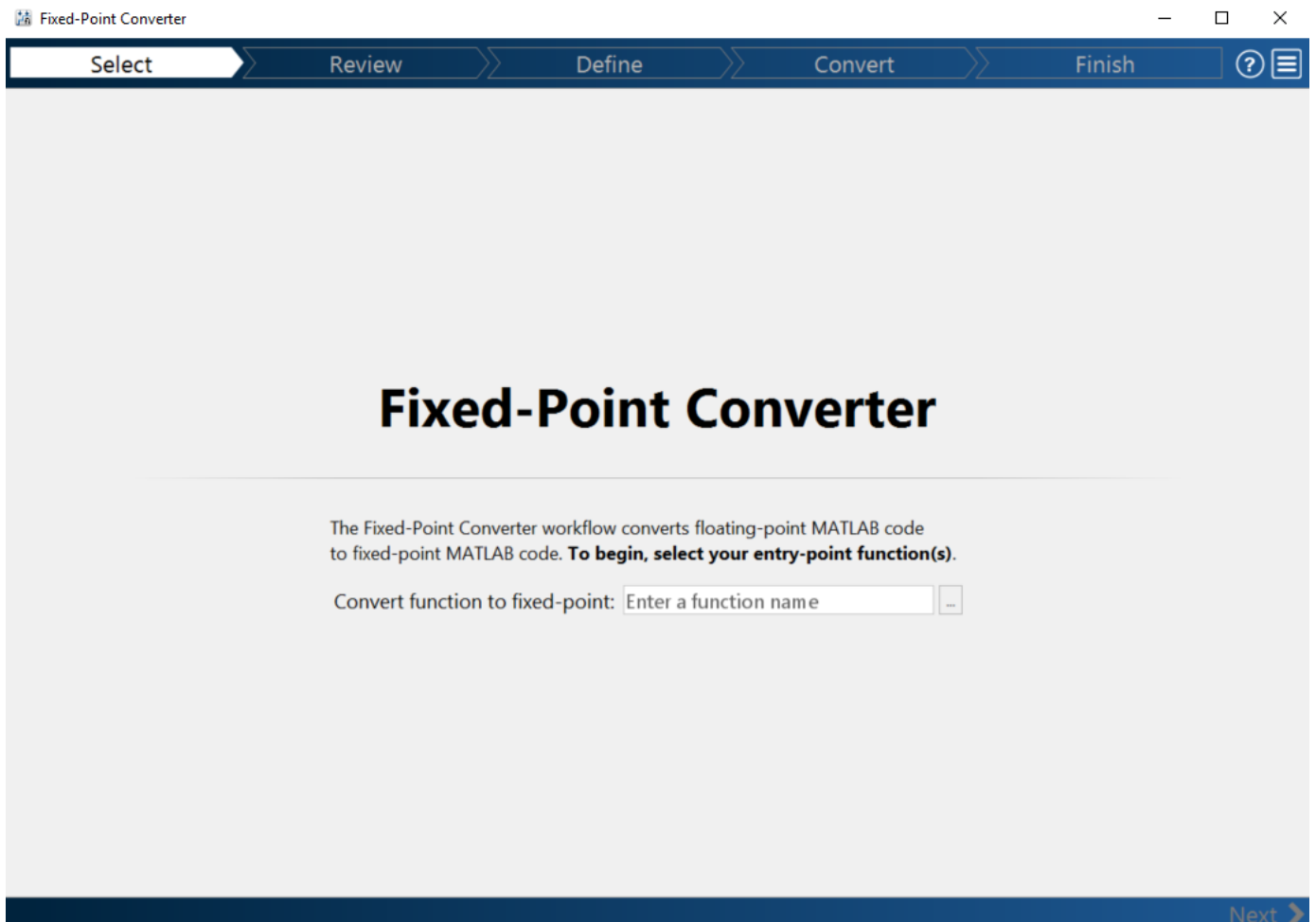
- 3 Create a test file, `custom_test.m`, that uses `call_custom_fcn`.

```
close all
clear all

x = linspace(-10,10,1e3);
for itr = 1e3:-1:1
    y(itr) = call_custom_fcn( x(itr) );
end
plot( x, y );
```

Open the Fixed-Point Converter App

- 1 Navigate to the work folder that contains the file for this example.
- 2 On the MATLAB Toolstrip **Apps** tab, under **Code Generation**, click the app icon.



Select Source Files

- 1 To add the entry-point function `call_custom_fcn` to the project, browse to the file `call_custom_fcn.m`, and then click **Open**. By default, the app saves information and settings for this project in the current folder in a file named `call_custom_fcn.prj`.
- 2 Click **Next** to go to the **Define Input Types** step.

The app screens `call_custom_fcn.m` for code violations and fixed-point conversion issues. The app opens the **Review Code Generation Readiness** page.

Review Code Generation Readiness

- 1 Click **Review Issues**. The app indicates that the `exp` function is not supported for fixed-point conversion. You can ignore this warning because you are going to replace `custom_fcn`, which is the function that calls `exp`.

Review Code Generation Readiness
REVIEW ISSUES

| | |
|-----------|--|
| Code | 1 <input type="checkbox"/> <code>function y = custom_fcn(x)</code> |
| tom_fcn.m | 2 <code>y = 1./(1+exp(-x));</code> |
| fcn.m | 3 <code>end</code> |
| | 4 |
| | 5 |

| Errors | | | |
|--------|------------|------|---|
| | Function | Line | Description |
| ⚠ | custom_fcn | 2 | exp is not supported for fixed-point conversion |

- 2 Click **Next** to go to the **Define Input Types** step.

Define Input Types

- 1 Add custom_test as a test file and then click **Autodefine Input Types**.

The test file runs. The app determines from the test file that x is a scalar double.

- 2 Click **Next** to go to the **Convert to Fixed Point** step.

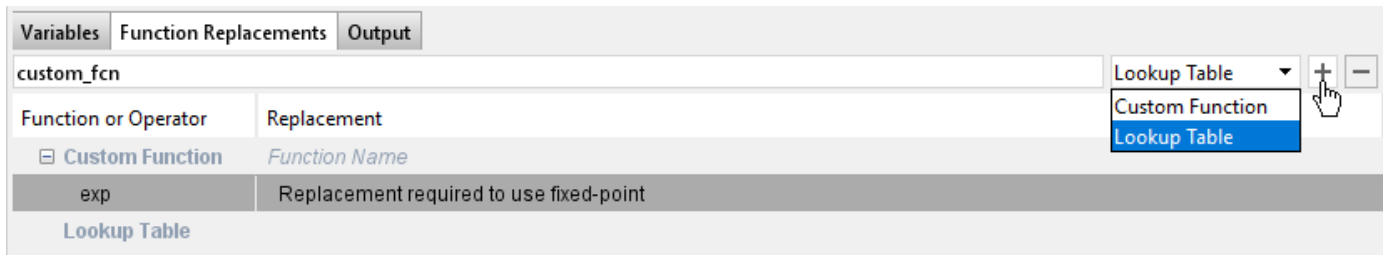
Replace custom_fcn with Lookup Table

- 1 Select the **Function Replacements** tab.

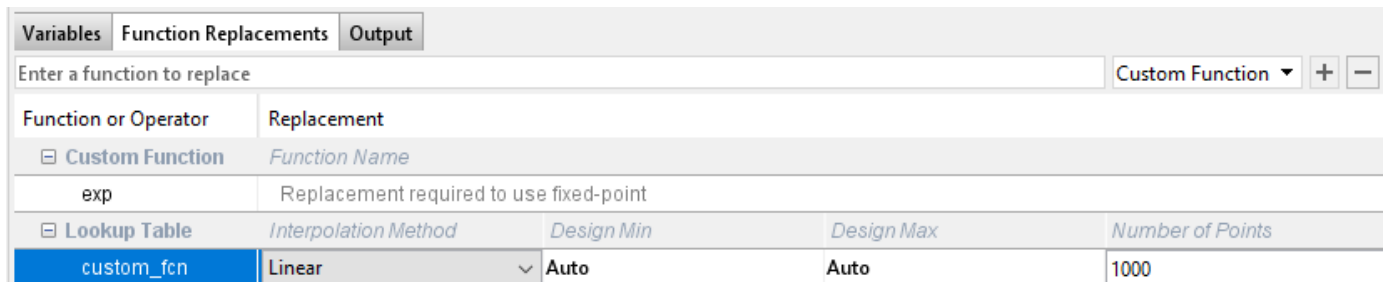
The app indicates that you must replace the exp function.


| Variables | Function Replacements | Output |
|--|---|-----------------------|
| Enter a function to replace | | Custom Function ▾ + - |
| Function or Operator | Replacement | |
| <input type="checkbox"/> Custom Function | <i>Function Name</i> | |
| exp | Replacement required to use fixed-point | |
| Lookup Table | | |

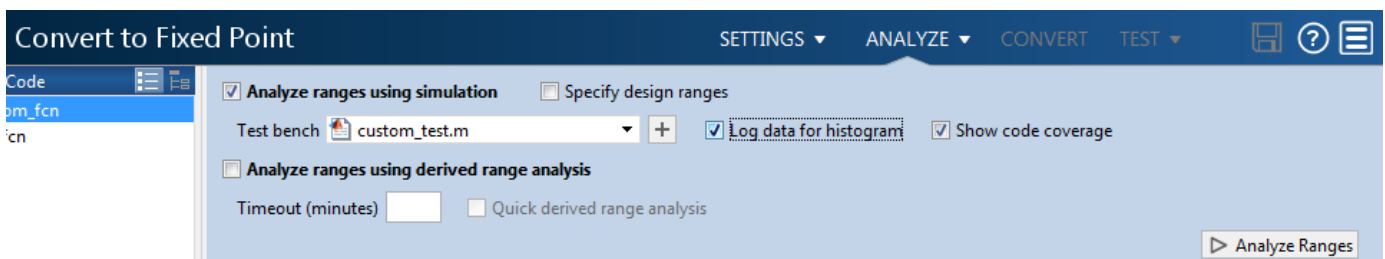
- 2 Enter the name of the function to replace, `custom_fcn`, select **Lookup Table**, and then click **+**.



The app adds `custom_fcn` to the list of functions that it will replace with a Lookup Table. By default, the lookup table uses linear interpolation and 1000 points. The app sets **Design Min** and **Design Max** to **Auto** which means that app uses the design minimum and maximum values that it detects by either running a simulation or computing derived ranges.



- 3 Click the **Analyze** arrow , select **Log data for histogram**, and verify that the test file is `call_custom_test`.



- 4 Click **Analyze**.

The simulation runs. The app displays simulation minimum and maximum ranges on the **Variables** tab. Using the simulation range data, the software proposes fixed-point types for each variable based on the default type proposal settings, and displays them in the **Proposed Type** column. The **Convert** option is now enabled.

- 5 Examine the proposed types and verify that they cover the full simulation range. To view logged histogram data for a variable, click its **Proposed Type** field. The histogram provides range information and the percentage of simulation range covered by the proposed data type.

| Variable | Type | Sim Min | Sim Max | Whole ... | Proposed Type | Log... | Max Diff |
|----------|--------|---------|---------|-----------|------------------------|--------|----------|
| Input | | | | | | | |
| x | double | -10 | 10 | No | numerictype(1, 16, 11) | | |
| Output | | | | | | | |
| y | double | 0 | | | | | |

Convert to Fixed Point

- 1 Click **Convert**.

The app validates the proposed types and generates a fixed-point version of the entry-point function, `call_custom_fcn_fixpt.m`.

- 2 In the Output Files list, select `call_custom_fcn_fixpt.m`.

The conversion process generates a lookup table approximation, `replacement_custom_fcn`, for the `custom_fcn` function. The fixed-point conversion process infers the ranges for the function and then uses an interpolated lookup table to replace the function. By default, the lookup table uses linear interpolation, 1000 points, and the minimum and maximum values detected by running the test file.

The generated fixed-point function, `call_custom_fcn_fixpt.m`, calls this approximation instead of calling `custom_fcn`.

```
function y = call_custom_fcn_fixpt(x)
    fm = get_fimath();

    y = fi(replacement_custom_fcn(x), 0, 16, 16, fm);
end
```

You can now test the generated fixed-point code and compare the results against the original MATLAB function. If the behavior of the generated fixed-point code does not match the behavior of the original code closely enough, modify the interpolation method or number of points used in the lookup table and then regenerate code.

See Also

More About

- “Replacing Functions Using Lookup Table Approximations” on page 7-50

Visualize Differences Between Floating-Point and Fixed-Point Results

This example shows how to configure the Fixed-Point Converter app to use a custom plot function to compare the behavior of the generated fixed-point code against the behavior of the original floating-point MATLAB code.

By default, when the **Log inputs and outputs for comparison plots** option is enabled, the conversion process uses a time series based plotting function to show the floating-point and fixed-point results and the difference between them. However, during fixed-point conversion you might want to visualize the numerical differences in a view that is more suitable for your application domain. This example shows how to customize plotting and produce scatter plots at the test numerics step of the fixed-point conversion.

Copy Relevant Files

Copy the `myFilter.m`, `myFilterTest.m`, `plotDiff.m`, and `filterData.mat` files to a local working folder.

Prerequisites

This example requires the following products:

- MATLAB
- Fixed-Point Designer
- C compiler

See https://www.mathworks.com/support/compilers/current_release/.

You can use `mex -setup` to change the default compiler. See “Change Default Compiler”.

Inspect Example Files

It is a best practice is to create a separate test script to do pre- and post-processing, such as:

- Loading inputs.
- Setting up input values.
- Outputting test results.

For more information, see “Create a Test File” on page 11-3.

| Type | Name | Description |
|-------------------|-----------------------------|--|
| Function code | <code>myFilter.m</code> | Entry-point MATLAB function |
| Test file | <code>myFilterTest.m</code> | MATLAB script that tests <code>myFilter.m</code> |
| Plotting function | <code>plotDiff.m</code> | Custom plot function |
| MAT-file | <code>filterData.mat</code> | Data to filter. |

The myFilter Function

```
function [y, ho] = myFilter(in)

persistent b h;
if isempty(b)
    b = complex(zeros(1,16));
    h = complex(zeros(1,16));
    h(8) = 1;
end

b = [in, b(1:end-1)];
y = b*h.';

errf = 1-sqrt(real(y)*real(y) + imag(y)*imag(y));
update = 0.001*conj(b)*y*errf;

h = h + update;
h(8) = 1;
ho = h;

end
```

The myFilterTest File

```
% load data
data = load('filterData.mat');
d = data.symbols;

for idx = 1:4000
    y = myFilter(d(idx));
end
```

The plotDiff Function

```
% varInfo - structure with information about the variable.
% It has the following fields
%         i) name
%        ii) functionName
% floatVals - cell array of logged original values for
% the 'varInfo.name' variable
% fixedVals - cell array of logged values for the
% 'varInfo.name' variable after Fixed-Point conversion.
function plotDiff(varInfo, floatVals, fixedVals)
    varName = varInfo.name;
    fcnName = varInfo.functionName;

    % escape the '_'s because plot titles treat these as subscripts
    escapedVarName = regexp(varName, '_', '\\_');
    escapedFcnName = regexp(fcnName, '_', '\\_');

    % flatten the values
    flatFloatVals = floatVals(1:end);
    flatFixedVals = fixedVals(1:end);

    % build Titles
    floatTitle = [escapedFcnName ' > ' 'float : ' escapedVarName];
    fixedTitle = [escapedFcnName ' > ' 'fixed : ' escapedVarName];
```

```
data = load('filterData.mat');

switch varName
    case 'y'
        x_vec = data.symbols;

        figure('Name','Comparison plot','NumberTitle','off');

        % plot floating point values
        y_vec = flatFloatVals;
        subplot(1, 2, 1);
        plotScatter(x_vec, y_vec, 100, floatTitle);

        % plot fixed point values
        y_vec = flatFixedVals;
        subplot(1, 2, 2);
        plotScatter(x_vec, y_vec, 100, fixedTitle);

    otherwise
        % Plot only output 'y' for this example, skip the rest
end

end

function plotScatter(x_vec, y_vec, n, figTitle)
    % plot the last n samples
    x_plot = x_vec(end-n+1:end);
    y_plot = y_vec(end-n+1:end);

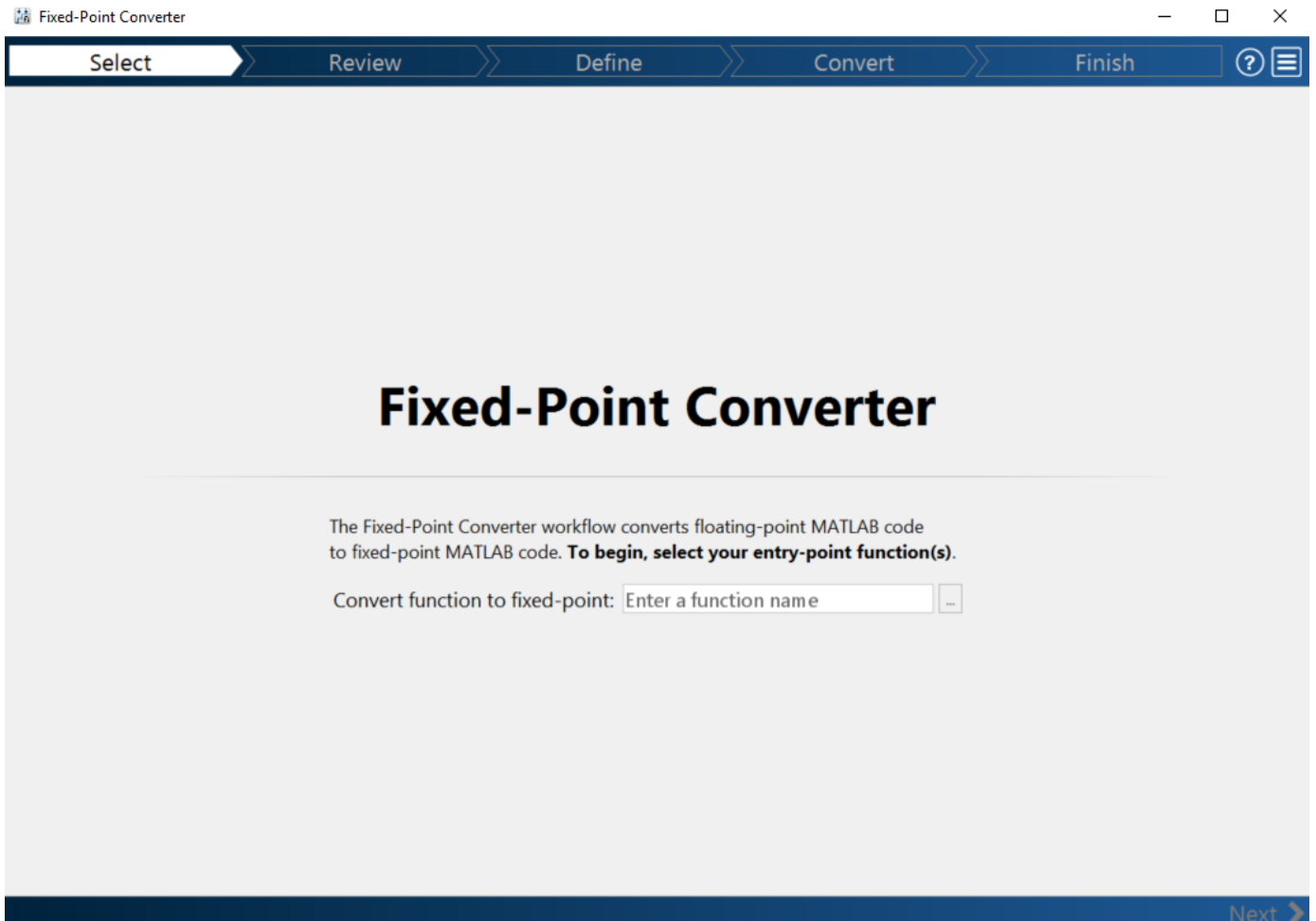
    hold on
    scatter(real(x_plot),imag(x_plot), 'bo');

    hold on
    scatter(real(y_plot),imag(y_plot), 'rx');

    title(figTitle);
end
```

Open the Fixed-Point Converter App

- 1 Navigate to the folder that contains the files for this example.
- 2 On the MATLAB Toolstrip **Apps** tab, under **Code Generation**, click the app icon.



Select Source Files

- 1 To add the entry-point function `myFilter` to the project, browse to the file `myFilter.m`, and then click **Open**.

By default, the app saves information and settings for this project in the current folder in a file named `myFilter.prj`.

- 2 Click **Next** to go to the **Define Input Types** step.

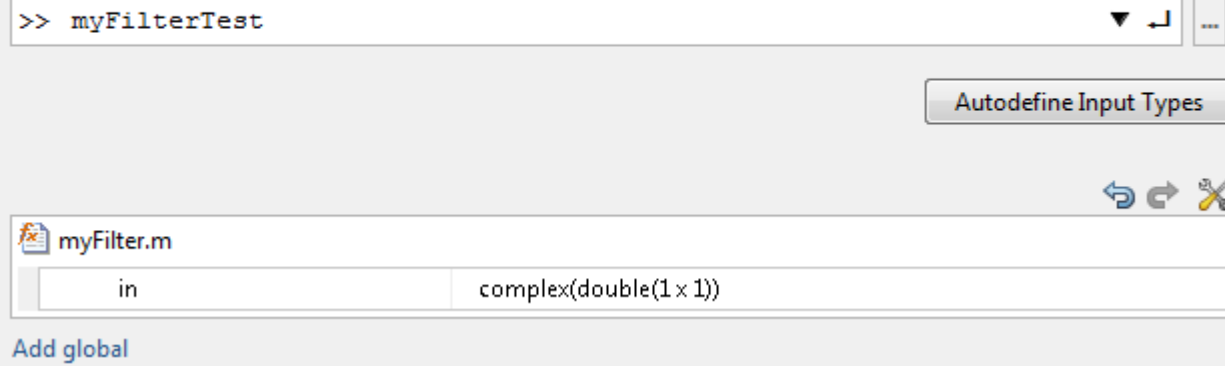
The app screens `myFilter.m` for code violations and fixed-point conversion readiness issues. The app does not find issues in `myFilter.m`.

Define Input Types

- 1 On the **Define Input Types** page, to add `myFilterTest` as a test file, browse to `myFilterTest.m`, and then click **Open**.
- 2 Click **Autodefine Input Types**.

The app determines from the test file that the input type of `in` is `complex(double(1x1))`.

To **automatically define input types**, call `myFilter` or enter a script that calls `myFilter` in the MATLAB prompt below:



- 3 Click **Next** to go to the **Convert to Fixed Point** step.

Convert to Fixed Point

- 1 The app generates an instrumented MEX function for your entry-point MATLAB function. The app displays compiled information for variables in your code. For more information, see “View and Modify Variable Information” on page 8-35.

Convert to Fixed Point

SETTINGS ▾ ANALYZE ▾ CONVERT TEST ▾

```

1 function [y, ho] = myFilter(in)
2
3 persistent b h;
4 if isempty(b)
5     b = complex(zeros(1,16));
6     h = complex(zeros(1,16));
7     h(8) = 1;
8 end
9
10 b = [in, b(1:end-1)];
11 y = b*h.';
12
13 errf = 1-sqrt(real(y)*real(y) + imag(y)*imag(y));
14 update = 0.001*conj(b)*y*errf;
15
16 h = h + update;
17 h(8) = 1;
18 ho = h;
19
20 end

```

| Variable | Type | Sim Min | Sim Max | Whole Number | Proposed Type |
|------------|-------------------|---------|---------|--------------|---------------|
| Input | | | | | |
| in | complex double | | | No | |
| Output | | | | | |
| y | complex double | | | No | |
| ho | 1 x 16 complex... | | | No | |
| Persistent | | | | | |
| b | 1 x 16 complex... | | | No | |

- 2 To open the settings dialog box, click the **Settings** arrow ▾.
 - a Verify that **Default word length** is set to 16.
 - b Under **Advanced**, set **Signedness** to Signed
 - c Under **Plotting and Reporting**, set **Custom plot function** to plotDiff.
- 3 Click the **Analyze** arrow ▾. Verify that the test file is myFilterTest.
- 4 Click **Analyze**.

The test file, myFilterTest, runs and the app displays simulation minimum and maximum ranges on the **Variables** tab. Using the simulation range data, the software proposes fixed-point types for each variable based on the default type proposal settings, and displays them in the **Proposed Type** column.

Convert to Fixed Point

SETTINGS ▾ ANALYZE ▾ CONVERT TEST ▾

```

1 function [y, ho] = myFilter(in)
2
3 persistent b h;
4 if isempty(b)
5     b = complex(zeros(1,16));
6     h = complex(zeros(1,16));
7     h(8) = 1;
8 end
9
10 b = [in, b(1:end-1)];
11 y = b*h.';
12
13 errf = 1-sqrt(real(y)*real(y) + imag(y)*imag(y));
14 update = 0.001*conj(b)*y*errf;
15
16 h = h + update;
17 h(8) = 1;
18 ho = h;
19
20 end

```

| Variable | Type | Sim Min | Sim Max | Whole Number | Proposed Type |
|-------------------|-------------------|---------|---------|--------------|------------------------|
| Input | | | | | |
| in | complex double | -0.95 | 0.95 | No | numerictype(1, 16, 15) |
| Output | | | | | |
| y | complex double | -0.95 | 0.95 | No | numerictype(1, 16, 15) |
| ho | 1 x 16 complex... | -0.13 | 1 | No | numerictype(1, 16, 14) |
| Persistent | | | | | |
| b | 1 x 16 complex... | -0.95 | 0.95 | No | numerictype(1, 16, 15) |

Next >

- To convert the floating-point algorithm to fixed point, click **Convert**.

The software validates the proposed types and generates a fixed-point version of the entry-point function.

```

7 function [y, ho] = myFilter_fixpt(in)
8
9 fm = get_fimath();
10
11 persistent b h;
12 if isempty(b)
13     b = fi(complex(zeros(1,16)), 1, 16, 15, fm);
14     h = fi(complex(zeros(1,16)), 1, 16, 14, fm);
15     h(8) = 1;
16 end
17
18 b(:) = [fi(in, 1, 16, 15, fm), b(1:end-1)];
19 y = fi(b*h.', 1, 16, 15, fm);
20
21 errf = fi(fi_signed(fi(1, 1, 2, 0, fm))-sqrt(real(y)*real(y) + imag(y)*imag(y)), 1, 16, 14, fm);
22 update = fi(fi(0.001, 1, 16, 24, fm)*conj(b)*y*errf, 1, 16, 25, fm);
23
24 h(:) = h + update;
25 h(8) = 1;
26 ho = fi(h, 1, 16, 14, fm);
27
28 end
29


```

| Variable | Type | Size | Signed | Word Length | Fraction Length |
|-------------------|-------------|--------|--------|-------------|-----------------|
| Input | | | | | |
| in | embedded.fi | 1 x 1 | Yes | 16 | 15 |
| Output | | | | | |
| y | embedded.fi | 1 x 1 | Yes | 16 | 15 |
| ho | embedded.fi | 1 x 16 | Yes | 16 | 14 |
| Persistent | | | | | |
| b | embedded | 16 | Yes | 16 | 15 |

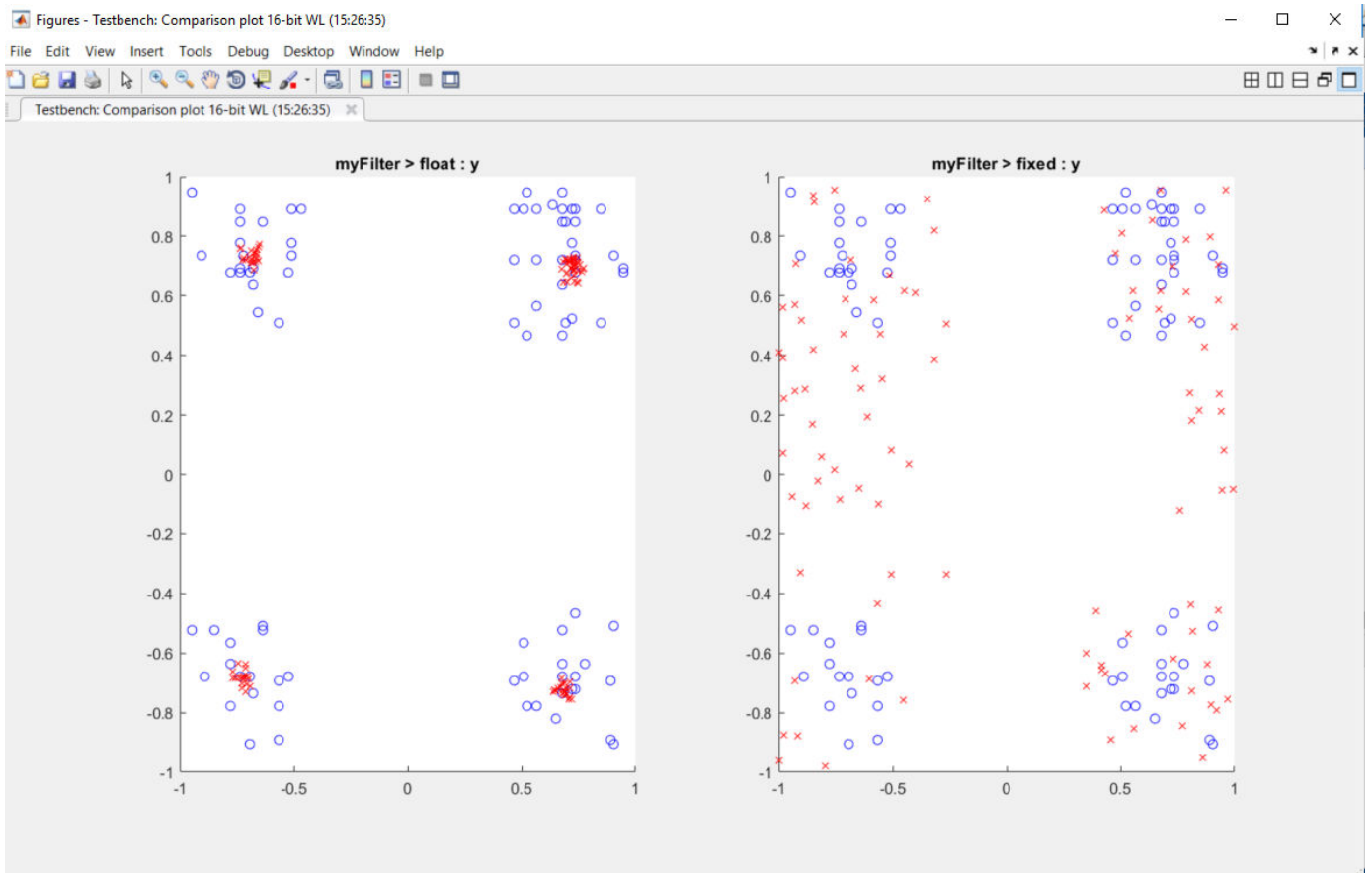
Validation succeeded

Next

Test Numerics and View Comparison Plots

- 1 Click **Test** arrow , select **Log inputs and outputs for comparison plots**, and then click **Test**.

The app runs the test file that you used to define input types to test the fixed-point MATLAB code. Because you selected to log inputs and outputs for comparison plots and to use the custom plotting function, `plotDiff.m`, for these plots, the app uses this function to generate the comparison plot. The plot shows that the fixed-point results do not closely match the floating-point results.

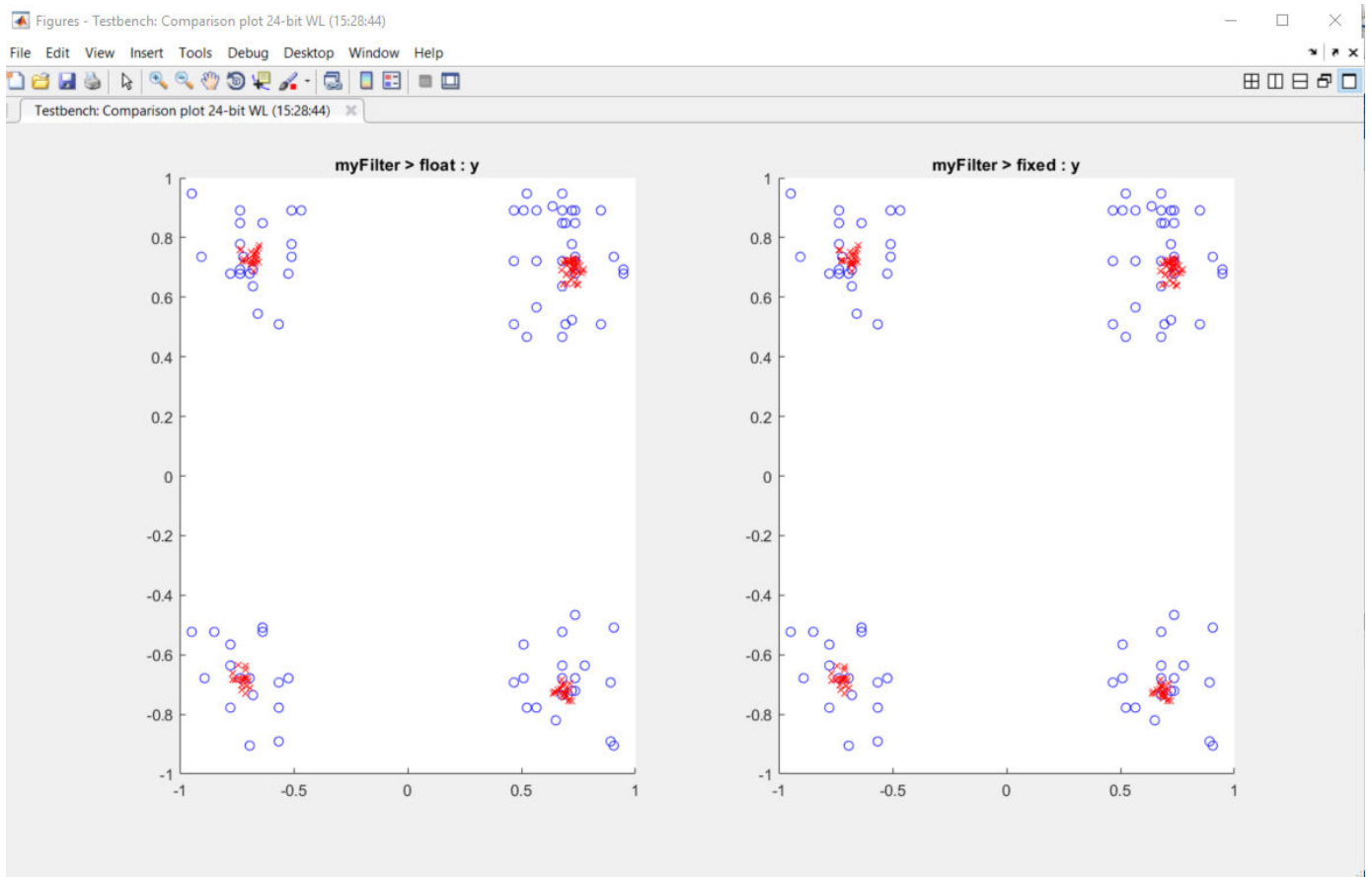


- 2 In the settings, increase the **DefaultWordLength** to 24 and then convert to fixed point again.

The app converts `myFilter.m` to fixed point and proposes fixed-point data types using the new default word length.

- 3 Run the test numerics step again.

The increased word length improves the results. This time, the plot shows that the fixed-point results match the floating-point results.



See Also

More About

- "Custom Plot Functions" on page 7-51

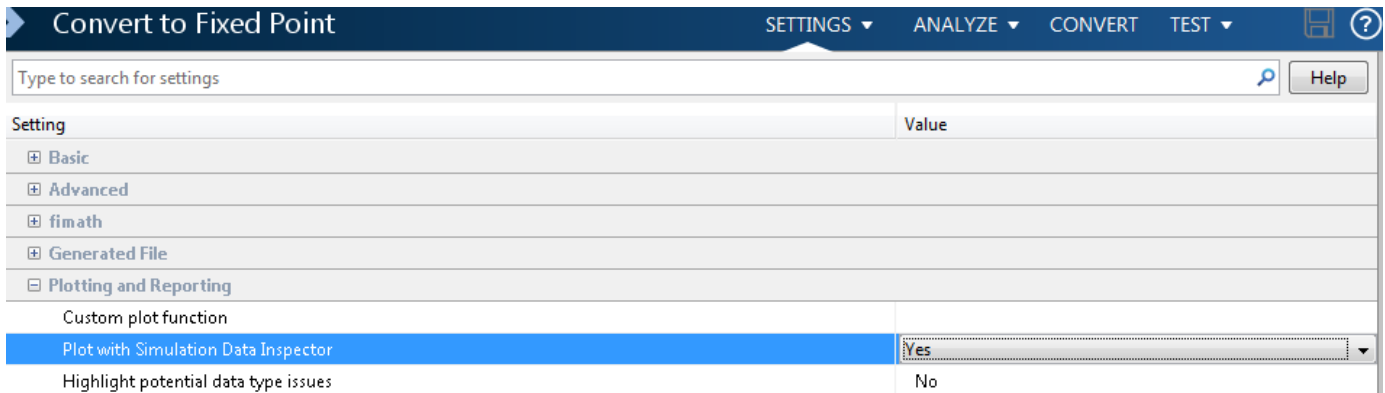
Enable Plotting Using the Simulation Data Inspector


You can use the Simulation Data Inspector with the Fixed-Point Converter app to inspect and compare floating-point and fixed-point logged input and output data.

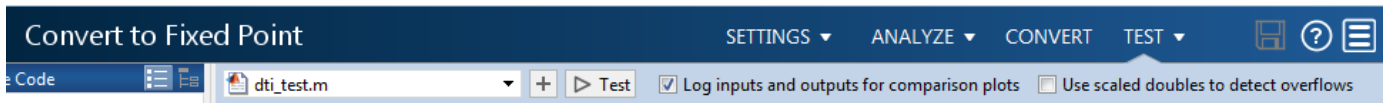
- 1 On the **Convert to Fixed Point** page,

Click the **Settings** arrow .

- 2 Expand the **Plotting and Reporting** settings and set **Plot with Simulation Data Inspector** to **Yes**.



- 3 Click the **Test** arrow . Select **Log inputs and outputs for comparison plots**, and then click **Test**.



For an example, see “Propose Data Types Based on Derived Ranges” on page 8-24.

Add Global Variables by Using the App

To add global variables to the project:

- 1** On the **Define Input Types** page, automatically define input types or click **Let me enter input or global types directly**.

The app displays a table of entry-point inputs.

- 2** To add a global variable, click **Add global**.

By default, the app names the first global variable in a project **g**, and subsequent global variables **g1**, **g2**, and so on.

- 3** Under **Global variables**, enter a name for the global variable.
- 4** After adding a global variable, but before generating code, specify its type and initial value. Otherwise, you must create a variable with the same name in the global workspace. See “Specify Global Variable Type and Initial Value Using the App” on page 8-80.

Automatically Define Input Types by Using the App

If you specify a test file that calls the project entry-point functions, the Fixed-Point Converter app can infer the input argument types by running the test file. If a test file calls an entry-point function multiple times with different size inputs, the app takes the union of the inputs. The app infers that the inputs are variable size, with an upper bound equal to the size of the largest input.

Before using the app to automatically define function input argument types, you must add at least one entry-point file to your project. You must also specify code that calls your entry-point functions with the expected input types. It is a best practice to provide a test file that calls your entry-point functions. The test file can be either a MATLAB function or a script. The test file must call the entry-point function at least once.

To automatically define input types:

- 1 On the **Define Input Types** page, specify a test file. Alternatively, you can enter code directly.
- 2 Click **Autodefine Input Types**.

The app runs the test file and infers the types for entry-point input arguments. The app displays the inferred types.

Note If you automatically define the input types, the entry-point functions must be in a writable folder.

If your test file does not call an entry-point function with different size inputs, the resulting type dimensions are fixed-size. After you define the input types, you can specify and apply rules for making type dimensions variable-size when they meet a size threshold. See “Make Dimensions Variable-Size When They Meet Size Threshold” (MATLAB Coder).

Define Constant Input Parameters Using the App

- 1 On the **Define Input Types** page, click **Let me enter input or global types directly**.
- 2 Click the field to the right of the input parameter name.
- 3 Select **Define Constant**.
- 4 In the field to the right of the parameter name, enter the value of the constant or a MATLAB expression that represents the constant.

The app uses the value of the specified MATLAB expression as a compile-time constant.

Define or Edit Input Parameter Type by Using the App

| In this section... |
|---|
| "Define or Edit an Input Parameter Type" on page 8-66 |
| "Specify a String Scalar Input Parameter" on page 8-67 |
| "Specify an Enumerated Type Input Parameter" on page 8-67 |
| "Specify a Fixed-Point Input Parameter" on page 8-68 |
| "Specify a Structure Input Parameter" on page 8-68 |
| "Specify a Cell Array Input Parameter" on page 8-69 |

Define or Edit an Input Parameter Type

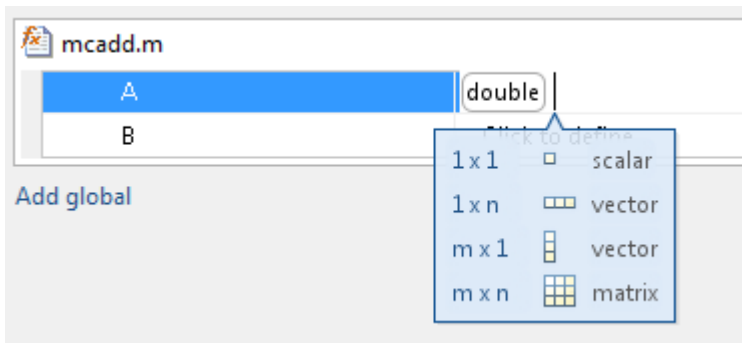
The following procedure shows you how to define or edit `double`, `single`, `int64`, `int32`, `int16`, `int8`, `uint64`, `uint32`, `uint16`, `uint8`, `logical`, and `char` types.

For more information about defining other types, see the information in this table.

| Input Type | Link |
|---|---|
| A string scalar (1-by-1 string array) | "Specify a String Scalar Input Parameter" on page 8-67 |
| A structure (struct) | "Specify a Structure Input Parameter" on page 8-68 |
| A cell array (cell (Homogeneous) or cell (Heterogeneous)) | "Specify a Cell Array Input Parameter" on page 8-69 |
| A fixed-point data type (embedded.fi) | "Specify a Fixed-Point Input Parameter" on page 8-68 |
| An input by example (Define by Example) | "Define Input Parameter by Example by Using the App" on page 8-73 |
| A constant (Define Constant) | "Define Constant Input Parameters Using the App" on page 8-65 |

- 1 Click the field to the right of the input parameter name.
- 2 Optionally, for numeric types, to make the parameter a complex type, select the **Complex number** check box.
- 3 Select the input type.

The app displays the selected type. It displays and the size options.



- 4 From the list, select whether your input is a scalar, a $1 \times n$ vector, a $m \times 1$ vector, or a $m \times n$ matrix. By default, if you do not select a size option, the app defines inputs as scalars.
- 5 Optionally, if your input is not scalar, enter sizes m and n . You can specify:
 - Fixed size, for example, 10.
 - Variable size, up to a specified limit, by using the `:` prefix. For example, to specify that your input can vary in size up to 10, enter `:10`.
 - Unbounded variable size by entering `:Inf`.

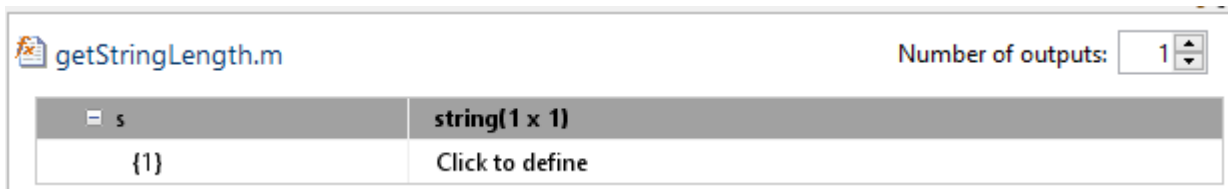
You can edit the size of each dimension.

Specify a String Scalar Input Parameter

To specify that an input is a string scalar:

- 1 On the **Define Input Types** page, click **Let me enter input or global types directly**.
- 2 Click the field to the right of the input parameter that you want to define.
- 3 Select **string**. Then select `1x1 scalar`.

The type is a 1-by-1 string array (string scalar) that contains a character vector.



- 4 To specify the size of the character vector, click the field to the right of the string array element `{1}`. Select **char**. Then, select `1xn vector` and enter the size.
- 5 To make the string variable-size, click the second dimension.
 - To specify that the second dimension is unbounded, select `:Inf`.
 - To specify that the second dimension has an upper bound, enter the upper bound, for example 8. Then, select `:8`.

Specify an Enumerated Type Input Parameter

To specify that an input uses the enumerated type `MyColors`:

- 1 Suppose that the enumeration `MyColors` is on the MATLAB path.

```
classdef MyColors < int32
    enumeration
        green(1),
        red(2),
    end
end
```


- 2 On the **Define Input Types** page, click **Let me enter input or global types directly**.
- 3 In the field to the right of the input parameter, enter `MyColors`.

Specify a Fixed-Point Input Parameter

To specify fixed-point inputs, Fixed-Point Designer software must be installed.

- 1 On the **Define Input Types** page, click **Let me enter input or global types directly**.
- 2 Click the field to the right of the input parameter that you want to define.
- 3 Select **embedded.fi**.
- 4 Select the size. If you do not specify the size, the size defaults to 1x1.
- 5 Specify the input parameter `numerictype` and `fimath` properties.

If you do not specify a local `fimath`, the app uses the default `fimath`. See “Default `fimath` Usage to Share Arithmetic Rules” on page 3-17.

To modify the `numerictype` or `fimath` properties, open the properties dialog box. To open the properties dialog box, click to the right of the fixed-point type definition. Optionally, click .

Specify a Structure Input Parameter

When a primary input is a structure, the app treats each field as a separate input. Therefore, you must specify properties for all fields of a primary structure input in the order that they appear in the structure definition:

- For each field of an input structure, specify class, size, and complexity.
- For each field that is a fixed-point class, also specify `numerictype`, and `fimath`.

Specify Structures by Type

- 1 On the **Define Input Types** page, click **Let me enter input or global types directly**.
- 2 Click the field to the right of the input parameter that you want to define.
- 3 Select **struct**.

The app displays the selected type, `struct`. The app displays the size options.

- 4 Specify that your structure is a scalar, 1 × n vector, m × 1 vector, or m × n matrix. By default, if you do not select a size option, the app defines inputs as scalars.
- 5 If your input is not scalar, enter sizes for each dimension. Click the dimension. Enter the size. Select from the size options. For example, for size 10:
 - To specify fixed size, select 10.

- To specify variable size with an upper bound of 10, select :10.
 - To specify unbounded variable size, select :Inf.
- 6 Add fields to the structure. Specify the class, size, and complexity of the fields. See “Add a Field to a Structure” on page 8-69.

Rename a Field in a Structure

Select the name field of the structure that you want to rename. Enter the new name.

Add a Field to a Structure

- 1 To the right of the structure, click **+**
- 2 Enter the field name. Specify the class, size, and complexity of the field.

Insert a Field into a Structure

- 1 Select the structure field below which you want to add another field.
- 2 Right-click the structure field.
- 3 Select **Insert Field Below**.

The app adds the field after the field that you selected.

- 4 Enter the field name. Specify the class, size, and complexity of the field.

Remove a Field from a Structure

- 1 Right-click the field that you want to remove.
- 2 Select **Remove Field**.

Specify a Cell Array Input Parameter

Note The Fixed-Point Converter app does not support cell arrays.

For code generation, cell arrays are homogeneous or heterogeneous. . A homogeneous cell array is represented as an array in the generated code. All elements have the same properties. A heterogeneous cell array is represented as a structure in the generated code. Elements can have different properties.

Specify a Homogeneous Cell Array

- 1 On the **Define Input Types** page, click **Let me enter input or global types directly**.
- 2 Click the field to the right of the input parameter that you want to define.
- 3 Select **cell (Homogeneous)**.

The app displays the selected type, `cell`. The app displays the size options.

- 4 From the list, select whether your input is a scalar, a $1 \times n$ vector, a $m \times 1$ vector, or a $m \times n$ matrix. By default, if you do not select a size option, the app defines inputs as scalars.
- 5 If your input is not scalar, enter sizes for each dimension. Click the dimension. Enter the size. Select from the size options. For example, for size 10:

- To specify fixed size, select `10`.
- To specify variable size with an upper bound of `10`, select `:10`.
- To specify unbounded variable size, select `:Inf`.

Below the cell array variable, a colon inside curly braces `{:}` indicates that the cell array elements have the same properties (class, size, and complexity).

- 6 To specify the class, size, and complexity of the elements in the cell array, click the field to the right of `{:}`.

Specify a Heterogeneous Cell Array

- 1 On the **Define Input Types** page, click **Let me enter input or global types directly**.
- 2 Click the field to the right of the input parameter that you want to define.
- 3 Select **cell (Heterogeneous)**.

The app displays the selected type, `cell`. The app displays the size options.


- 4 Specify that your structure is a scalar, `1 x n` vector, `m x 1` vector, or `m x n` matrix. By default, if you do not select a size option, the app defines inputs as scalars.
- 5 Optionally, if your input is not scalar, enter sizes `m` and `n`. A heterogeneous cell array is fixed size.

The app lists the cell array elements. It uses indexing notation to specify each element. For example, `{1,2}` indicates the element in row 1, column 2.

- 6 Specify the class, size, and complexity for each cell array element.
- 7 Optionally, add elements. See “Add an Element to a Heterogeneous Cell Array” on page 8-72

Set Structure Properties for a Heterogeneous Cell Array

A heterogeneous cell array is represented as a structure in the generated code. You can specify the properties for the structure that represents the cell array.

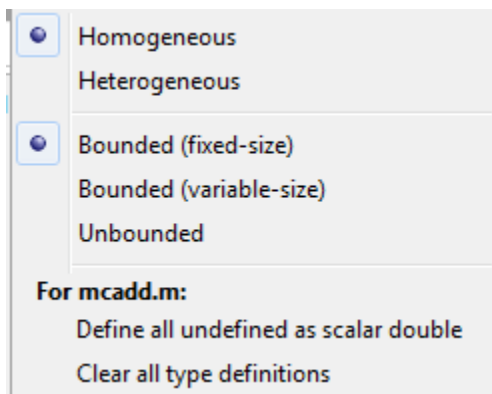
- 1 Click to the right of the cell array definition. Optionally click .
- 2 In the dialog box, specify properties for the structure in the generated code.

| Property | Description |
|---------------------------------------|--|
| C type definition name | Name for the structure type in the generated code. |
| Type definition is externally defined | <p>Default: No — type definition is not externally defined.</p> <p>If you select Yes to declare an externally defined structure, the app does not generate the definition of the structure type. You must provide it in a custom include file.</p> <p>Dependency: C type definition name enables this option.</p> |

| Property | Description |
|-------------------------------|---|
| C type definition header file | <p>Name of the header file that contains the external definition of the structure, for example, "mystruct.h". Specify the path to the file using the Additional include directories parameter on the project settings dialog box Custom Code tab.</p> <p>By default, the generated code contains <code>#include</code> statements for custom header files after the standard header files. If a standard header file refers to the custom structure type, then the compilation fails. If you specify the C type definition header file, the app includes that header file exactly at the point where it is required.</p> <p>Dependency: When Type definition is externally defined is set to Yes, this option is enabled.</p> |
| Data alignment boundary | <p>The run-time memory alignment of structures of this type in bytes.</p> <p>Alignment must be either -1 or a power of 2 that is no more than 128.</p> <p>Default: 0</p> <p>Dependency: When Type definition is externally defined is set to Yes, this option is enabled.</p> |

Change Classification as Homogeneous or Heterogeneous

To change the classification as homogeneous or heterogeneous, right-click the variable. Select **Homogeneous** or **Heterogeneous**.



The app clears the definitions of the elements.

Change the Size of the Cell Array

- 1 In the definition of the cell array, click a dimension. Specify the size.
- 2 For a homogeneous cell array, specify whether the dimension is variable size and whether the dimension is bounded or unbounded. Alternatively, right-click the variable. Select **Bounded (fixed-size)**, **Bounded (variable-size)**, or **Unbounded**
- 3 For a heterogeneous cell array, the app adds elements so that the cell array has the specified size and shape.

Add an Element to a Heterogeneous Cell Array

- 1 In the definition of the cell array, click a dimension. Specify the size. For example, enter 1 for the first dimension and 4 for the second dimension.

The app adds elements so that the cell array has the specified size and shape. For example for a 1x4 heterogeneous cell array, the app lists four elements: {1, 1}, {1, 2}, {1, 3}, and {1, 4}.

- 2 Specify the properties of the new elements.

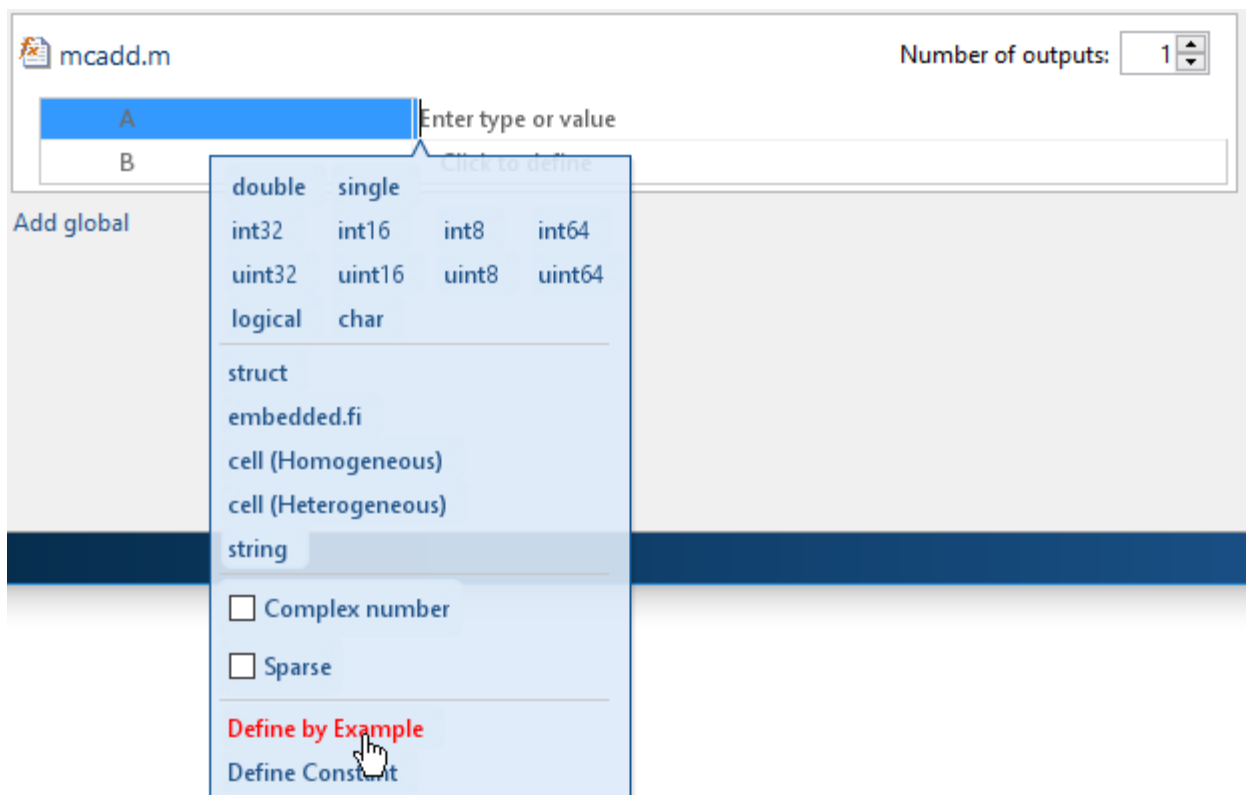
Define Input Parameter by Example by Using the App

In this section...

- “Define an Input Parameter by Example” on page 8-73
- “Specify Input Parameters by Example” on page 8-74
- “Specify a String Scalar Input Parameter by Example” on page 8-75
- “Specify a Structure Type Input Parameter by Example” on page 8-75
- “Specify a Cell Array Type Input Parameter by Example” on page 8-76
- “Specify an Enumerated Type Input Parameter by Example” on page 8-77
- “Specify a Fixed-Point Input Parameter by Example” on page 8-78
- “Specify an Input from an Entry-Point Function Output Type” on page 8-79

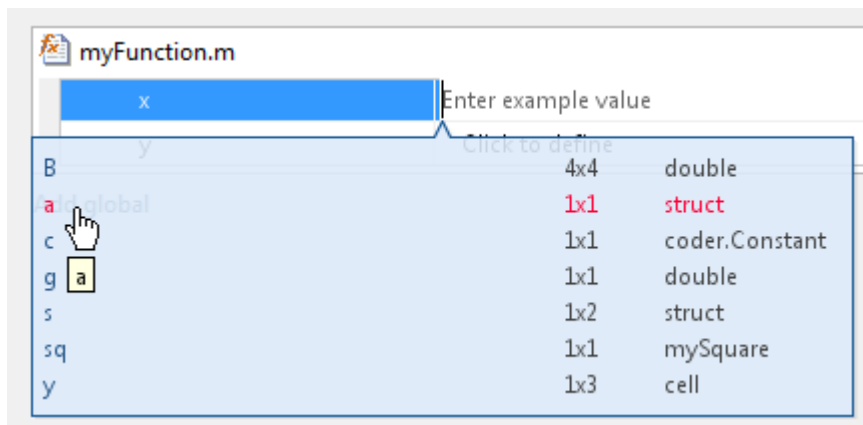
Define an Input Parameter by Example

- 1 On the **Define Input Types** page, click **Let me enter input or global types directly**.
- 2 Click the field to the right of the input parameter that you want to define.



- 3 Select **Define by Example**.
- 4 In the field to the right of the parameter, enter a MATLAB expression. The variable has the class, size, and complexity of the value of the expression.

Alternatively, you can select a variable from the list of workspace variables that displays.



Specify Input Parameters by Example

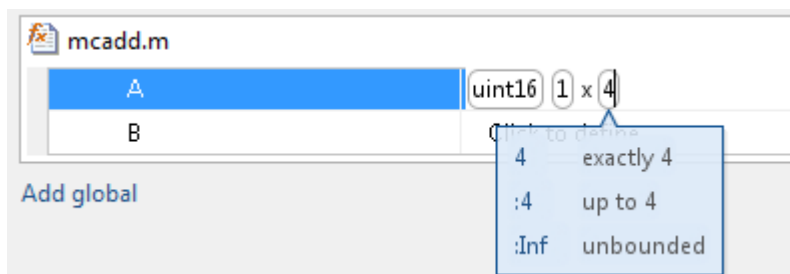
This example shows how to specify a 1-by-4 vector of unsigned 16-bit integers.

- 1 On the **Define Input Types** page, click **Let me enter input or global types directly**.
- 2 Click the field to the right of the input parameter that you want to define.
- 3 Select **Define by Example**.
- 4 In the field to the right of the parameter, enter:

```
zeros(1,4,'uint16')
```

The input type is `uint16(1x4)`.

- 5 Optionally, after you specify the input type, you can specify that the input is variable size. For example, select the second dimension.



- 6 To specify that the second dimension is variable size with an upper bound of 4, select `:4`. Alternatively, to specify that the second dimension is unbounded, select `:Inf`.

Alternatively, you can specify that the input is variable size by using the `coder.newtype` function. Enter the MATLAB expression:

```
coder.newtype('uint16',[1 4],[0 1])
```

Note To specify that an input is a double-precision scalar, enter `0`.

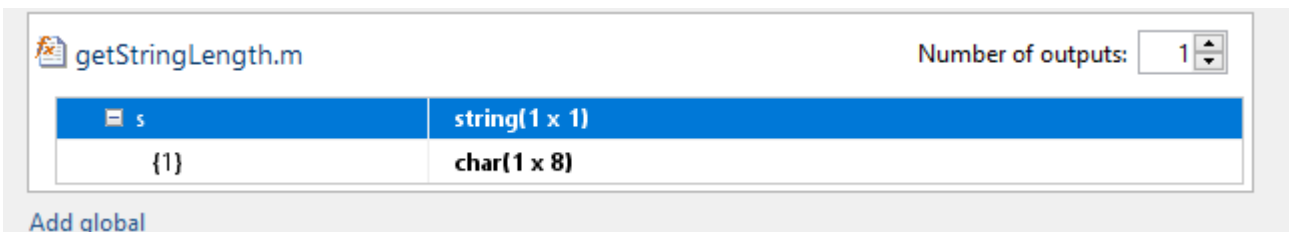
Specify a String Scalar Input Parameter by Example

This example shows how to specify a string scalar type by providing an example string.

- 1 On the **Define Input Types** page, click **Let me enter input or global types directly**.
- 2 Click the field to the right of the input parameter that you want to define.
- 3 Select **Define by Example**.
- 4 In the field to the right of the parameter, enter:

```
"mystring"
```

The input parameter is a 1-by-1 string array (string scalar) that contains a 1-by-8 character vector.



- 5 To make the string variable-size, click the second dimension.
 - To specify that the second dimension is unbounded, select `:Inf`.
 - To specify that the second dimension has an upper bound, enter the upper bound, for example 8. Then, select `:8`.

Specify a Structure Type Input Parameter by Example

This example shows how to specify a structure with two fields, a and b. The input type of a is scalar double. The input type of b is scalar char.

- 1 On the **Define Input Types** page, click **Let me enter input or global types directly**.
- 2 Click the field to the right of the input parameter that you want to define.
- 3 Select **Define by Example**.
- 4 In the field to the right of the parameter, enter:

```
struct('a', 1, 'b', 'x')
```

The type of the input parameter is `struct(1x1)`. The type of field a is `double(1x1)`. The type of field b is `char(1x1)`.

- 5 For an array of structures, to specify the size of each dimension, click the dimension and specify the size. For example, enter 4 for the first dimension.
- 6 To specify that the second dimension is variable size with an upper bound of 4, select `:4`. Alternatively, to specify that the second dimension is unbounded select `:Inf`.

Alternatively, specify the size of the array of structures in the `struct` function call. For example, `struct('a', { 1 2}, 'b', {'x', 'y'})` specifies a 1x2 array of structures with fields a and b. The type of field a is `double(1x1)`. The type of field b is `char(1x1)`.

To modify the type definition, see “Specify a Structure Input Parameter” (MATLAB Coder).

Specify a Cell Array Type Input Parameter by Example

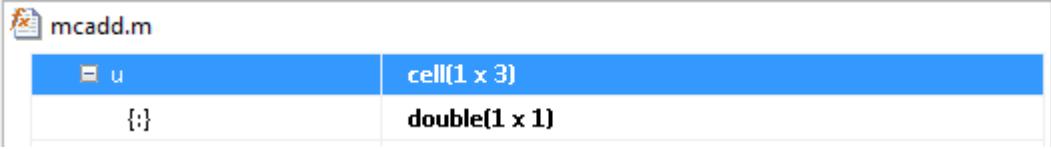
Note The Fixed-Point Converter app does not support cell arrays.

This example shows how to specify a cell array input by example. When you define a cell array by example, the app determines whether the cell array is homogeneous or heterogeneous. . If you want to control whether the cell array is homogeneous or heterogeneous, specify the cell array by type. See “Specify a Cell Array Input Parameter” (MATLAB Coder).

- 1 On the **Define Input Types** page, click **Let me enter input or global types directly**.
- 2 Click the field to the right of the input parameter that you want to define.
- 3 Select **Define by Example**.
- 4 In the field to the right of the parameter, enter an example cell array.
 - If all cell array elements have the same properties, the cell array is homogeneous. For example, enter:

```
{1 2 3}
```

The input is a 1x3 cell array. The type of each element is `double(1x1)`.



| mcadd.m | |
|---------|----------------------|
| u | cell(1 x 3) |
| {:} | double(1 x 1) |

The colon inside curly braces{: } indicates that all elements have the same properties.

- If elements of the cell array have different classes, the cell array is heterogeneous. For example, enter:

```
{'a', 1}
```

The input is a 1x2 cell array. For a heterogeneous cell array, the app lists each element. The type of the first element is `char(1x1)`. The type of the second element is `double(1x1)`.

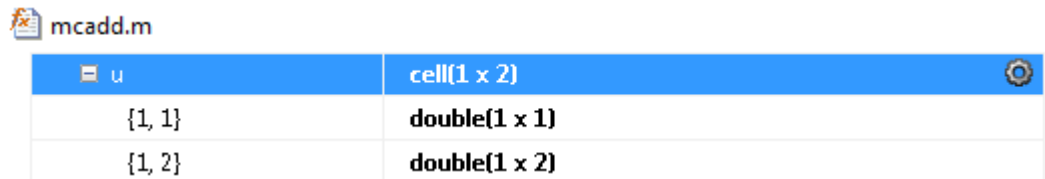


| mcadd.m | |
|---------|----------------------|
| u | cell(1 x 2) |
| {1, 1} | char(1 x 1) |
| {1, 2} | double(1 x 1) |

- For some example cell arrays, the classification as homogeneous or heterogeneous is ambiguous. For these cell arrays, the app uses heuristics to determine whether the cell array is homogeneous or heterogeneous. For example, for the example cell array, enter:

```
{1 [2 3]}
```

The elements have the same class, but different sizes. The app determines that the input is a 1x2 heterogeneous cell array. The type of the first element is `double(1x1)`. The type of the second element is `double(1x2)`.



| u | cell(1 x 2) |
|--------|---------------|
| {1, 1} | double(1 x 1) |
| {1, 2} | double(1 x 2) |

However, the example cell array, `{1 [2 3]}`, can also be a homogeneous cell array whose elements are `1x:2 double`. If you want this cell array to be homogeneous, do one of the following:

- Specify the cell array input by type. Specify that the input is a homogeneous cell array. Specify that the elements are `1x:2 double`. See “Specify a Cell Array Input Parameter” (MATLAB Coder).
- Right-click the variable. Select **Homogeneous**. Specify that the elements are `1x:2 double`.

If you use `coder.typeof` to specify that the example cell array is variable size, the app makes the cell array homogeneous. For example, for the example input, enter:

```
coder.typeof({1 [2 3]}, [1 3], [0 1])
```

The app determines that the input is a `1x:3 homogeneous cell array` whose elements are `1x:2 double`.

To modify the type definition, see “Specify a Cell Array Input Parameter” (MATLAB Coder).

Specify an Enumerated Type Input Parameter by Example

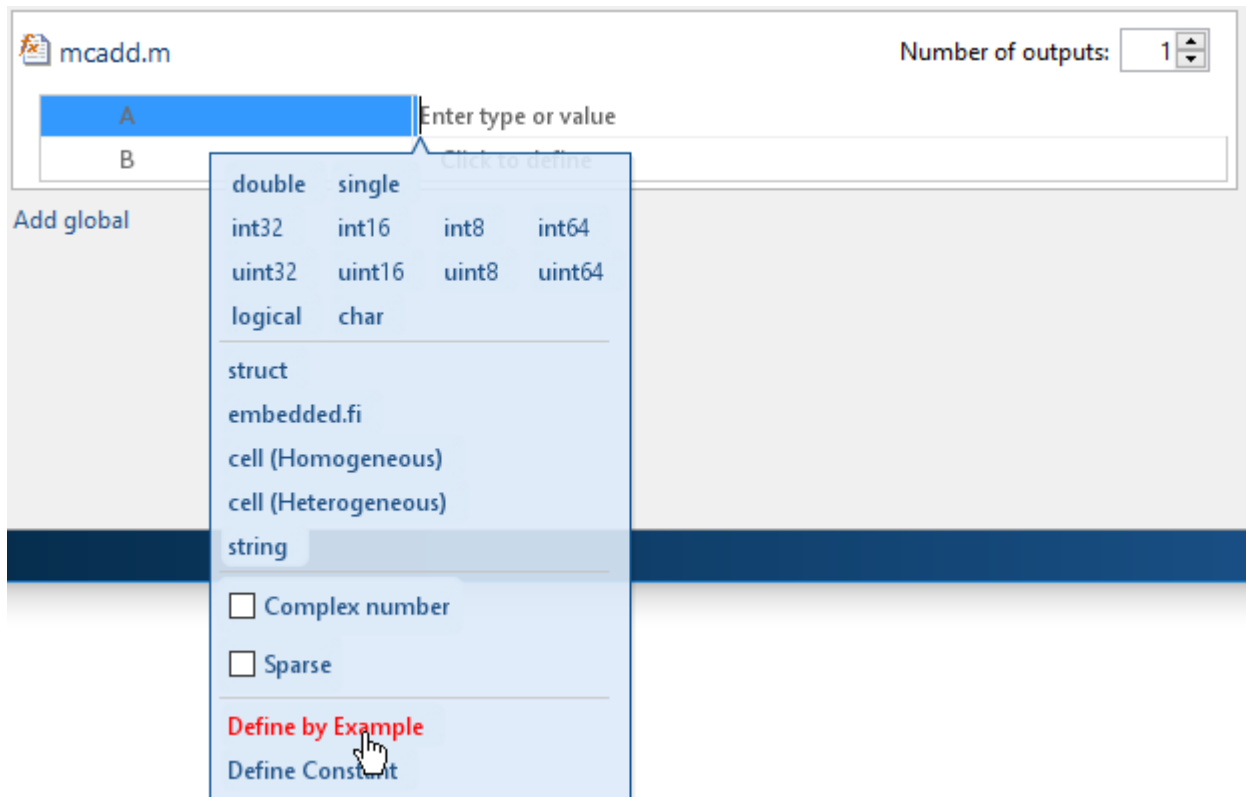
This example shows how to specify that an input uses the enumerated type `MyColors`.

Suppose that `MyColors.m` is on the MATLAB path.

```
classdef MyColors < int32
    enumeration
        green(1),
        red(2),
    end
end
```

To specify that an input has the enumerated type `MyColors`:

- 1 On the **Define Input Types** page, click **Let me enter input or global types directly**.
- 2 Click the field to the right of the input parameter that you want to define.



- 3 Select **Define by Example**.
- 4 In the field to the right of the parameter, enter the MATLAB expression:

```
MyColors.red
```

Specify a Fixed-Point Input Parameter by Example

To specify fixed-point inputs, Fixed-Point Designer software must be installed.

This example shows how to specify a signed fixed-point type with a word length of eight bits, and a fraction length of three bits.

- 1 On the **Define Input Types** page, click **Let me enter input or global types directly**.
- 2 Click the field to the right of the input parameter that you want to define.
- 3 Select **Define by Example**.
- 4 In the field to the right of the parameter, enter:

```
fi(10, 1, 8, 3)
```

The app sets the type of input `u` to `fi(1x1)`. By default, if you do not specify a local `fimath`, the app uses the default `fimath`. See “`fimath` for Sharing Arithmetic Rules” on page 3-17.

Optionally, modify the fixed-point properties or the size of the input. See “Specify a Fixed-Point Input Parameter” on page 8-68 and “Define or Edit Input Parameter Type by Using the App” on page 8-66.

Specify an Input from an Entry-Point Function Output Type

When generating code for multiple entry-point functions, you can use the output type from one entry-point function as the input type to another entry-point function. For more information, see “Pass an Entry-Point Function Output as an Input” (MATLAB Coder).

- 1 On the **Define Input Types** page, click **Let me enter input or global types directly**.
- 2 Click the field to the right of the input parameter that you want to define and select **Use Output**.

The screenshot shows the MATLAB Coder interface for defining input types. It features two function panels: 'makeSparse.m' and 'useSparse.m'. The 'useSparse.m' panel has an input parameter 'in' selected. A dropdown menu is open, listing various data types such as 'double', 'single', 'int32', 'int16', 'int8', 'int64', 'uint32', 'uint16', 'uint8', 'uint64', 'logical', 'char', 'struct', 'embedded.fi', 'cell (Homogeneous)', 'cell (Heterogeneous)', and 'string'. Below these are checkboxes for 'Complex number' and 'Sparse'. At the bottom of the menu, there are three options: 'Define by Example', 'Define Constant', and 'Use Output', which is highlighted in red and has a mouse cursor pointing to it.

- 3 Select the name of the entry-point function and the corresponding output parameter from which to define the input type.

Specify Global Variable Type and Initial Value Using the App

In this section...

“Why Specify a Type Definition for Global Variables?” on page 8-80

“Specify a Global Variable Type” on page 8-80

“Define a Global Variable by Example” on page 8-80

“Define or Edit Global Variable Type” on page 8-81

“Define Global Variable Initial Value” on page 8-81

“Define Global Variable Constant Value” on page 8-82

“Remove Global Variables” on page 8-82

Why Specify a Type Definition for Global Variables?

If you use global variables in your MATLAB algorithm, before building the project, you must add a global type definition and initial value for each global variable. If you do not initialize the global data, the app looks for the variable in the MATLAB global workspace. If the variable does not exist, the app generates an error.

For MEX functions, if you use global data, you must also specify whether to synchronize this data between MATLAB and the MEX function.

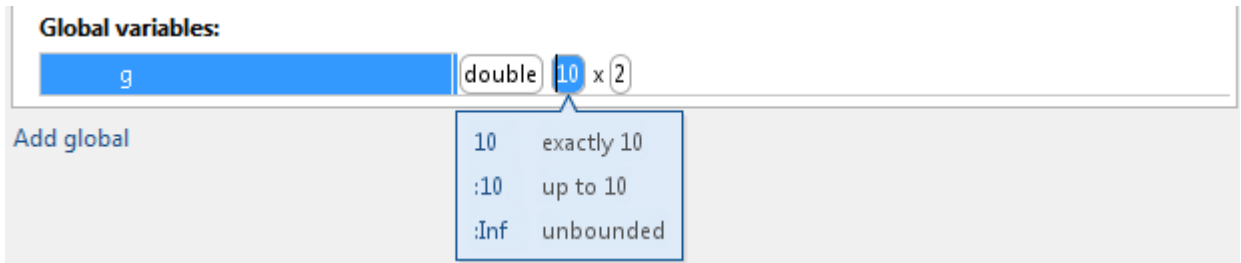
Specify a Global Variable Type

- 1 Specify the type of each global variable using one of the following methods:
 - Define by example on page 8-80
 - Define type on page 8-81
- 2 Define an initial value on page 8-81 for each global variable.

If you do not provide a type definition and initial value for a global variable, create a variable with the same name and suitable class, size, complexity, and value in the MATLAB workspace.

Define a Global Variable by Example

- 1 Click the field to the right of the global variable that you want to define.
- 2 Select **Define by Example**.
- 3 In the field to the right of the global name, enter a MATLAB expression that has the required class, size, and complexity. MATLAB Coder software uses the class, size, and complexity of the value of this expression as the type for the global variable.
- 4 Optionally, change the size of the global variable. Click the dimension that you want to change and enter the size, for example, 10.



You can specify:

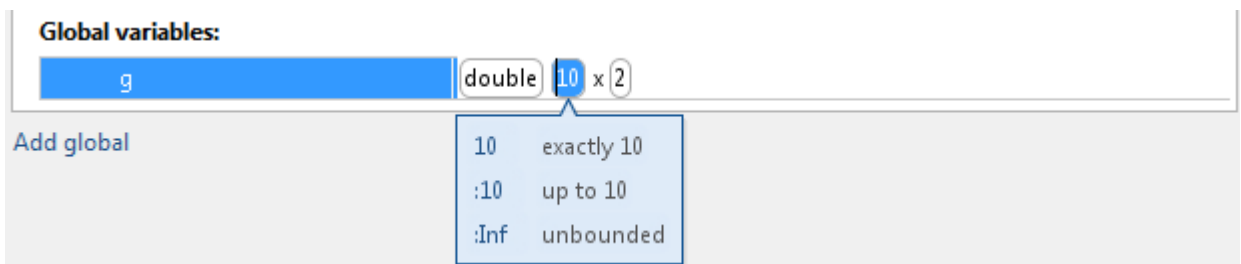
- Fixed size. In this example, select 10.
- Variable size, up to a specified limit, by using the : prefix. In this example, to specify that your input can vary in size up to 10, select :10.
- Unbounded variable size by selecting :Inf.

Define or Edit Global Variable Type

- 1 Click the field to the right of the global variable that you want to define.
- 2 Optionally, for numeric types, select **Complex** to make the parameter a complex type. By default, inputs are real.
- 3 Select the type for the global variable. For example, double.

By default, the global variable is a scalar.

- 4 Optionally, change the size of the global variable. Click the dimension that you want to change and enter the size, for example, 10.



You can specify:

- Fixed size. In this example, select 10.
- Variable size, up to a specified limit, by using the : prefix. In this example, to specify that your input can vary in size up to 10, select :10.
- Unbounded variable size by selecting :Inf.

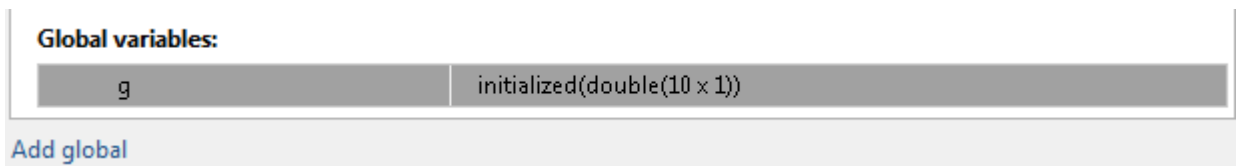
Define Global Variable Initial Value

- “Define Initial Value Before Defining Type” on page 8-82
- “Define Initial Value After Defining Type” on page 8-82

Define Initial Value Before Defining Type

- 1 Click the field to the right of the global variable.
- 2 Select **Define Initial Value**.
- 3 Enter a MATLAB expression. MATLAB Coder software uses the value of the specified MATLAB expression as the value of the global variable. Because you did not define the type of the global variable before you defined its initial value, MATLAB Coder uses the initial value type as the global variable type.

The project shows that the global variable is initialized.

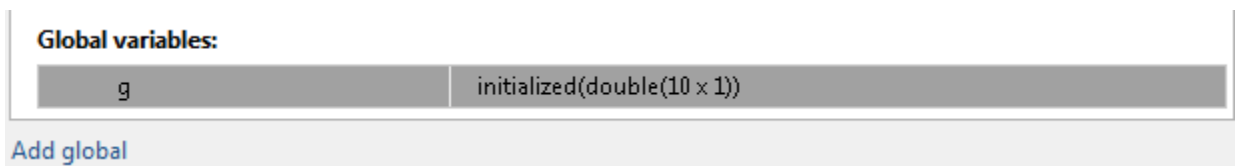


If you change the type of a global variable after defining its initial value, you must redefine the initial value.

Define Initial Value After Defining Type

- Click the type field of a predefined global variable.
- Select **Define Initial Value**.
- Enter a MATLAB expression. MATLAB Coder software uses the value of the specified MATLAB expression as the value of the global variable.

The project shows that the global variable is initialized.



Define Global Variable Constant Value

- 1 Click the field to the right of the global variable.
- 2 Select **Define Constant Value**.
- 3 In the field to the right of the global variable, enter a MATLAB expression.

Remove Global Variables

- 1 Right-click the global variable.
- 2 From the menu, select **Remove Global**.

Specify Properties of Entry-Point Function Inputs Using the App

Why Specify Input Properties?

Fixed-Point Designer must determine the properties of all variables in the MATLAB files. To infer variable properties in MATLAB files, Fixed-Point Designer must identify the properties of the inputs to the *primary* function, also known as the *top-level* or *entry-point* function. Therefore, if your primary function has inputs, you must specify the properties of these inputs to Fixed-Point Designer. If your primary function has no input parameters, you do not need to specify properties of inputs to local functions or external functions called by the primary function.

Unless you use the tilde (~) character to specify unused function inputs, you must specify the same number and order of inputs as the MATLAB function. If you use the tilde character, the inputs default to real, scalar doubles.

See Also

- “Properties to Specify” on page 31-2

Specify an Input Definition Using the App

Specify an input definition using one of the following methods:

- Autodefine Input Types on page 8-64
- Define Type on page 8-66
- Define by Example on page 8-73
- Define Constant on page 8-65

Detect Unexecuted and Constant-Folded Code

During the simulation of your test file, the Fixed-Point Converter app detects unexecuted code or code that is constant folded. Code that is not executed by the test bench may be unreachable code or dead code. The app uses the code coverage information when translating your code from floating-point MATLAB code to fixed-point MATLAB code. Reviewing code coverage results helps you to verify that your test file is exercising the algorithm adequately.

The app inserts inline comments in the fixed-point code to mark the unexecuted and untranslated regions. It includes the code coverage information in the generated fixed-point conversion HTML report. The app editor displays a color-coded bar to the left of the code. This table describes the color coding.

| Coverage Bar Color | Indicates |
|--------------------|---|
| Green | One of the following situations: <ul style="list-style-type: none"> The entry-point function executes multiple times and the code executes more than one time. The entry-point function executes one time and the code executes one time. Different shades of green indicate different ranges of line execution counts. The darkest shade of green indicates the highest range. |
| Orange | The entry-point function executes multiple times, but the code executes one time. |
| Red | Code does not execute. |

What Is Unexecuted Code?

Unexecuted code is code that is not executed by the test bench during simulation. Unexecuted code can result from these scenarios:

- Defensive code containing intended corner cases that are not reached
- Human error in the code, resulting in code that cannot be reached by any execution path, sometimes referred to as unreachable code or dead code
- Inadequate test bench range which does not provide inputs that execute all paths in the code
- Constant folding

Detect Unexecuted Code

This example shows how to detect code in your algorithm that is not executed by the test bench by using the Fixed-Point Converter .

- In a local writable folder, create the function `myFunction.m`.

```
function y = myFunction(u,v)
    %#codegen
    for i = 1:length(u)
        if u(i) > v(i)
            y=bar(u,v);
        end
    end
end
```

```

        else
            tmp = u;
            v = tmp;
            y = baz(u,v);
        end
    end
end

```

```

function y = bar(u,v)
    y = u+v;
end

```

```

function y = baz(u,v)
    y = u-v;
end

```

- 2 In the same folder, create a test file, myFunction_tb.

```

u = 1:100;
v = 101:200;

```

```

myFunction(u,v);

```

- 3 From the apps gallery, open the Fixed-Point Converter .
- 4 On the **Select Source Files** page, browse to the myFunction file, and click **Open**.
- 5 Click **Next**. On the **Define Input Types** page, browse to select the test file that you created, myFunction_tb. Click **Autodefine Input Types**.
- 6 Click **Next**. On the **Convert to Fixed-Point** page, click **Analyze** to simulate the entry-point functions, gather range information, and get proposed data types.

The color-coded bar on the left side of the edit window indicates whether the code executes. The code in the first condition of the if-statement does not execute during simulation because u is never greater than v . The `bar` function never executes because the if-statement never executes. These parts of the algorithm are marked with a red bar, indicating that they are not executed by the test bench.

Convert to Fixed Point

SETTINGS ▾ ANALYZE ▾ CONVERT TEST ▾

```

1 function y = myFunction(u,v)
2     %#codegen
3     for i = 1:length(u)
4         if u(i) > v(i)
5             y=bar(u,v);
6         else
7             tmp = u;
8             v = tmp;
9             y = baz(u,v);
10        end
11    end
12 end
13
14 function y = bar(u,v)
15     y = u+v;
16 end
17
18 function y = baz(u,v)
19     y = u-v;
20 end

```

Variables Function Replacements Output

| Variable | Type | Sim Min | Sim Max | Whole Number | Proposed Type |
|----------|--------------|---------|---------|--------------|----------------------|
| Input | | | | | |
| u | 1x100 double | 1 | 100 | Yes | numerictype(0, 7, 0) |
| v | 1x100 double | 1 | 200 | Yes | numerictype(0, 8, 0) |
| Output | | | | | |
| y | 1x100 double | 0 | 0 | Yes | numerictype(0, 1, 0) |
| Local | | | | | |
| i | double | 1 | 100 | Yes | numerictype(0, 7, 0) |

- 7 To apply the proposed data types to the function, click **Convert**.

The Fixed-Point Converter generates a fixed-point function, `myFunction_fixpt`. The generated fixed-point code contains comments around the pieces of code identified as not being executed by the test bench. The **Validation Results** pane proposes that you use a more thorough test bench.

The screenshot shows the MATLAB Fixed-Point Converter interface. The main window displays a code editor with the following code:

```

1 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
2 %
3 %           Generated by MATLAB 9.1 and Fixed-Point Designer 5.3
4 %
5 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
6 function y = myFunction_fixpt(u,v_1)
7     %#codegen
8     fm = get_fimath();
9     v = fi(v_1, 0, 8, 0, fm);
10
11     for i = 1:length(u)
12         if u(i) > v(i)
13             %F2F: No information found for converting the following block of code
14             %F2F: Start block
15             y=fi(bar_dead(u,v), 0, 1, 0, fm);
16             %F2F: End block
17         else
18             tmp = fi(u, 0, 7, 0, fm);
19             v = fi(tmp, 0, 8, 0, fm);
20             y = fi(baz(u,v), 0, 1, 0, fm);
21         end
22     end
23 end
24

```

Below the code editor, there is a table with the following data:

| Function | Line | Description |
|------------|------|---|
| myFunction | 5 | The expression 'y=bar(u,v)' was not executed during simulation. Consider using a more thorough testbench. |
| bar | 14 | The function 'bar' was not executed during simulation. Consider using a more thorough testbench. |

When the Fixed-Point Converter detects unexecuted code, consider editing your test file so that your algorithm is exercised over its full range. If your test file already reflects the full range of the input variables, consider editing your algorithm to eliminate the unreachable code.

- 8 Close the Fixed-Point Converter .

Fix Unexecuted Code

- 1 Edit the test file `myFunction_tb.m` to include a wider range of inputs.

```

u = 1:100;
v = -50:2:149;

```

```
myFunction(u,v);
```

- 2 Reopen the Fixed-Point Converter app.
- 3 Using the same function and the edited test file, go through the conversion process again.
- 4 After you click **Analyze**, this time the code coverage bar shows that all parts of the algorithm execute with the new test file input ranges.

Convert to Fixed Point

SETTINGS ▾ ANALYZE ▾ CONVERT TEST ▾

Code

```

1 function y = myFunction(u,v)
2     %#codegen
3     for i = 1:length(u)
4         if u(i) > v(i)
5             y=bar(u,v);
6         else
7             tmp = u;
8             v = tmp;
9             y = baz(u,v);
10        end
11    end
12 end
13
14 function y = bar(u,v)
15     y = u+v;
16 end
17
18 function y = baz(u,v)
19     y = u-v;
20 end

```

Files

- tion_fixpt.m
- tion_wrapper_fixpt.m
- tion_report.html
- tion_fixpt_args.mat
- tion_float_mex.mexw64
- tion_wrapper_fixpt_mex.m
- tion_fixpt_log.txt

| Variable | Type | Sim Min | Sim Max | Whole Number | Proposed Type |
|----------|--------------|---------|---------|--------------|------------------|
| Input | | | | | |
| u | 1x100 double | 1 | 100 | Yes | numeric(0, 7, 0) |
| v | 1x100 double | -50 | 200 | Yes | numeric(1, 9, 0) |
| Output | | | | | |
| y | 1x100 double | -49 | 248 | Yes | numeric(1, 9, 0) |
| Local | | | | | |
| i | double | 1 | 100 | Yes | numeric(0, 7, 0) |

Next >

To finish the conversion process and convert the function to fixed point, click **Convert**.

Automated Conversion Using Programmatic Workflow

- “Propose Data Types Based on Simulation Ranges” on page 9-2
- “Propose Data Types Based on Derived Ranges” on page 9-6
- “Detect Overflows” on page 9-12
- “Replace the exp Function with a Lookup Table” on page 9-16
- “Replace a Custom Function with a Lookup Table” on page 9-18
- “Visualize Differences Between Floating-Point and Fixed-Point Results” on page 9-20
- “Enable Plotting Using the Simulation Data Inspector” on page 9-25

Propose Data Types Based on Simulation Ranges

This example shows how to propose fixed-point data types based on simulation range data using the `fiaccel` function.

Prerequisites

To complete this example, you must install the following products:

- MATLAB
- Fixed-Point Designer
- C compiler

See https://www.mathworks.com/support/compilers/current_release/.

You can use `mex -setup` to change the default compiler. See “Change Default Compiler”.

Create a New Folder and Copy Relevant Files

- 1 In a local, writable folder, create a function `ex_2ndOrder_filter.m`.

```
function y = ex_2ndOrder_filter(x) %#codegen
    persistent z
    if isempty(z)
        z = zeros(2,1);
    end
    % [b,a] = butter(2, 0.25)
    b = [0.0976310729378175, 0.195262145875635, 0.0976310729378175];
    a = [1, -0.942809041582063, 0.333333333333333];

    y = zeros(size(x));
    for i = 1:length(x)
        y(i) = b(1)*x(i) + z(1);
        z(1) = b(2)*x(i) + z(2) - a(2) * y(i);
        z(2) = b(3)*x(i) - a(3) * y(i);
    end
end
```

- 2 Create a test file, `ex_2ndOrder_filter_test.m`, to exercise the `ex_2ndOrder_filter` algorithm.

It is best practice to create a separate test script to do all the pre- and post-processing such as loading inputs, setting up input values, calling the function under test, and outputting test results.

To cover the full intended operating range of the system, the test script runs the `ex_2ndOrder_filter` function with three input signals: chirp, step, and impulse. The script then plots the outputs.

```
% ex_2ndOrder_filter_test
%
% Define representative inputs
N = 256; % Number of points
t = linspace(0,1,N); % Time vector from 0 to 1 second
f1 = N/2; % Target frequency of chirp set to Nyquist
```

```

x_chirp = sin(pi*f1*t.^2); % Linear chirp from 0 to Fs/2 Hz in 1 second
x_step = ones(1,N);      % Step
x_impulse = zeros(1,N);  % Impulse
x_impulse(1) = 1;

% Run the function under test
x = [x_chirp;x_step;x_impulse];
y = zeros(size(x));
for i = 1:size(x,1)
    y(i,:) = ex_2ndOrder_filter(x(i,:));
end

% Plot the results
titles = {'Chirp','Step','Impulse'}
clf
for i = 1:size(x,1)
    subplot(size(x,1),1,i)
    plot(t,x(i,:),t,y(i,:))
    title(titles{i})
    legend('Input','Output')
end
xlabel('Time (s)')
figure(gcf)

disp('Test complete.')

```

| Type | Name | Description |
|---------------|---------------------------|---|
| Function code | ex_2ndOrder_filter.m | Entry-point MATLAB function |
| Test file | ex_2ndOrder_filter_test.m | MATLAB script that tests ex_2ndOrder_filter.m |

Set Up the Fixed-Point Configuration Object

Create a fixed-point configuration object and configure the test file name.

```

cfg = coder.config('fixpt');
cfg.TestBenchName = 'ex_2ndOrder_filter_test';

```

Collect Simulation Ranges and Generate Fixed-Point Code

Use the `fiaccel` function to convert the floating-point MATLAB function, `ex_2ndOrder_filter`, to fixed-point MATLAB code. Set the default word length for the fixed-point data types to 16.

```

cfg.ComputeSimulationRanges = true;
cfg.DefaultWordLength = 16;

```

```

% Derive ranges and generate fixed-point code
fiaccel -float2fixed cfg ex_2ndOrder_filter

```

`fiaccel` analyzes the floating-point code. Because you did not specify the input types for the `ex_2ndOrder_filter` function, the conversion process infers types by simulating the test file. The conversion process then derives ranges for variables in the algorithm. It uses these derived ranges to propose fixed-point types for these variables. When the conversion is complete, it generates a type proposal report.

View Range Information

Click the link to the type proposal report for the `ex_2ndOrder_filter` function, `ex_2ndOrder_filter_report.html`.

The report opens in a web browser.

Fixed-Point Report `ex_2ndOrder_filter`

| Simulation Coverage | Code |
|---------------------|---|
| 100% | <code>function y = ex_2ndOrder_filter(x) %#codegen</code> |
| Once | <code>persistent z</code> |
| | <code>if isempty(z)</code> |
| | <code>z = zeros(2,1);</code> |
| | <code>end</code> |
| 100% | <code>% [b,a] = butter(2, 0.25)</code> |
| | <code>b = [0.0976310729378175, 0.195262145875635, 0.0976310729378175];</code> |
| | <code>a = [1, -0.942809041582063, 0.3333333333333333];</code> |
| | <code>y = zeros(size(x));</code> |
| | <code>for i=1:length(x)</code> |
| | <code>y(i) = b(1)*x(i) + z(1);</code> |
| | <code>z(1) = b(2)*x(i) + z(2) - a(2) * y(i);</code> |
| | <code>z(2) = b(3)*x(i) - a(3) * y(i);</code> |
| | <code>end</code> |
| | <code>end</code> |

| Variable Name | Type | Sim Min | Sim Max | Static Min | Static Max | Whole Number | ProposedType (Best For WL = 16) |
|---------------|----------------|---------------------|--------------------|------------|------------|--------------|---------------------------------|
| a | double 1 x 3 | -0.942809041582063 | 1 | | | No | numericType(1, 16, 14) |
| b | double 1 x 3 | 0.0976310729378175 | 0.195262145875635 | | | No | numericType(0, 16, 18) |
| i | double | 1 | 256 | | | Yes | numericType(0, 9, 0) |
| x | double 1 x 256 | -0.9999756307053946 | 1 | | | No | numericType(1, 16, 14) |
| y | double 1 x 256 | -0.9696817930434206 | 1.0553496057969345 | | | No | numericType(1, 16, 14) |
| z | double 2 x 1 | -0.8907046852192462 | 0.957718532859117 | | | No | numericType(1, 16, 15) |

View Generated Fixed-Point MATLAB Code

`fiaccel` generates a fixed-point version of the `ex_2ndOrder_filter.m` function, `ex_2ndOrder_filter_fixpt.m`, and a wrapper function that calls `ex_2ndOrder_filter_fixpt`. These files are generated in the `codegen\ex_2ndOrder_filter\fixpt` folder in your local working folder.

```
function y = ex_2ndOrder_filter_fixpt(x) %#codegen
    fm = get_fimath();

    persistent z
    if isempty(z)
        z = fi(zeros(2,1),1,16,15, fm);
    end
    % [b,a] = butter(2,0.25)
    b = fi([0.0976310729378175, 0.195262145875635, ...
        0.0976310729378175],0,16,18, fm);
```

```
a = fi([1,-0.942809041582063,...
0.3333333333333333],1,16,14, fm);

y = fi(zeros(size(x)),1,16,14, fm);
for i=1:length(x)
    y(i) = b(1)*x(i) + z(1);
    z(1) = fi_signed(b(2)*x(i) + z(2)) - a(2) * y(i);
    z(2) = fi_signed(b(3)*x(i))          - a(3) * y(i);
end
end

function y = fi_signed(a)
    coder.inline( 'always' );
    if isfi(a) && ~(issigned(a))
        nt = numerictype(a);
        new_nt = numerictype(1,nt.WordLength + 1,nt.FractionLength);
        y = fi(a,new_nt, fimath(a));
    else
        y = a;
    end
end

function fm = get_fimath()
    fm = fimath('RoundingMethod','Floor',...
        'OverflowAction','Wrap',...
        'ProductMode','FullPrecision',...
        'MaxProductWordLength',128,...
        'SumMode','FullPrecision',...
        'MaxSumWordLength',128);
end
```

Propose Data Types Based on Derived Ranges

This example shows how to propose fixed-point data types based on static ranges using the `fiaccel` function. The advantage of proposing data types based on derived ranges is that you do not have to provide test files that exercise your algorithm over its full operating range. Running such test files often takes a very long time so you can save time by deriving ranges instead.

Note Derived range analysis is not supported for non-scalar variables.

Prerequisites

To complete this example, you must install the following products:

- MATLAB
- Fixed-Point Designer
- C compiler

See https://www.mathworks.com/support/compilers/current_release/.

You can use `mex -setup` to change the default compiler. See “Change Default Compiler”.

Create a New Folder and Copy Relevant Files

- 1 In a local, writable folder, create a function `dti.m`.

The `dti` function implements a Discrete Time Integrator in MATLAB.

```
function [y, clip_status] = dti(u_in) %#codegen
% Discrete Time Integrator in MATLAB
%
% Forward Euler method, also known as Forward
% Rectangular, or left-hand approximation.
% The resulting expression for the output of
% the block at step 'n' is
% y(n) = y(n-1) + K * u(n-1)
%
init_val = 1;
gain_val = 1;
limit_upper = 500;
limit_lower = -500;

% Variable to hold state between
% consecutive calls to this block
persistent u_state;
if isempty(u_state)
    u_state = init_val+1;
end

% Compute Output
if (u_state > limit_upper)
    y = limit_upper;
    clip_status = -2;
elseif (u_state >= limit_upper)
```

```

        y = limit_upper;
        clip_status = -1;
elseif (u_state < limit_lower)
    y = limit_lower;
    clip_status = 2;
elseif (u_state <= limit_lower)
    y = limit_lower;
    clip_status = 1;
else
    y = u_state;
    clip_status = 0;
end

```

```

% Update State
tprod = gain_val * u_in;
u_state = y + tprod;

```

- 2 Create a test file, `dti_test.m`, to exercise the `dti` algorithm.

The test script runs the `dti` function with a sine wave input. The script then plots the input and output signals.

```

% dti_test
% cleanup
clear dti

% input signal
x_in = sin(2.*pi.*(0:0.001:2)).';

pause(10);

len = length(x_in);
y_out = zeros(1,len);
is_clipped_out = zeros(1,len);

for ii=1:len
    data = x_in(ii);
    % call to the dti function
    init_val = 0;
    gain_val = 1;
    upper_limit = 500;
    lower_limit = -500;

    % call to the design that does DTI
    [y_out(ii), is_clipped_out(ii)] = dti(data);
end

figure('Name', [mfilename, '_plot']);
subplot(2,1,1)
plot(1:len,x_in)
xlabel('Time')
ylabel('Amplitude')
title('Input Signal (Sin)')

subplot(2,1,2)
plot(1:len,y_out)
xlabel('Time')

```

```

ylabel('Amplitude')
title('Output Signal (DTI)')

disp('Test complete.');
```

It is a best practice is to create a separate test script to do pre- and post-processing, such as:

- Loading inputs.
- Setting up input values.
- Outputting test results.

| Type | Name | Description |
|---------------|------------|--------------------------------|
| Function code | dti.m | Entry-point MATLAB function |
| Test file | dti_test.m | MATLAB script that tests dti.m |

Set Up the Fixed-Point Configuration Object

Create a fixed-point configuration object and configure the test file name.

```

fixptcfg = coder.config('fixpt');
fixptcfg.TestBenchName = 'dti_test';
```

Specify Design Ranges

Specify design range information for the dti function input parameter u_in.

```

fixptcfg.addDesignRangeSpecification('dti', 'u_in', -1.0, 1.0)
```

Enable Plotting Using the Simulation Data Inspector

Select to run the test file to verify the generated fixed-point MATLAB code. Log inputs and outputs for comparison plotting and select to use the Simulation Data Inspector to plot the results.

```

fixptcfg.TestNumerics = true;
fixptcfg.LogIOForComparisonPlotting = true;
fixptcfg.PlotWithSimulationDataInspector = true;
```

Derive Ranges and Generate Fixed-Point Code

Use the fiaccel function to convert the floating-point MATLAB function, dti, to fixed-point MATLAB code. Set the default word length for the fixed-point data types to 16.

```

fixptcfg.ComputeDerivedRanges = true;
fixptcfg.ComputeSimulationRanges = false;
fixptcfg.DefaultWordLength = 16;

% Derive ranges and generate fixed-point code
fiaccel -float2fixed fixptcfg dti
```

fiaccel analyzes the floating-point code. Because you did not specify the input types for the dti function, the conversion process infers types by simulating the test file. The conversion process then derives ranges for variables in the algorithm. It uses these derived ranges to propose fixed-point types for these variables. When the conversion is complete, it generates a type proposal report.

View Derived Range Information

Click the link to the type proposal report for the `dti` function, `dti_report.html`.

The report opens in a web browser.

Fixed Point Report dti

```
function [y,clip_status] = dti(u_in) %#codegen
% Discrete Time Integrator in MATLAB
%
% Forward Euler method, also known as Forward Rectangular, or left-hand
% approximation. The resulting expression for the output of the block at
% step 'n' is y(n) = y(n-1) + K * u(n-1)
%
init_val = 1;
gain_val = 1;
limit_upper = 500;
limit_lower = -500;
% variable to hold state between consecutive calls to this block
persistent u_state
if isempty( u_state )
    u_state = init_val + 1;
end
% Compute Output
if (u_state>limit_upper)
    y = limit_upper;
    clip_status = -2;
elseif (u_state>=limit_upper)
    y = limit_upper;
    clip_status = -1;
elseif (u_state
```

| Variable Name | Type | Sim Min | Sim Max | Static Min | Static Max | Whole Number | ProposedType (Best For WL = 16) |
|---------------|--------|---------|---------|------------|------------|--------------|------------------------------------|
| clip_status | double | | | -2 | 2 | No | numerictype(1, 16, 13) |
| gain_val | double | | | 1 | 1 | Yes | numerictype(0, 1, 0) |
| init_val | double | | | 1 | 1 | Yes | numerictype(0, 1, 0) |
| limit_lower | double | | | -500 | -500 | Yes | numerictype(1, 10, 0) |
| limit_upper | double | | | 500 | 500 | Yes | numerictype(0, 9, 0) |
| tprod | double | | | -1 | 1 | No | numerictype(1, 16, 14) |
| u_in | double | | | -1 | 1 | No | numerictype(1, 16, 14) |
| u_state | double | | | -501 | 501 | No | numerictype(1, 16, 6) |
| y | double | | | -500 | 500 | No | numerictype(1, 16, 6) |

View Generated Fixed-Point MATLAB Code

`fiaccl` generates a fixed-point version of the `dti` function, `dti_fxpt.m`, and a wrapper function that calls `dti_fxpt`. These files are generated in the `codegen\dti\fixpt` folder in your local working folder.

```
function [y, clip_status] = dti_fxpt(u_in) %#codegen
% Discrete Time Integrator in MATLAB
%
% Forward Euler method, also known as
% Forward Rectangular, or left-hand
% approximation. The resulting expression
% for the output of the block at
% step 'n' is y(n) = y(n-1) + K * u(n-1)
```

```

%
fm = get_fimath();

init_val = fi(1, 0, 1, 0, fm);
gain_val = fi(1, 0, 1, 0, fm);
limit_upper = fi(500, 0, 9, 0, fm);
limit_lower = fi(-500, 1, 10, 0, fm);

% variable to hold state between
% consecutive calls to this block
persistent u_state;
if isempty(u_state)
    u_state = fi(init_val+fi(1,0,1,0, fm),1,16,6, fm);
end

% Compute Output
if (u_state > limit_upper)
    y = fi(limit_upper, 1, 16, 6, fm);
    clip_status = fi(-2, 1, 16, 13, fm);
elseif (u_state >= limit_upper)
    y = fi(limit_upper, 1, 16, 6, fm);
    clip_status = fi(-1, 1, 16, 13, fm);
elseif (u_state < limit_lower)
    y = fi(limit_lower, 1, 16, 6, fm);
    clip_status = fi(2, 1, 16, 13, fm);
elseif (u_state <= limit_lower)
    y = fi(limit_lower, 1, 16, 6, fm);
    clip_status = fi(1, 1, 16, 13, fm);
else
    y = fi(u_state, 1, 16, 6, fm);
    clip_status = fi(0, 1, 16, 13, fm);
end

% Update State
tprod = fi(gain_val * u_in, 1, 16, 14, fm);
u_state(:) = y + tprod;
end

function fm = get_fimath()
    fm = fimath('RoundingMethod','Floor',...
        'OverflowAction','Wrap',...
        'ProductMode','FullPrecision',...
        'MaxProductWordLength',128,...
        'SumMode','FullPrecision',...
        'MaxSumWordLength',128);
end

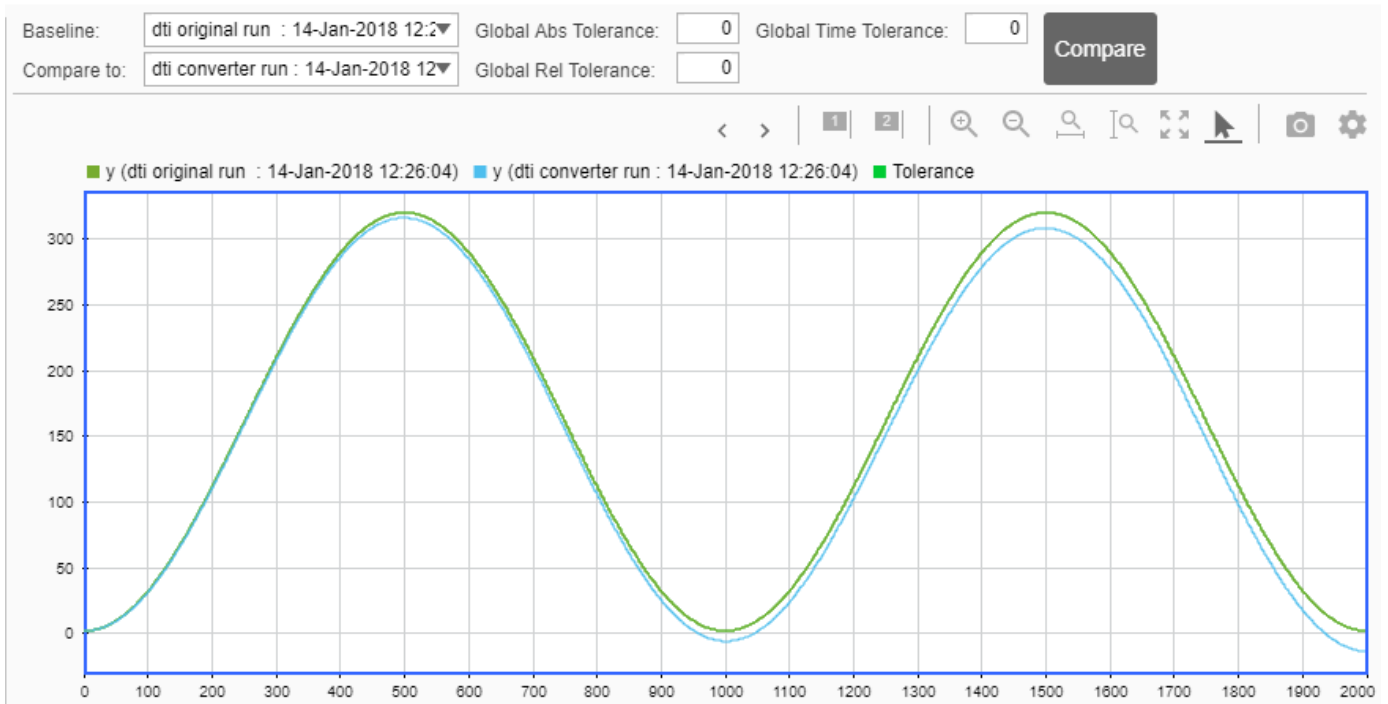
```

Compare Floating-Point and Fixed-Point Runs

Because you selected to log inputs and outputs for comparison plots and to use the Simulation Data Inspector for these plots, the Simulation Data Inspector opens.

You can use the Simulation Data Inspector to view floating-point and fixed-point run information and compare results. For example, to compare the floating-point and fixed-point values for the output *y*, on the **Compare** tab, select *y*, and then click **Compare Runs**.

The Simulation Data Inspector displays a plot of the baseline floating-point run against the fixed-point run and the difference between them.



Detect Overflows

This example shows how to detect overflows using the `fiaccel` function. At the numerical testing stage in the conversion process, the tool simulates the fixed-point code using scaled doubles. It then reports which expressions in the generated code produce values that would overflow the fixed-point data type.

Prerequisites

To complete this example, you must install the following products:

- MATLAB
- Fixed-Point Designer
- C compiler

See https://www.mathworks.com/support/compilers/current_release/.

You can use `mex -setup` to change the default compiler. See “Change Default Compiler”.

Create a New Folder and Copy Relevant Files

- 1 In a local, writable folder, create a function `overflow.m`.

```
function y = overflow(b,x,reset)
    if nargin<3, reset = true; end
    persistent z p
    if isempty(z) || reset
        p = 0;
        z = zeros(size(b));
    end
    [y,z,p] = fir_filter(b,x,z,p);
end
function [y,z,p] = fir_filter(b,x,z,p)
    y = zeros(size(x));
    nx = length(x);
    nb = length(b);
    for n = 1:nx
        p=p+1; if p>nb, p=1; end
        z(p) = x(n);
        acc = 0;
        k = p;
        for j=1:nb
            acc = acc + b(j)*z(k);
            k=k-1; if k<1, k=nb; end
        end
        y(n) = acc;
    end
end
```

- 2 Create a test file, `overflow_test.m`, to exercise the overflow algorithm.

```
function overflow_test
    % The filter coefficients were computed using
    % the FIR1 function from Signal Processing Toolbox.
    % b = fir1(11,0.25);
    b = [-0.004465461051254
        -0.004324228005260
```

```

+0.012676739550326
+0.074351188907780
+0.172173206073645
+0.249588554524763
+0.249588554524763
+0.172173206073645
+0.074351188907780
+0.012676739550326
-0.004324228005260
-0.004465461051254]';

% Input signal
nx = 256;
t = linspace(0,10*pi,nx)';

% Impulse
x_impulse = zeros(nx,1); x_impulse(1) = 1;

% Max Gain
% The maximum gain of a filter will occur
% when the inputs line up with the
% signs of the filter's impulse response.
x_max_gain = sign(b)';
x_max_gain = repmat(x_max_gain,ceil(nx/length(b)),1);
x_max_gain = x_max_gain(1:nx);

% Sums of sines
f0=0.1; f1=2;
x_sines = sin(2*pi*t*f0) + 0.1*sin(2*pi*t*f1);

% Chirp
f_chirp = 1/16; % Target frequency
x_chirp = sin(pi*f_chirp*t.^2); % Linear chirp

x = [x_impulse,x_max_gain,x_sines,x_chirp];
titles = {'Impulse','Max gain','Sum of sines','Chirp'};
y = zeros(size(x));

for i=1:size(x,2)
    reset = true;
    y(:,i) = overflow(b,x(:,i),reset);
end

test_plot(1,titles,t,x,y)

end
function test_plot(fig,titles,t,x,y1)
figure(fig)
clf
sub_plot = 1;
font_size = 10;
for i=1:size(x,2)
    subplot(4,1,sub_plot)
    sub_plot = sub_plot+1;
    plot(t,x(:,i),'c',t,y1(:,i),'k')
    axis('tight')
    xlabel('t','FontSize',font_size);
end

```

```

        title(titles{i}, 'FontSize', font_size);
        ax = gca;
        ax.FontSize = 10;
    end
    figure(gcf)
end

```

It is best practice to create a separate test script to do all the pre- and post-processing such as loading inputs, setting up input values, calling the function under test, and outputting test results.

| Type | Name | Description |
|---------------|-----------------|-------------------------------------|
| Function code | overflow.m | Entry-point MATLAB function |
| Test file | overflow_test.m | MATLAB script that tests overflow.m |

Set Up Configuration Object

- 1 Create a `coder.FixptConfig` object, `fixptcfg`, with default settings.

```
fixptcfg = coder.config('fixpt');
```

- 2 Set the test bench name. In this example, the test bench function name is `overflow_test`.

```
fixptcfg.TestBenchName = 'overflow_test';
```

- 3 Set the default word length to 16.

```
fixptcfg.DefaultWordLength = 16;
```

Enable Overflow Detection

```
fixptcfg.TestNumerics = true;
fixptcfg.DetectFixptOverflows = true;
```

Set fimath Options

Set the `fimath` `Product` mode and `Sum` mode to `KeepLSB`. These settings models the behavior of integer operations in the C language.

```
fixptcfg.fimath = ...
'fimath( 'RoundingMethod', 'Floor', 'OverflowAction', 'Wrap', ...
'ProductMode', 'KeepLSB', 'SumMode', 'KeepLSB');
```

Convert to Fixed Point

Convert the floating-point MATLAB function, `overflow`, to fixed-point MATLAB code. You do not need to specify input types for the `fiaccl` command because it infers the types from the test file.

```
fiaccl -float2fixed fixptcfg overflow
```

The numerics testing phase reports an overflow.

```
Overflow error in expression 'acc + b( j ) * z( k )'.
Percentage of Current Range = 104%.
```

Review Results

Determine if the addition or the multiplication in this expression overflowed. Set the `fimath` `ProductMode` to `FullPrecision` so that the multiplication will not overflow, and then run the `fiaccl` command again.

```
fixptcfg.fimath = ...  
'fimath('RoundingMethod','Floor','OverflowAction','Wrap',...  
        'ProductMode','FullPrecision','SumMode','KeepLSB');  
fiaccel -float2fixed fixptcfg overflow
```

The numerics testing phase still reports an overflow, indicating that it is the addition in the expression that is overflowing.

Replace the exp Function with a Lookup Table

This example shows how to replace the `exp` function with a lookup table approximation in the generated fixed-point code using the `fiaccel` function.

Prerequisites

To complete this example, you must install the following products:

- MATLAB
- Fixed-Point Designer
- C compiler

See https://www.mathworks.com/support/compilers/current_release/.

You can use `mex -setup` to change the default compiler. See “Change Default Compiler”.

Create Algorithm and Test Files

- 1 Create a MATLAB function, `my_fcn.m`, that calls the `exp` function.

```
function y = my_fcn(x)
    y = exp(x);
end
```

- 2 Create a test file, `my_fcn_test.m`, that uses `my_fcn.m`.

```
close all

x = linspace(-10,10,1e3);
for itr = 1e3:-1:1
    y(itr) = my_fcn( x(itr) );
end
plot( x, y );
```

Configure Approximation

Create a function replacement configuration object to approximate the `exp` function, using the default settings of linear interpolation and 1000 points in the lookup table.

```
q = coder.approximation('exp');
```

Set Up Configuration Object

Create a `coder.FixptConfig` object, `fixptcfg`. Specify the test file name and enable numerics testing. Associate the function replacement configuration object with the fixed-point configuration object.

```
fixptcfg = coder.config('fixpt');
fixptcfg.TestBenchName = 'my_fcn_test';
fixptcfg.TestNumerics = true;
fixptcfg.DefaultWordLength = 16;
fixptcfg.addApproximation(q);
```

Convert to Fixed Point

Generate fixed-point MATLAB code.


```
fiaccel -float2fixed fixptcfg my_fcn
```

View Generated Fixed-Point Code

To view the generated fixed-point code, click the link to `my_fcn_fixpt`.

The generated code contains a lookup table approximation, `replacement_exp`, for the `exp` function. The fixed-point conversion process infers the ranges for the function and then uses an interpolated lookup table to replace the function. By default, the lookup table uses linear interpolation, 1000 points, and the minimum and maximum values detected by running the test file.

The generated fixed-point function, `my_fcn_fixpt`, calls this approximation instead of calling `exp`.

```
function y = my_fcn_fixpt(x)
    fm = get_fimath();

    y = fi(replacement_exp(x), 0, 16, 1, fm);
end
```

You can now test the generated fixed-point code and compare the results against the original MATLAB function. If the behavior of the generated fixed-point code does not match the behavior of the original code closely enough, modify the interpolation method or number of points used in the lookup table and then regenerate code.

See Also

More About

- “Replacing Functions Using Lookup Table Approximations” on page 7-50

Replace a Custom Function with a Lookup Table

This example shows how to replace a custom function with a lookup table approximation function using the `fiaccl` function.

Prerequisites

To complete this example, you must install the following products:

- MATLAB
- Fixed-Point Designer
- C compiler

See https://www.mathworks.com/support/compilers/current_release/.

You can use `mex -setup` to change the default compiler. See “Change Default Compiler”.

Create a MATLAB function, `custom_fcn.m`. This is the function that you want to replace.

```
function y = custom_fcn(x)
    y = 1./(1+exp(-x));
end
```

Create a wrapper function that calls `custom_fcn.m`.

```
function y = call_custom_fcn(x)
    y = custom_fcn(x);
end
```

Create a test file, `custom_test.m`, that uses `call_custom_fcn.m`.

```
close all

x = linspace(-10,10,1e3);
for itr = 1e3:-1:1
    y(itr) = call_custom_fcn( x(itr) );
end
plot( x, y );
```

Create a function replacement configuration object to approximate `custom_fcn`. Specify the function handle of the custom function and set the number of points to use in the lookup table to 50.

```
q = coder.approximation('Function','custom_fcn',...
    'CandidateFunction',@custom_fcn,...
    'NumberOfPoints',50);
```

Create a `coder.FixptConfig` object, `fixptcfg`. Specify the test file name and enable numerics testing. Associate the function replacement configuration object with the fixed-point configuration object.

```
fixptcfg = coder.config('fixpt');
fixptcfg.TestBenchName = 'custom_test';
fixptcfg.TestNumerics = true;
fixptcfg.addApproximation(q);
```

Generate fixed-point MATLAB code.

```
fiaccel -float2fixed fixptcfg call_custom_fcn
```

fiaccel generates fixed-point MATLAB code in `call_custom_fcn_fixpt.m`.

To view the generated fixed-point code, click the link to `call_custom_fcn_fixpt`.

The generated code contains a lookup table approximation, `replacement_custom_fcn`, for the `custom_fcn` function. The fixed-point conversion process infers the ranges for the function and then uses an interpolated lookup table to replace the function. The lookup table uses 50 points as specified. By default, it uses linear interpolation and the minimum and maximum values detected by running the test file.

The generated fixed-point function, `call_custom_fcn_fixpt`, calls this approximation instead of calling `custom_fcn`.

```
function y = call_custom_fcn_fixpt(x)
    fm = get_fimath();

    y = fi(replacement_custom_fcn(x), 0, 14, 14, fm);
end
```

You can now test the generated fixed-point code and compare the results against the original MATLAB function. If the behavior of the generated fixed-point code does not match the behavior of the original code closely enough, modify the interpolation method or number of points used in the lookup table and then regenerate code.

See Also

More About

- “Replacing Functions Using Lookup Table Approximations” on page 7-50

Visualize Differences Between Floating-Point and Fixed-Point Results

This example shows how to configure the `fiaccel` function to use a custom plot function to compare the behavior of the generated fixed-point code against the behavior of the original floating-point MATLAB code.

By default, when the `LogIOForComparisonPlotting` option is enabled, the conversion process uses a time series based plotting function to show the floating-point and fixed-point results and the difference between them. However, during fixed-point conversion you might want to visualize the numerical differences in a view that is more suitable for your application domain. This example shows how to customize plotting and produce scatter plots at the test numerics step of the fixed-point conversion.

Copy Relevant Files

Copy the `myFilter.m`, `myFilterTest.m`, `plotDiff.m`, and `filterData.mat` files to a local working folder.

Prerequisites

To complete this example, you must install the following products:

- MATLAB
- Fixed-Point Designer
- C compiler

See https://www.mathworks.com/support/compilers/current_release/.

You can use `mex -setup` to change the default compiler. See “Change Default Compiler”.

Inspect Example Files

It is best practice to create a separate test script to do all the pre- and post-processing such as loading inputs, setting up input values, calling the function under test, and outputting test results.

| Type | Name | Description |
|-------------------|-----------------------------|--|
| Function code | <code>myFilter.m</code> | Entry-point MATLAB function |
| Test file | <code>myFilterTest.m</code> | MATLAB script that tests <code>myFilter.m</code> |
| Plotting function | <code>plotDiff.m</code> | Custom plot function |
| MAT-file | <code>filterData.mat</code> | Data to filter. |

The `myFilter` Function

```
function [y, ho] = myFilter(in)
persistent b h;
```

```

if isempty(b)
    b = complex(zeros(1,16));
    h = complex(zeros(1,16));
    h(8) = 1;
end

b = [in, b(1:end-1)];
y = b*h.';

errf = 1-sqrt(real(y)*real(y) + imag(y)*imag(y));
update = 0.001*conj(b)*y*errf;

h = h + update;
h(8) = 1;
ho = h;

end

```

The myFilterTest File

```

% load data
data = load('filterData.mat');
d = data.symbols;

for idx = 1:4000
    y = myFilter(d(idx));
end

```

The plotDiff Function

```

% varInfo - structure with information about
% the variable. It has the following fields
%     i) name
%     ii) functionName
% floatVals - cell array of logged original values
% for the 'varInfo.name' variable
% fixedVals - cell array of logged values for
% the 'varInfo.name' variable after Fixed-Point conversion.
function plotDiff(varInfo, floatVals, fixedVals)
    varName = varInfo.name;
    fcnName = varInfo.functionName;

    % escape the '_'s because plot titles treat these as subscripts
    escapedVarName = regexp(varName, '_', '\\_');
    escapedFcnName = regexp(fcnName, '_', '\\_');

    % flatten the values
    flatFloatVals = floatVals(1:end);
    flatFixedVals = fixedVals(1:end);

    % build Titles
    floatTitle = [escapedFcnName ' > ' 'float : ' escapedVarName];
    fixedTitle = [escapedFcnName ' > ' 'fixed : ' escapedVarName];

    data = load('filterData.mat');

    switch varName
        case 'y'

```

```

x_vec = data.symbols;

figure('Name','Comparison plot','NumberTitle','off');

% plot floating point values
y_vec = flatFloatVals;
subplot(1, 2, 1);
plotScatter(x_vec, y_vec, 100, floatTitle);

% plot fixed point values
y_vec = flatFixedVals;
subplot(1, 2, 2);
plotScatter(x_vec, y_vec, 100, fixedTitle);

otherwise
    % Plot only output 'y' for this example, skip the rest
end

end

function plotScatter(x_vec, y_vec, n, figTitle)
    % plot the last n samples
    x_plot = x_vec(end-n+1:end);
    y_plot = y_vec(end-n+1:end);

    hold on
    scatter(real(x_plot),imag(x_plot), 'bo');

    hold on
    scatter(real(y_plot),imag(y_plot), 'rx');

    title(figTitle);
end

```

Set Up Configuration Object

- 1 Create a `coder.FixptConfig` object.

```
fxptcfg = coder.config('fixpt');
```

- 2 Specify the test file name and custom plot function name. Enable logging and numerics testing.

```

fxptcfg.TestBenchName = 'myFilterTest';
fxptcfg.PlotFunction = 'plotDiff';
fxptcfg.TestNumerics = true;
fxptcfg.LogIOForComparisonPlotting = true;
fxptcfg.DefaultWordLength = 16;

```

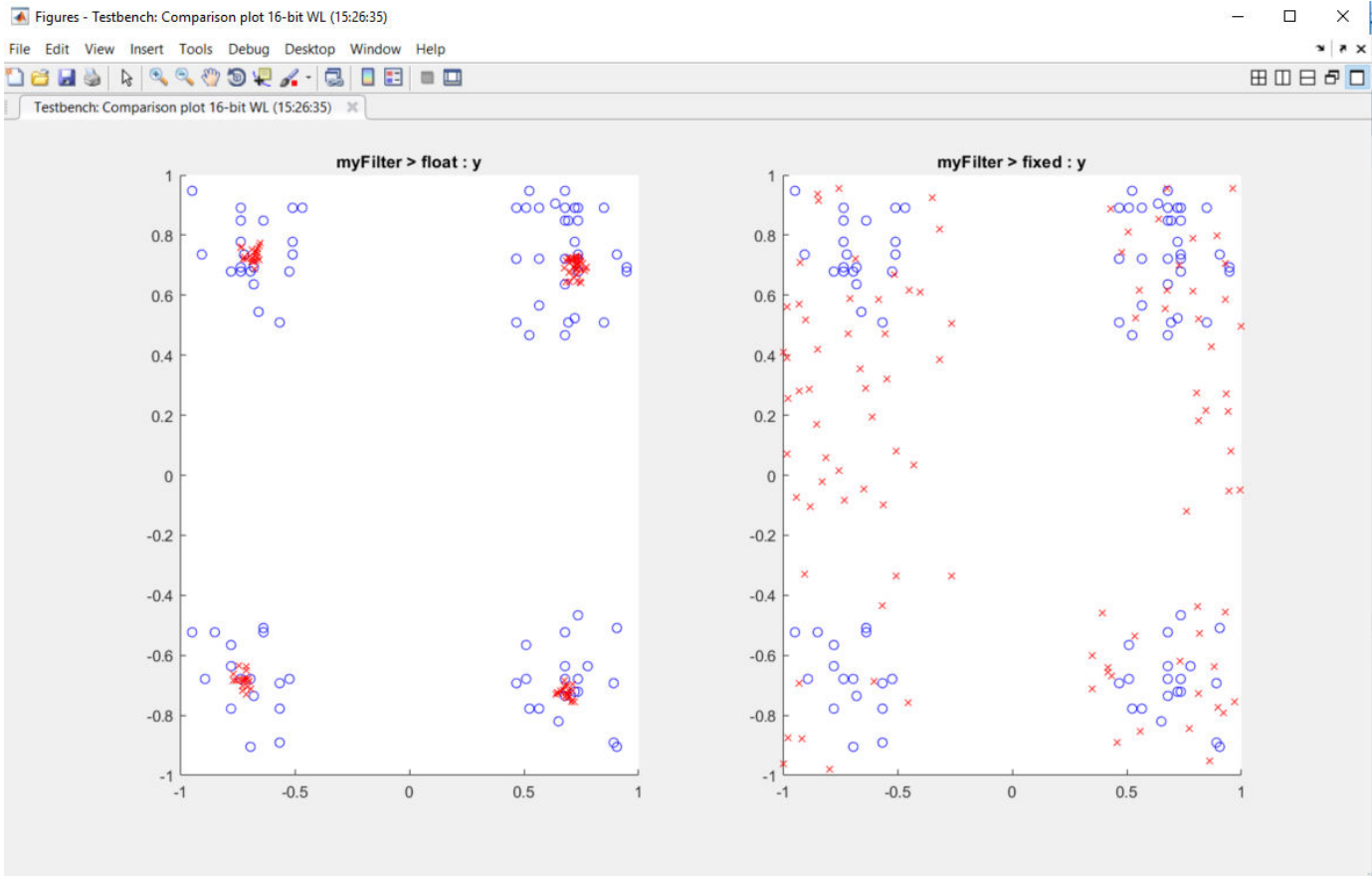
Convert to Fixed Point

Convert the floating-point MATLAB function, `myFilter`, to fixed-point MATLAB code. You do not need to specify input types for the `fiaccl` command because it infers the types from the test file.

```
fiaccl -args {complex(0, 0)} -float2fixed fxptcfg myFilter
```

The conversion process generates fixed-point code using a default word length of 16 and then runs a fixed-point simulation by running the `myFilterTest.m` function and calling the fixed-point version of `myFilter.m`.

Because you selected to log inputs and outputs for comparison plots and to use the custom plotting function, `plotDiff.m`, for these plots, the conversion process uses this function to generate the comparison plot.

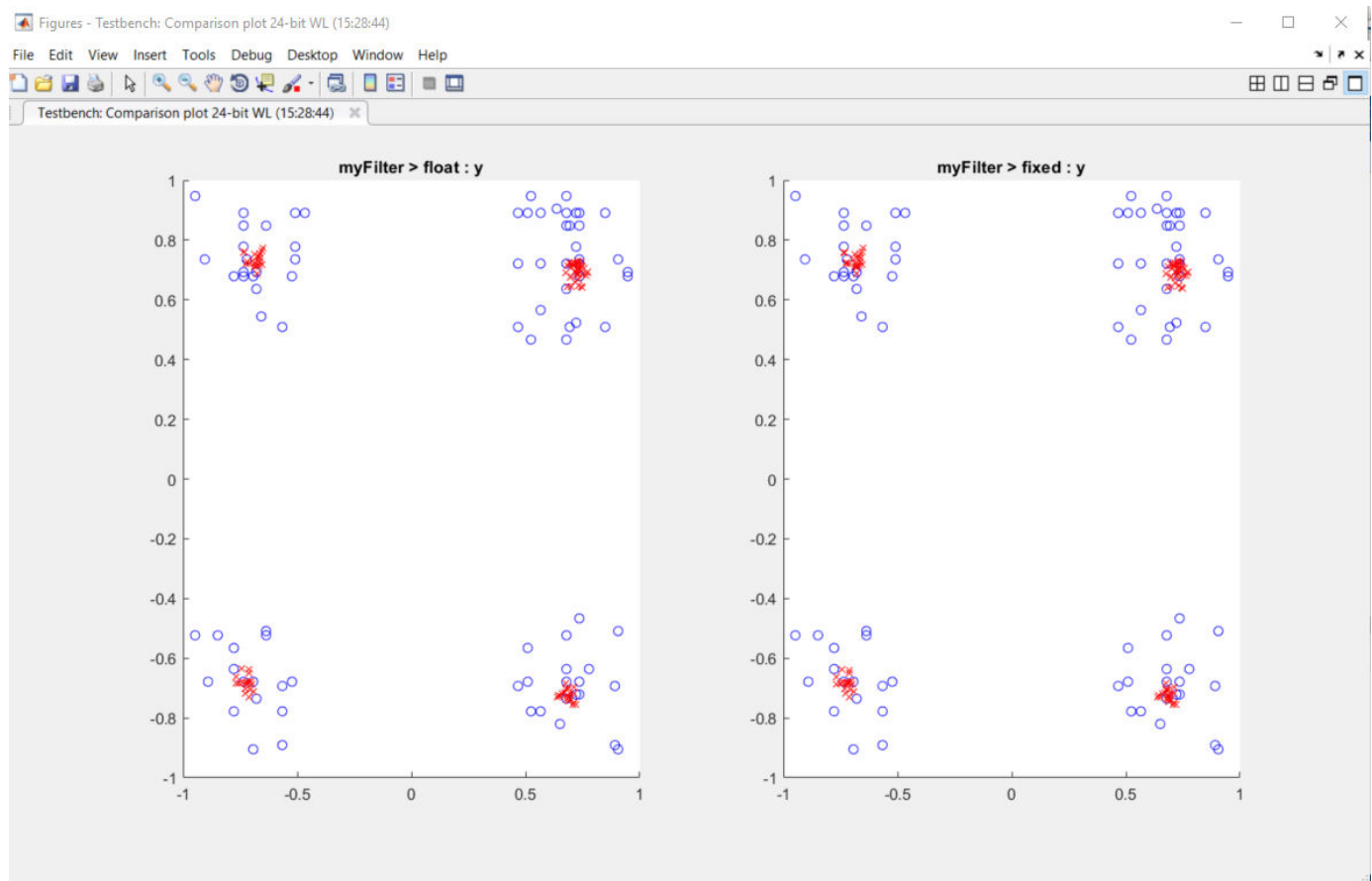


The plot shows that the fixed-point results do not closely match the floating-point results.

Increase the word length to 24 and then convert to fixed point again.

```
fxptcfg.DefaultWordLength = 24;
fiaccl -args {complex(0, 0)} -float2fixed fxptcfg myFilter
```

The increased word length improved the results. This time, the plot shows that the fixed-point results match the floating-point results.



See Also

More About

- "Custom Plot Functions" on page 7-51

Enable Plotting Using the Simulation Data Inspector

You can use the Simulation Data Inspector to inspect and compare floating-point and fixed-point input and output data logged using the `fiaccel` function. At the MATLAB command line:

- 1 Create a fixed-point configuration object and configure the test file name.

```
fixptcfg = coder.config('fixpt');  
fixptcfg.TestBenchName = 'dti_test';
```

- 2 Select to run the test file to verify the generated fixed-point MATLAB code. Log inputs and outputs for comparison plotting and select to use the Simulation Data Inspector to plot the results.

```
fixptcfg.TestNumerics = true;  
fixptcfg.LogIOForComparisonPlotting = true;  
fixptcfg.PlotWithSimulationDataInspector = true;
```

- 3 Generate fixed-point MATLAB code using `fiaccel`.

```
fiaccel -float2fixed fixptcfg dti
```

For an example, see “Propose Data Types Based on Derived Ranges” on page 9-6.

Single-Precision Conversion

- “Generate Single-Precision MATLAB Code” on page 10-2
- “MATLAB Language Features Supported for Single-Precision Conversion” on page 10-8
- “Single-Precision Conversion Best Practices” on page 10-10

Generate Single-Precision MATLAB Code

This example shows how to generate single-precision MATLAB code from double-precision MATLAB code.

Prerequisites

To complete this example, install the following products:

- MATLAB
- Fixed-Point Designer
- C compiler

See https://www.mathworks.com/support/compilers/current_release/.

You can use `mex -setup` to change the default compiler. See “Change Default Compiler”.

Create a Folder and Copy Relevant Files

- 1 In a local, writable folder, create a function `ex_2ndOrder_filter.m`.

```
function y = ex_2ndOrder_filter(x) %#codegen
    persistent z
    if isempty(z)
        z = zeros(2,1);
    end
    % [b,a] = butter(2, 0.25)
    b = [0.0976310729378175, 0.195262145875635, 0.0976310729378175];
    a = [1, -0.942809041582063, 0.333333333333333];

    y = zeros(size(x));
    for i = 1:length(x)
        y(i) = b(1)*x(i) + z(1);
        z(1) = b(2)*x(i) + z(2) - a(2) * y(i);
        z(2) = b(3)*x(i) - a(3) * y(i);
    end
end
```

- 2 Create a test file, `ex_2ndOrder_filter_test.m`, to exercise the `ex_2ndOrder_filter` algorithm.

It is a best practice to create a separate test script for preprocessing and postprocessing such as:

- Setting up input values.
- Calling the function under test.
- Outputting the test results.

To cover the full intended operating range of the system, the test script runs the `ex_2ndOrder_filter` function with three input signals: chirp, step, and impulse. The script then plots the outputs.

```
% ex_2ndOrder_filter_test
%
```

```

% Define representative inputs
N = 256; % Number of points
t = linspace(0,1,N); % Time vector from 0 to 1 second
f1 = N/2; % Target frequency of chirp set to Nyquist
x_chirp = sin(pi*f1*t.^2); % Linear chirp from 0 to Fs/2 Hz in 1 second
x_step = ones(1,N); % Step
x_impulse = zeros(1,N); % Impulse
x_impulse(1) = 1;

% Run the function under test
x = [x_chirp;x_step;x_impulse];
y = zeros(size(x));
for i = 1:size(x,1)
    y(i,:) = ex_2ndOrder_filter(x(i,:));
end

% Plot the results
titles = {'Chirp','Step','Impulse'}
clf
for i = 1:size(x,1)
    subplot(size(x,1),1,i)
    plot(t,x(i,:),t,y(i,:))
    title(titles{i})
    legend('Input','Output')
end
xlabel('Time (s)')
figure(gcf)

disp('Test complete.')

```

| Type | Name | Description |
|---------------|---------------------------|---|
| Function code | ex_2ndOrder_filter.m | Entry-point MATLAB function |
| Test file | ex_2ndOrder_filter_test.m | MATLAB script that tests ex_2ndOrder_filter.m |

Set Up the Single-Precision Configuration Object

Create a single-precision configuration object. Specify the test file name. Verify the single-precision code using the test file. Plot the error between the double-precision code and single-precision code. Use the default values for the other properties.

```

scfg = coder.config('single');
scfg.TestBenchName = 'ex_2ndOrder_filter_test';
scfg.TestNumerics = true;
scfg.LogIOForComparisonPlotting = true;

```

Generate Single-Precision MATLAB Code

To convert the double-precision MATLAB function, `ex_2ndOrder_filter`, to single-precision MATLAB code, use the `convertToSingle`

```
convertToSingle -config scfg ex_2ndOrder_filter
```

`convertToSingle` analyzes the double-precision code. The conversion process infers types by running the test file because you did not specify the input types for the `ex_2ndOrder_filter`

function. The conversion process selects single-precision types for the double-precision variables. It selects int32 for index variables. When the conversion is complete, `convertToSingle` generates a type proposal report.

View the Type Proposal Report

To see the types that the conversion process selected for the variables, open the type proposal report for the `ex_2ndOrder_filter` function. Click the link `ex_2ndOrder_filter_report.html`.

The report opens in a web browser. The conversion process converted:

- Double-precision variables to `single`.
- The index `i` to `int32`. The conversion process casts index and dimension variables to `int32`.

Single-Precision Report `ex_2ndOrder_filter`

| Simulation Coverage | Code |
|---------------------|---|
| 100% | <code>function y = ex_2ndOrder_filter(x) %#codegen</code> |
| Once | <code>persistent z</code> |
| | <code>if isempty(z)</code> |
| | <code>z = zeros(2,1);</code> |
| | <code>end</code> |
| 100% | <code>% [b,a] = butter(2, 0.25)</code> |
| | <code>b = [0.0976310729378175, 0.195262145875635, 0.0976310729378175];</code> |
| | <code>a = [1, -0.942809041582063, 0.333333333333333];</code> |
| | |
| | <code>y = zeros(size(x));</code> |
| | <code>for i=1:length(x)</code> |
| | <code>y(i) = b(1)*x(i) + z(1);</code> |
| | <code>z(1) = b(2)*x(i) + z(2) - a(2) * y(i);</code> |
| | <code>z(2) = b(3)*x(i) - a(3) * y(i);</code> |
| | <code>end</code> |
| | <code>end</code> |

| Variable Name | Type | Sim Min | Sim Max | Whole Number | ProposedType |
|---------------|----------------|---------------------|--------------------|--------------|--------------|
| a | double 1 x 3 | -0.942809041582063 | 1 | No | single |
| b | double 1 x 3 | 0.0976310729378175 | 0.195262145875635 | No | single |
| i | double | 1 | 256 | Yes | int32 |
| x | double 1 x 256 | -0.9999756307053946 | 1 | No | single |
| y | double 1 x 256 | -0.9696817930434206 | 1.0553496057969345 | No | single |
| z | double 2 x 1 | -0.8907046852192462 | 0.957718532859117 | No | single |

View Generated Single-Precision MATLAB Code

To view the report for the generation of the single-precision MATLAB code, in the Command Window:

- 1 Scroll to the `Generate Single-Precision Code` step. Click the **View report** link.

2 In the **MATLAB Source** pane, click `ex_2ndOrder_filter_single`.

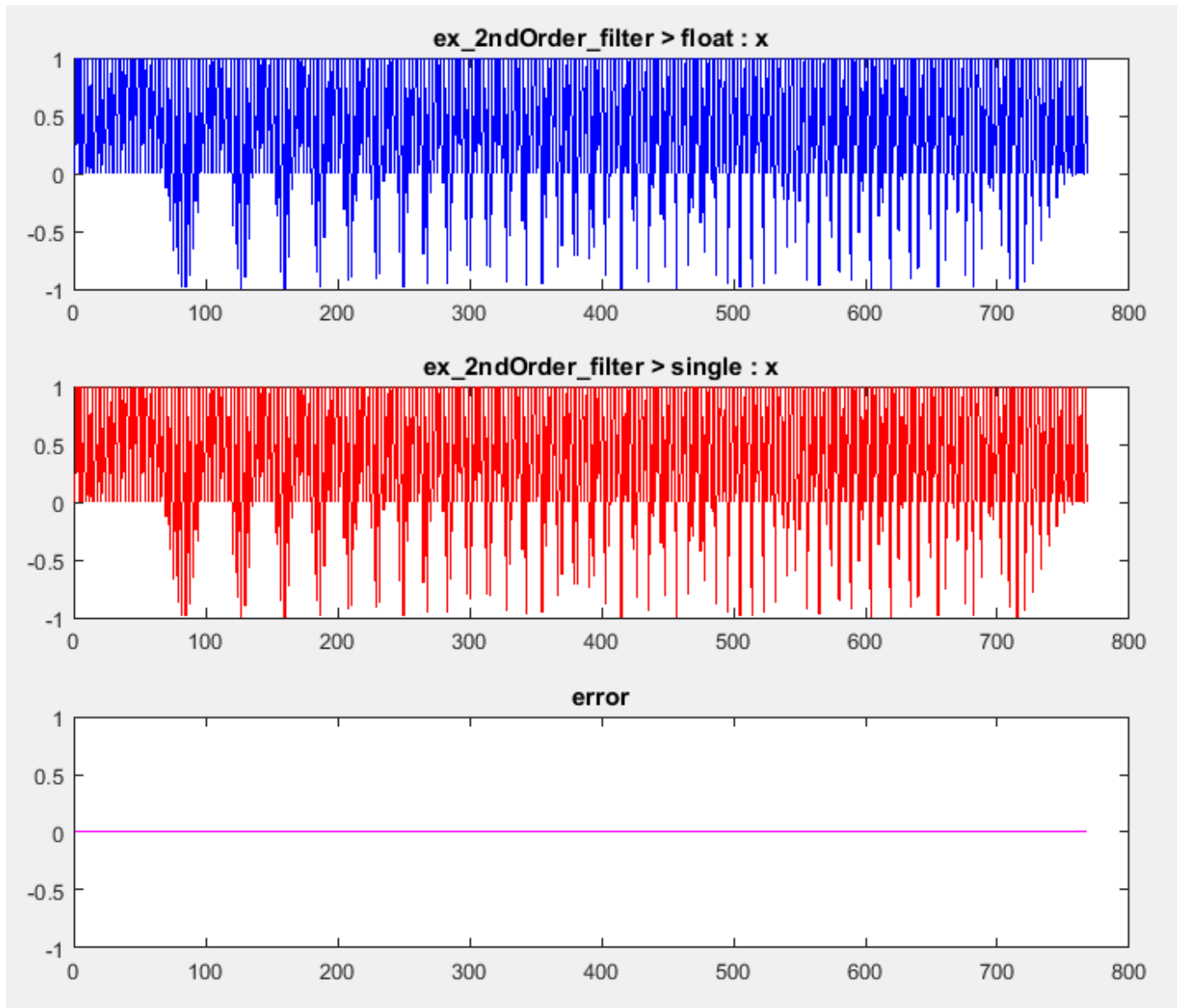
The code generation report displays the single-precision MATLAB code for `ex_2ndOrder_filter`.

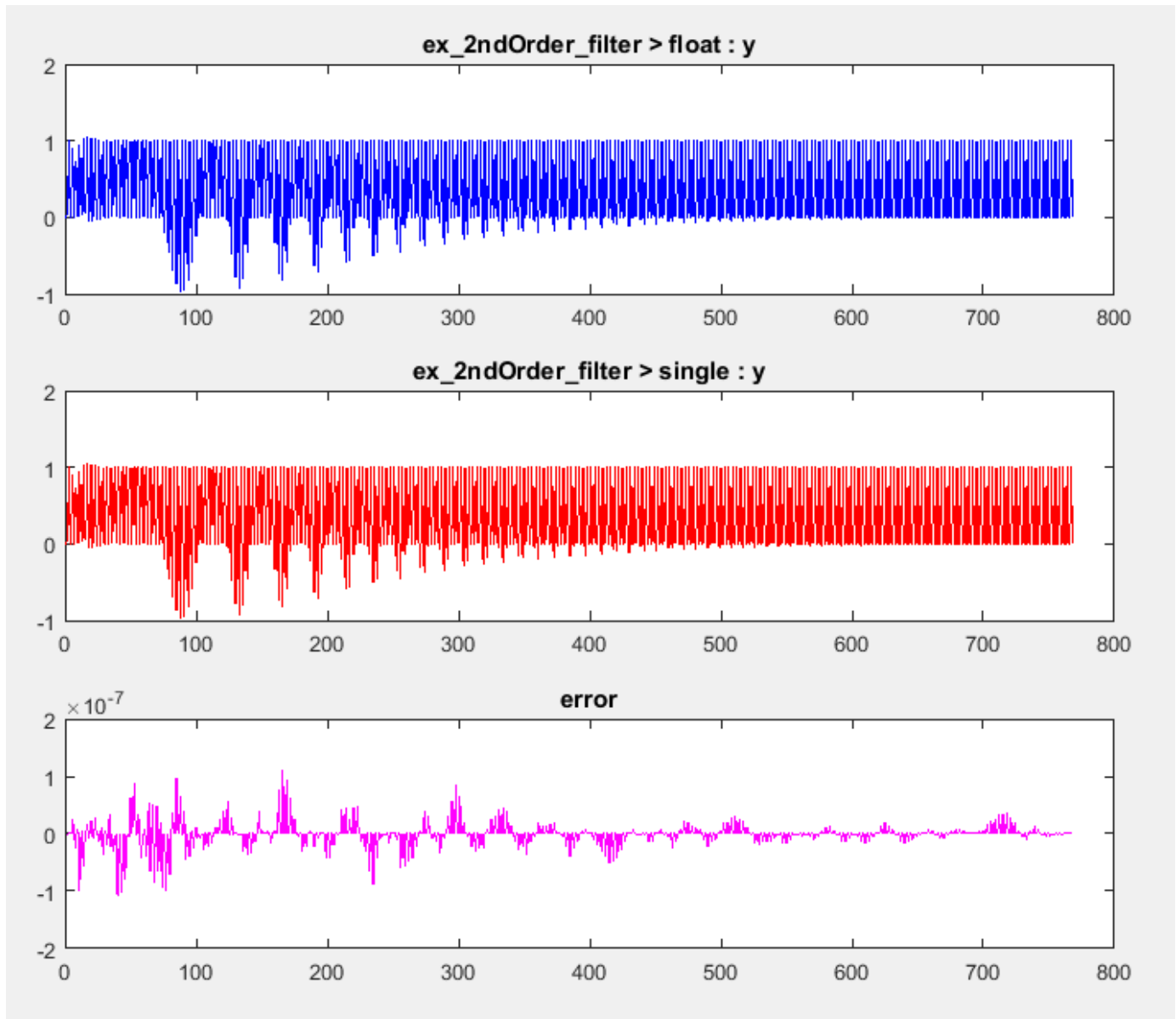
View Potential Data Type Issues

When you generate single-precision code, `convertToSingle` enables highlighting of potential data type issues in code generation reports. If `convertToSingle` cannot remove a double-precision operation, the report highlights the MATLAB expression that results in the operation. Click the **Code Insights** tab. The absence of potential data type issues indicates that no double-precision operations remain.

Compare the Double-Precision and Single-Precision Variables

You can see the comparison plots for the input `x` and output `y` because you selected to log inputs and outputs for comparison plots .





See Also

[coder.SingleConfig](#) | [coder.config](#) | [convertToSingle](#)

More About

- “Single-Precision Conversion Best Practices” on page 10-10

MATLAB Language Features Supported for Single-Precision Conversion

In this section...

“MATLAB Language Features Supported for Single-Precision Conversion” on page 10-8

“MATLAB Language Features Not Supported for Single-Precision Conversion” on page 10-9

MATLAB Language Features Supported for Single-Precision Conversion

Single-precision conversion supports the following MATLAB language features:

- N-dimensional arrays.
- Matrix operations, including deletion of rows and columns.
- Variable-size data. Comparison plotting does not support variable-size data.
- Subscripting (see “Incompatibility with MATLAB in Matrix Indexing Operations for Code Generation” on page 29-22).
- Complex numbers (see “Code Generation for Complex Data” on page 16-8).
- Numeric classes (see “Supported Variable Types” on page 18-13).
- Program control statements `if`, `switch`, `for`, `while`, and `break`.
- Arithmetic, relational, and logical operators.
- Local functions.
- Global variables.
- Persistent variables.
- Structures.
- Characters.

Single-precision conversion does not support the complete set of Unicode characters. Characters are restricted to 8 bits of precision in generated code. Many mathematical operations require more than 8 bits of precision. If you intend to convert your MATLAB algorithm to single precision, it is a best practice not to perform arithmetic with characters.

- MATLAB classes. Single-precision conversion supports:
 - Class properties
 - Constructors
 - Methods
 - Specializations

It does not support class inheritance or packages.

Single-precision conversion using `codegen` with the `-singleC` option does not support classes when the properties have default values. Property values must be initialized in the constructor. Constant properties cannot be initialized to double precision data types.

- Function calls (see “Resolution of Function Calls for Code Generation” on page 14-2)

MATLAB Language Features Not Supported for Single-Precision Conversion

Single-precision conversion does not support the following features:

- Anonymous functions
- Cell arrays
- String scalars
- Objects of value classes as entry-point function inputs or outputs
- Function handles
- Java
- Nested functions
- Recursion
- Sparse matrices
- `try/catch` statements
- `varargin` and `varargout`, or generation of fewer input or output arguments than an entry-point function defines

Single-Precision Conversion Best Practices

In this section...

“Use Integers for Index Variables” on page 10-10
“Limit Use of assert Statements” on page 10-10
“Initialize MATLAB Class Properties in Constructor” on page 10-10
“Provide a Test File That Calls Your MATLAB Function” on page 10-10
“Prepare Your Code for Code Generation” on page 10-11
“Use the -args Option to Specify Input Properties” on page 10-11
“Test Numerics and Log I/O Data” on page 10-11

Use Integers for Index Variables

In MATLAB code that you want to convert to single precision, it is a best practice to use integers for index variables. However, if the code does not use integers for index variables, when possible `convertToSingle` tries to detect the index variables and select `int32` types for them.

Limit Use of assert Statements

- Do not use `assert` statements to define the properties of input arguments.
- Do not use `assert` statements to test the type of a variable. For example, do not use

```
assert(isa(a, 'double'))
```

Initialize MATLAB Class Properties in Constructor

Do not initialize MATLAB class properties in the `properties` block. Instead, use the constructor to initialize the class properties.

Provide a Test File That Calls Your MATLAB Function

Separate your core algorithm from other code that you use to test and verify the results. Create a test file that calls your double-precision MATLAB algorithm. You can use the test file to:

- Automatically define properties of the top-level function inputs.
- Verify that the double-precision algorithm behaves as you expect. The double-precision behavior is the baseline against which you compare the behavior of the single-precision versions of your algorithm.
- Compare the behavior of the single-precision version of your algorithm to the double-precision baseline.

For best results, the test file must exercise the algorithm over its full operating range.

Prepare Your Code for Code Generation

MATLAB code that you want to convert to single precision must comply with code generation requirements. See “MATLAB Language Features Supported for C/C++ Code Generation” on page 19-29.

To help you identify unsupported functions or constructs in your MATLAB code, add the `%#codegen` pragma to the top of your MATLAB file. When you edit your code in the MATLAB editor, the MATLAB Code Analyzer flags functions and constructs that are not supported for code generation. See “Check Code Using the MATLAB Code Analyzer” on page 12-65. When you use the MATLAB Coder app, the app screens your code for code generation readiness. At the function line, you can use the Code Generation Readiness Tool. See “Check Code Using the Code Generation Readiness Tool” on page 12-64.

Use the `-args` Option to Specify Input Properties

When you generate single-precision MATLAB code, if you specify a test file, you do not have to specify argument properties with the `-args` option. In this case, the code generator runs the test file to determine the properties of the input types. However, running the test file can slow the code generation. It is a best practice to pass the properties to the `-args` option so that `convertToSingle` does not run the test file to determine the argument properties. If you have a MATLAB Coder license, you can use `coder.getArgTypes` to determine the argument properties. For example:

```
types = coder.getArgTypes('myfun_test', 'myfun');
scfg = coder.config('single');
convertToSingle -config scfg -args types myfun
```

Test Numerics and Log I/O Data

When you use the `convertToSingle` function to generate single-precision MATLAB code, enable numerics testing and I/O data logging for comparison plots. To use numerics testing, you must provide a test file that calls your MATLAB function. To enable numerics testing and I/O data logging, create a `coder.SingleConfig` object. Set the `TestBenchName`, `TestNumerics`, and `LogIOForComparisonPlotting` properties. For example:

```
scfg = coder.config('single');
scfg.TestBenchName = 'mytest';
scfg.TestNumerics = true;
scfg.LogIOForComparisonPlotting = true;
```


Fixed-Point Conversion — Manual Conversion

- “Manual Fixed-Point Conversion Workflow” on page 11-2
- “Manual Fixed-Point Conversion Best Practices” on page 11-3
- “Fixed-Point Design Exploration in Parallel” on page 11-15
- “Real-Time Image Acquisition, Image Processing, and Fixed-Point Blob Analysis for Target Practice Analysis” on page 11-20

Manual Fixed-Point Conversion Workflow

- 1 Implement your algorithm in MATLAB.
- 2 Write a test file that calls your original MATLAB algorithm to validate the behavior of your algorithm.

Create a test file to validate that the algorithm works as expected in floating point before converting it to fixed point. Use the same test file to propose fixed-point data types. After the conversion, use this test file to compare fixed-point results to the floating-point baseline.

- 3 Prepare algorithm for instrumentation.
- 4 Write an entry-point function.

For instrumentation and code generation, it is convenient to have an entry-point function that calls the function to be converted to fixed point. You can cast the function inputs to different data types, and add calls to different variations of the algorithm for comparison. By using an entry-point function, you can run both fixed-point and floating-point variants of your algorithm. You can also run different variants of fixed-point. This approach allows you to iterate on your code more quickly to arrive at the optimal fixed-point design.

- 5 Build instrumented MEX for original MATLAB algorithm.
- 6 Run your original MATLAB algorithm to log min/max data. View this data in the instrumentation report.
- 7 Separate data types from algorithm.

Convert functions to use types tables and update entry-point function.

- 8 Validate modified function.
 - a Create fixed-point types table based on proposed data types.
 - b Build MEX function.
 - c Run and compare MEX function behavior against baseline.
- 9 Use proposed fixed-point data types.

Create fixed-point types table based on proposed data types, build mex, run, and then compare against baseline.

- 10 Optionally, if have a MATLAB Coder license, generate code.

Start by testing native C-types.

- 11 Iterate, tune algorithm.

For example, tune the algorithm to avoid overflow or eliminate bias.

Manual Fixed-Point Conversion Best Practices

In this section...

“Create a Test File” on page 11-3

“Prepare Your Algorithm for Code Acceleration or Code Generation” on page 11-4

“Check for Fixed-Point Support for Functions Used in Your Algorithm” on page 11-5

“Manage Data Types and Control Bit Growth” on page 11-6

“Separate Data Type Definitions from Algorithm” on page 11-6

“Convert to Fixed Point” on page 11-7

“Optimize Data Types” on page 11-9

“Optimize Your Algorithm” on page 11-12

Fixed-Point Designer software helps you design and convert your algorithms to fixed point. Whether you are simply designing fixed-point algorithms in MATLAB or using Fixed-Point Designer in conjunction with MathWorks code generation products, these best practices help you get from generic MATLAB code to an efficient fixed-point implementation. These best practices are also covered in this webinar: Manual Fixed-Point Conversion Best Practices Webinar

Create a Test File

A best practice for structuring your code is to separate your core algorithm from other code that you use to test and verify the results. Create a test file to call your original MATLAB algorithm and fixed-point versions of the algorithm. For example, as shown in the following table, you might set up some input data to feed into your algorithm, and then, after you process that data, create some plots to verify the results. Since you need to convert only the algorithmic portion to fixed-point, it is more efficient to structure your code so that you have a test file, in which you create your inputs, call your algorithm, and plot the results, and one (or more) algorithmic files, in which you do the core processing.

| Original code | Best Practice | Modified code |
|---|--|--|
| <pre>% TEST INPUT x = randn(100,1); % ALGORITHM y = zeros(size(x)); y(1) = x(1); for n=2:length(x) y(n)=y(n-1) + x(n); end % VERIFY RESULTS yExpected=cumsum(x); plot(y-yExpected) title('Error')</pre> | <p>Issue</p> <p>Generation of test input and verification of results are intermingled with the algorithm code.</p> <p>Fix</p> <p>Create a test file that is separate from your algorithm. Put the algorithm in its own function.</p> | <p>Test file</p> <pre>% TEST INPUT x = randn(100,1); % ALGORITHM y = cumulative_sum(x); % VERIFY RESULTS yExpected = cumsum(x); plot(y-yExpected) title('Error')</pre> <p>Algorithm in its own function</p> <pre>function y = cumulative_sum(x) y = zeros(size(x)); y(1) = x(1); for n=2:length(x) y(n) = y(n-1) + x(n); end end</pre> |

You can use the test file to:

- Verify that your floating-point algorithm behaves as you expect before you convert it to fixed point. The floating-point algorithm behavior is the baseline against which you compare the behavior of the fixed-point versions of your algorithm.
- Propose fixed-point data types.
- Compare the behavior of the fixed-point versions of your algorithm to the floating-point baseline.

Your test files should exercise the algorithm over its full operating range so that the simulation ranges are accurate. For example, for a filter, realistic inputs are impulses, sums of sinusoids, and chirp signals. With these inputs, using linear theory, you can verify that the outputs are correct. Signals that produce maximum output are useful for verifying that your system does not overflow. The quality of the proposed fixed-point data types depends on how well the test files cover the operating range of the algorithm with the accuracy that you want.

Prepare Your Algorithm for Code Acceleration or Code Generation

Using Fixed-Point Designer, you can:

- Instrument your code and provide data type proposals to help you convert your algorithm to fixed point, using the following functions:
 - `buildInstrumentedMex`, which generates compiled C code that includes logging instrumentation.
 - `showInstrumentationResults`, which shows the results logged by the instrumented, compiled C code.
 - `clearInstrumentationResults`, which clears the logged instrumentation results from memory.

- Accelerate your fixed-point algorithms by creating a MEX file using the `fiaccl` function.

Any MATLAB algorithms that you want to instrument using `buildInstrumentedMex` and any fixed-point algorithms that you want to accelerate using `fiaccl` must comply with code generation requirements and rules. To view the subset of the MATLAB language that is supported for code generation, see “Functions and Objects Supported for C/C++ Code Generation” on page 26-2.

To help you identify unsupported functions or constructs in your MATLAB code, use one of the following tools.

- Add the `%#codegen` pragma to the top of your MATLAB file. The MATLAB code analyzer flags functions and constructs that are not available in the subset of the MATLAB language supported for code generation. This advice appears in real-time as you edit your code in the MATLAB editor.

For more information, see “Check Code Using the MATLAB Code Analyzer” on page 12-65.

- Use the Code Generation Readiness tool to generate a static report on your code. The report identifies calls to functions and the use of data types that are not supported for code generation. To generate a report for a function, `myFunction1`, at the command line, enter `coder.screener('myFunction1')`.

For more information, see “Check Code Using the Code Generation Readiness Tool” on page 12-64.

Check for Fixed-Point Support for Functions Used in Your Algorithm

Before you start your fixed-point conversion, identify which functions used in your algorithm are not supported for fixed point. Consider how you might replace them or otherwise modify your implementation to be more optimized for embedded targets. For example, you might need to find (or write your own) replacements for functions like `log2`, `fft`, and `exp`. Other functions like `sin`, `cos`, and `sqrt` may support fixed point, but for better efficiency, you may want to consider an alternative implementation like a lookup table or CORDIC-based algorithm.

If you cannot find a replacement immediately, you can continue converting the rest of your algorithm to fixed point by simply insulating any functions that don’t support fixed-point with a cast to double at the input, and a cast back to a fixed-point type at the output.

| Original Code | Best Practice | Modified Code |
|----------------------------|---|------------------------------------|
| <code>y = 1/exp(x);</code> | <p>Issue</p> <p>The <code>exp()</code> function is not defined for fixed-point inputs.</p> <p>Fix</p> <p>Cast the input to double until you have a replacement. In this case, <code>1/exp(x)</code> is more suitable for fixed-point growth than <code>exp(x)</code>, so replace the whole expression with a <code>1/exp</code> function, possibly as a lookup table.</p> | <code>y = 1/exp(double(x));</code> |

Manage Data Types and Control Bit Growth

The `(:)=` syntax is known as subscripted assignment. When you use this syntax, MATLAB overwrites the value of the left-hand side argument, but retains the existing data type and array size. This is particularly important in keeping fixed-point variables fixed point (as opposed to inadvertently turning them into doubles), and for preventing bit growth when you want to maintain a particular data type for the output.

| Original Code | Best Practice | Modified Code |
|--|---|---|
| <pre>acc = 0; for n = 1:numel(x) acc = acc + x(n); end</pre> | <p>Issue</p> <p><code>acc = acc + x(n)</code> overwrites <code>acc</code> with <code>acc + x(n)</code>. When you are using all double types, this behavior is fine. However, when you introduce fixed-point data types in your code, if <code>acc</code> is overwritten, the data type of <code>acc</code> might change.</p> <p>Fix</p> <p>To preserve the original data type of <code>acc</code>, assign into <code>acc</code> using <code>acc(:)=</code>. Using subscripted assignment casts the right-hand-side value into the same data type as <code>acc</code> and prevents bit growth.</p> | <pre>acc = 0; for n = 1:numel(x) acc(:) = acc + x(n); end</pre> |

For more information, see “Controlling Bit Growth”.

Separate Data Type Definitions from Algorithm

For instrumentation and code generation, create an entry-point function that calls the function that you want to convert to fixed point. You can then cast the function inputs to different data types. You can add calls to different variations of the function for comparison. By using an entry-point function, you can run both fixed-point and floating-point variants of your algorithm. You can also run different variants of fixed-point. This approach allows you to iterate on your code more quickly to arrive at the optimal fixed-point design.

This method of fixed-point conversion makes it easier for you to compare several different fixed-point implementations, and also allows you to easily retarget your algorithm to a different device.

To separate data type definitions from your algorithm:

- 1 When a variable is first defined, use `cast(x, 'like', y)` or `zeros(m, n, 'like', y)` to cast it to your desired data type.
- 2 Create a table of data type definitions, starting with original data types used in your code. Before converting to fixed point, create a data type table that uses all single data types to find type mismatches and other problems.
- 3 Run your code connected to each table and look at the results to verify the connection.

| Original Code | Best Practice | Modified Code |
|---|--|--|
| <pre>% Algorithm n = 128; y = zeros(size(n));</pre> | <p>Issue</p> <p>The default data type in MATLAB is double-precision floating-point.</p> <p>Fix</p> <ol style="list-style-type: none"> 1 Use <code>cast(..., 'like', ...)</code> and <code>zeros(... 'like', ...)</code> to programmatically specify types that are defined in a separate table. 2 Create an original types table, usually in a separate function. 3 Add single data types to your table to help verify the connection with your code. | <pre>% Algorithm T = mytypes('double'); n = cast(128, 'like', T.n); y = zeros(size(n), 'like', T.y); function T = mytypes(dt) switch(dt) case 'double' T.n = double([]); T.y = double([]); case 'single' T.n = single([]); T.y = single([]); end end</pre> |

Separating data type specifications from algorithm code enables you to:

- Reuse your algorithm code with different data types.
- Keep your algorithm uncluttered with data type specifications and switch statements for different data types.
- Improve readability of your algorithm code.
- Switch between fixed-point and floating-point data types to compare baselines.
- Switch between variations of fixed-point settings without changing the algorithm code.

Convert to Fixed Point

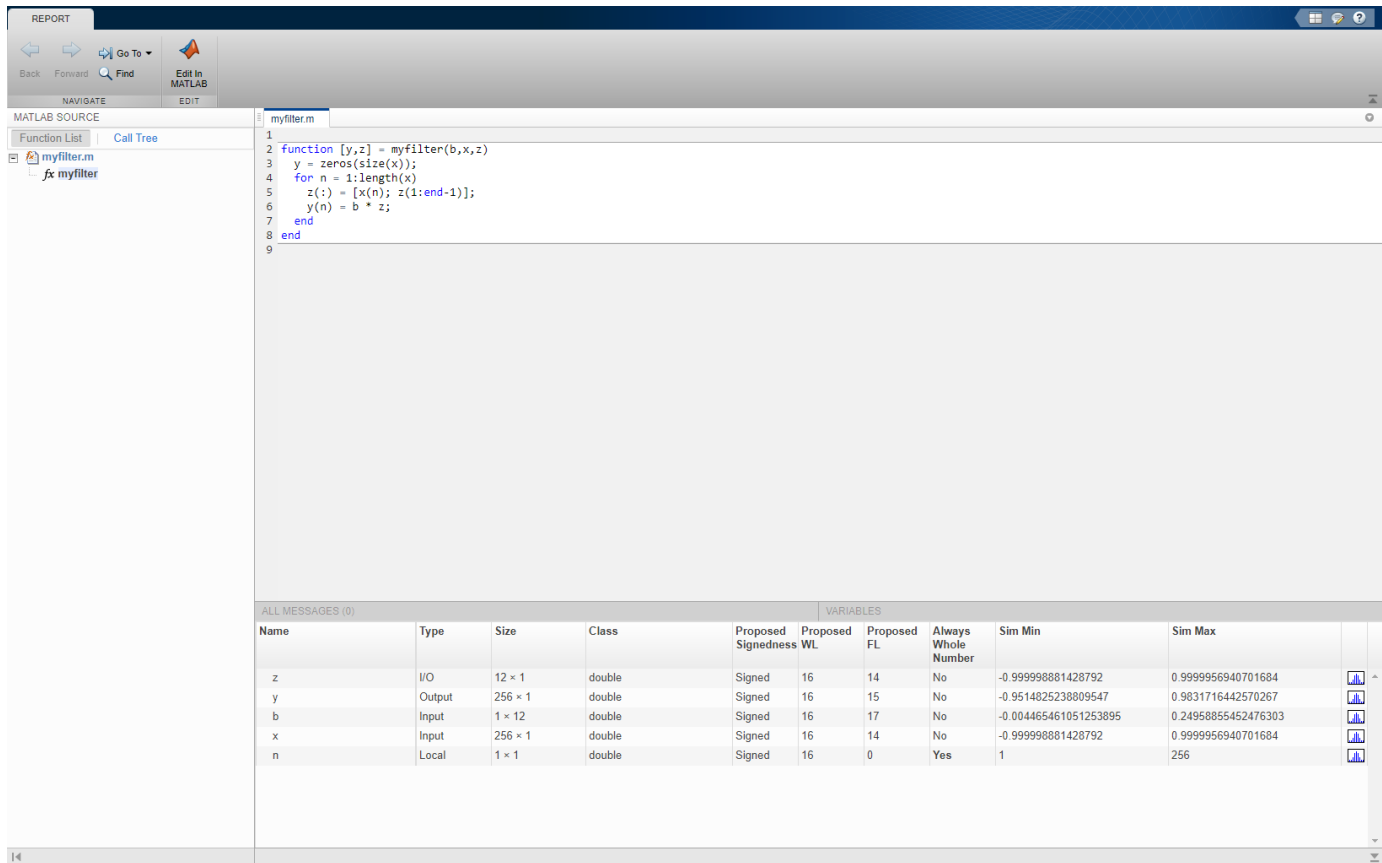
What Are Your Goals for Converting to Fixed Point?

Before you start the conversion, consider your goals for converting to fixed point. Are you implementing your algorithm in C or HDL? What are your target constraints? The answers to these questions determine many fixed-point properties such as the available word length, fraction length, and math modes, as well as available math libraries.

Build and Run an Instrumented MEX Function

Build and run an instrumented MEX function to get fixed-point types proposals using the `buildInstrumentedMex` and `showInstrumentationResults` functions. Test files should exercise your algorithm over its full operating range. The quality of the proposed fixed-point data types depends on how well the test file covers the operating range of the algorithm with the accuracy that you want. A simple set of test vectors may not exercise the full range of types, so use the proposals as a guideline for choosing an initial set of fixed-point types, and use your best judgement and experience in adjusting the types. If loop indices are used only as index variables, they are automatically converted to integer types, so you do not have to explicitly convert them to fixed point.

| Algorithm Code | Test File |
|--|--|
| <pre>function [y,z] = myfilter(b,x,z) y = zeros(size(x)); for n = 1:length(x) z(:) = [x(n); z(1:end-1)]; y(n) = b * z; end end</pre> | <pre>% Test inputs b = fir1(11,0.25); t = linspace(0,10*pi,256)'; x = sin((pi/16)*t.^2); % Linear chirp z = zeros(size(b')); % Build buildInstrumentedMex myfilter ... -args {b,x,z} -histogram % Run [y,z] = myfilter_mex(b,x,z); % Show showInstrumentationResults myfilter_mex ... -defaultDT numerictype(1,16) -proposeFL</pre> |



Create a Types Table

Create a types table using a structure with prototypes for the variables. The proposed types are computed from the simulation runs. A long simulation run with a wide range of expected data produces better proposals. You can use the proposed types or use your knowledge of the algorithm and implementation constraints to improve the proposals.

Because the data types, not the values, are used, specify the prototype values as empty ([]).

In some cases, it might be more efficient to leave some parts of the code in floating point. For example, when there is high dynamic range or that part of the code is sensitive to round-off errors.

Algorithm Code

```
function [y,z]=myfilter(b,x,z,T)
    y = zeros(size(x),'like',T.y);
    for n = 1:length(x)
        z(:) = [x(n); z(1:end-1)];
        y(n) = b * z;
    end
end
```

Types Tables

```
function T = mytypes(dt)
    switch dt
        case 'double'
            T.b = double([]);
            T.x = double([]);
            T.y = double([]);

        case 'fixed16'
            T.b = fi([],true,16,15);
            T.x = fi([],true,16,15);
            T.y = fi([],true,16,14);
    end
end
```

Test File

```
% Test inputs
b = firl(11,0.25);
t = linspace(0,10*pi,256)';
x = sin((pi/16)*t.^2);
% Linear chirp

% Cast inputs
T=mytypes('fixed16');
b=cast(b,'like',T.b);
x=cast(x,'like',T.x);
z=zeros(size(b),'like',T.x);

% Run
[y,z] = myfilter(b,x,z,T);
```

Run With Fixed-Point Types and Compare Results

Create a test file to validate that the floating-point algorithm works as expected before converting it to fixed point. You can use the same test file to propose fixed-point data types, and to compare fixed-point results to the floating-point baseline after the conversion.

Optimize Data Types

Use Scaled Doubles

Use scaled doubles to detect potential overflows. Scaled doubles are a hybrid between floating-point and fixed-point numbers. Fixed-Point Designer stores them as doubles with the scaling, sign, and


word length information retained. To use scaled doubles, you can use the data type override (DTO) property or you can set the 'DataType' property to 'ScaledDouble' in the fi or numerictype constructor.

| To... | Use... | Example |
|---------------------------------|-------------------------------------|--|
| Set data type override locally | numerictype DataType property | <pre>T.a = fi([],1,16,13,'DataType', 'ScaledDouble'); a = cast(pi, 'like', T.a)</pre> <p>a =</p> <p>3.1416</p> <p>DataTypeMode: Scaled double: binary point scaling Signedness: Signed WordLength: 16 FractionLength: 13</p> |
| Set data type override globally | fipref DataTypeOverride property | <pre>fipref('DataTypeOverride', 'ScaledDoubles') T.a = fi([],1,16,13);</pre> <p>a =</p> <p>3.1416</p> <p>DataTypeMode:Scaled double: binary point scaling Signedness: Signed WordLength:16 FractionLength:13</p> |

For more information, see “Scaled Doubles” on page 35-16.

Use the Histogram to Fine-Tune Data Type Settings

To fine-tune fixed-point type settings, run the `buildInstrumentedMex` function with the `-histogram` flag and then run the generated MEX function with your desired test inputs. When you use the `showInstrumentationResults` to display the code generation report, the report displays a Histogram icon. Click the icon to open the `NumericTypeScope` and view the distribution of values observed in your simulation for the selected variable.

Overflows indicated in red in the Code Generation Report show in the "outside range" bin in the `NumericTypeScope`. Launch the `NumericTypeScope` for an associated variable or expression by clicking on the histogram view icon .

Explore Design Tradeoffs

Once you have your first set of fixed-point data types, you can then add different variations of fixed-point values to your types table. You can modify and iterate to avoid overflows, adjust fraction lengths, and change rounding methods to eliminate bias.

| Algorithm Code |
|---|
| <pre>function [y,z] = myfilter(b,x,z,T) y = zeros(size(x),'like',T.y); for n = 1:length(x) z(:) = [x(n); z(1:end-1)]; y(n) = b * z; end end</pre> |

Types Tables

```
function T = mytypes(dt)
switch dt
case 'double'
    T.b = double([]);
    T.x = double([]);
    T.y = double([]);

case 'fixed8'
    T.b = fi([],true,8,7);
    T.x = fi([],true,8,7);
    T.y = fi([],true,8,6);

case 'fixed16'
    T.b = fi([],true,16,15);
    T.x = fi([],true,16,15);
    T.y = fi([],true,16,14);
end
end
```

Test File

```

function mytest
% Test inputs
b = fir1(11,0.25);
t = linspace(0,10*pi,256)';
x = sin((pi/16)*t.^2); % Linear chirp

% Run
y0 = entrypoint('double',b,x);
y8 = entrypoint('fixed8',b,x);
y16 = entrypoint('fixed16',b,x);

% Plot
subplot(3,1,1)
plot(t,x,'c',t,y0,'k')
legend('Input','Baseline output')
title('Baseline')

subplot(3,2,3)
plot(t,y8,'k')
title('8-bit fixed-point output')
subplot(3,2,4)
plot(t,y0-double(y8),'r')
title('8-bit fixed-point error')

subplot(3,2,5)
plot(t,y16,'k')
title('16-bit fixed-point output')
xlabel('Time (s)')
subplot(3,2,6)
plot(t,y0-double(y16),'r')
title('16-bit fixed-point error')
xlabel('Time (s)')
end

function [y,z] = entrypoint(dt,b,x)
T = mytypes(dt);
b = cast(b,'like',T.b);
x = cast(x,'like',T.x);
z = zeros(size(b),'like',T.x);
[y,z] = myfilter(b,x,z,T);
end

```

Optimize Your Algorithm

Use `fimath` to Get Natural Types for C or HDL

`fimath` properties define the rules for performing arithmetic operations on `fi` objects, including math, rounding, and overflow properties. You can use the `fimath` `ProductMode` and `SumMode` properties to retain natural data types for C and HDL. The `KeepLSB` setting for `ProductMode` and `SumMode` models the behavior of integer operations in the C language, while `KeepMSB` models the behavior of many DSP devices. Different rounding methods require different amounts of overhead code. Setting the `RoundingMethod` property to `Floor`, which is equivalent to two's complement truncation, provides the most efficient rounding implementation. Similarly, the standard method for

handling overflows is to wrap using modulo arithmetic. Other overflow handling methods create costly logic. Whenever possible, set the `OverflowAction` to `Wrap`.

| MATLAB Code | Best Practice | Generated C Code |
|---|--|--|
| <pre>% Code being compiled function y = adder(a,b) y = a + b; end With types defined with default fimath settings: T.a = fi([],1,16,0); T.b = fi([],1,16,0); a = cast(0,'like',T.a); b = cast(0,'like',T.b);</pre> | <p>Issue</p> <p>Additional code is generated to implement saturation overflow, nearest rounding, and full-precision arithmetic.</p> | <pre>int adder(short a, short b) { int y; int i; int i1; int i2; int i3; i = a; i1 = b; if ((i & 65536) != 0) { i2 = i -65536; } else { i2 = i & 65535; } if ((i1 & 65536) != 0) { i3 = i1 -65536; } else { i3 = i1 & 65535; } i = i2 + i3; if ((i & 65536) != 0) { y = i -65536; } else { y = i & 65535; } return y; }</pre> |
| <pre>Code being compiled function y = adder(a,b) y = a + b; end With types defined with fimath settings that match your processor types: F = fimath(... 'RoundingMethod','Floor', ... 'OverflowAction','Wrap', ... 'ProductMode','KeepLSB', ... 'ProductWordLength',32, ... 'SumMode','KeepLSB', ... 'SumWordLength',32); T.a = fi([],1,16,0,F); T.b = fi([],1,16,0,F); a = cast(0,'like',T.a); b = cast(0,'like',T.b);</pre> | <p>Fix</p> <p>To make the generated code more efficient, choose fixed-point math settings that match your processor types.</p> | <pre>int adder(short a, short b) { return a + b; }</pre> |

Replace Built-in Functions With More Efficient Fixed-Point Implementations

Some MATLAB built-in functions can be made more efficient for fixed-point implementation. For example, you can replace a built-in function with a Lookup table implementation, or a CORDIC implementation, which requires only iterative shift-add operations.

Re-implement Division Operations Where Possible

Often, division is not fully supported by hardware and can result in slow processing. When your algorithm requires a division, consider replacing it with one of the following options:

- Use bit shifting when the denominator is a power of two. For example, `bitsra(x,3)` instead of $x/8$.
- Multiply by the inverse when the denominator is constant. For example, $x*0.2$ instead of $x/5$.

Eliminate Floating-Point Variables

For more efficient code, eliminate floating-point variables. The one exception to this is loop indices because they usually become integer types.

Fixed-Point Design Exploration in Parallel

This example shows how to explore and test fixed-point designs by distributing tests across many computers in parallel. The example uses a `parfor` loop to test the accuracy of a QRS detector algorithm.

Running parallel simulations requires a Parallel Computing Toolbox™ license.

Using Parallel for-Loops for Design Exploration

Like a standard `for`-loop, a `parfor`-loop executes a series of statements over a range of values. Using the `parfor` command, you can set up a parallel `for`-loop in your code to explore fixed-point designs by distributing the tests across many computers. In a `parfor` loop, loop iterations execute in parallel which can provide better performance than standard `for`-loops.

The `test_heart_rate_detector_in_parallel` script sets up the system under test and initializes the arrays that will contain the results outside of the `parfor`-loop. It then uses a `parfor` loop to test each record in parallel. The `parfor`-loop loads the data, runs the system, then classifies and saves the results in parallel. When the `parfor`-loop finishes, the script displays the results.

```

%% Run test of all records in a database in parallel
record_names = {'ecg_01','ecg_02','ecg_03','ecg_04','ecg_05','ecg_06',...
    'ecg_07','ecg_08','ecg_09','ecg_10','ecg_11','ecg_12','ecg_13'};

%% Set up the system under test
data_type = 'fixedwrap';
T = heart_rate_detector_types(data_type);
[mex_function_name,Fs_target] = setup_heart_rate_detector(record_names,data_type,T);

%% Initialize array to contain results
results_file_names = cell(size(record_names));

%% Test each record in the database in parallel
parfor record_number = 1:length(record_names);
    % Load data
    record_name = record_names{record_number};
    [ecg,tm,ann,Fs] = load_ecg_data(record_name,Fs_target);

    % Run system under test
    detector_outputs = run_heart_rate_detector(mex_function_name,ecg,T);

    % Classify results
    [qrs_struct,qrs_stats] = classify_qrs(ann, Fs, detector_outputs);

    % Save results
    results_file_names{record_number} = save_heart_rate_data(...
        mex_function_name,record_name,...
        data_type,ecg,tm,ann,Fs,...
        detector_outputs,...
        qrs_struct,qrs_stats);

end

%% Display results
display_ecg_results(record_names, results_file_names);

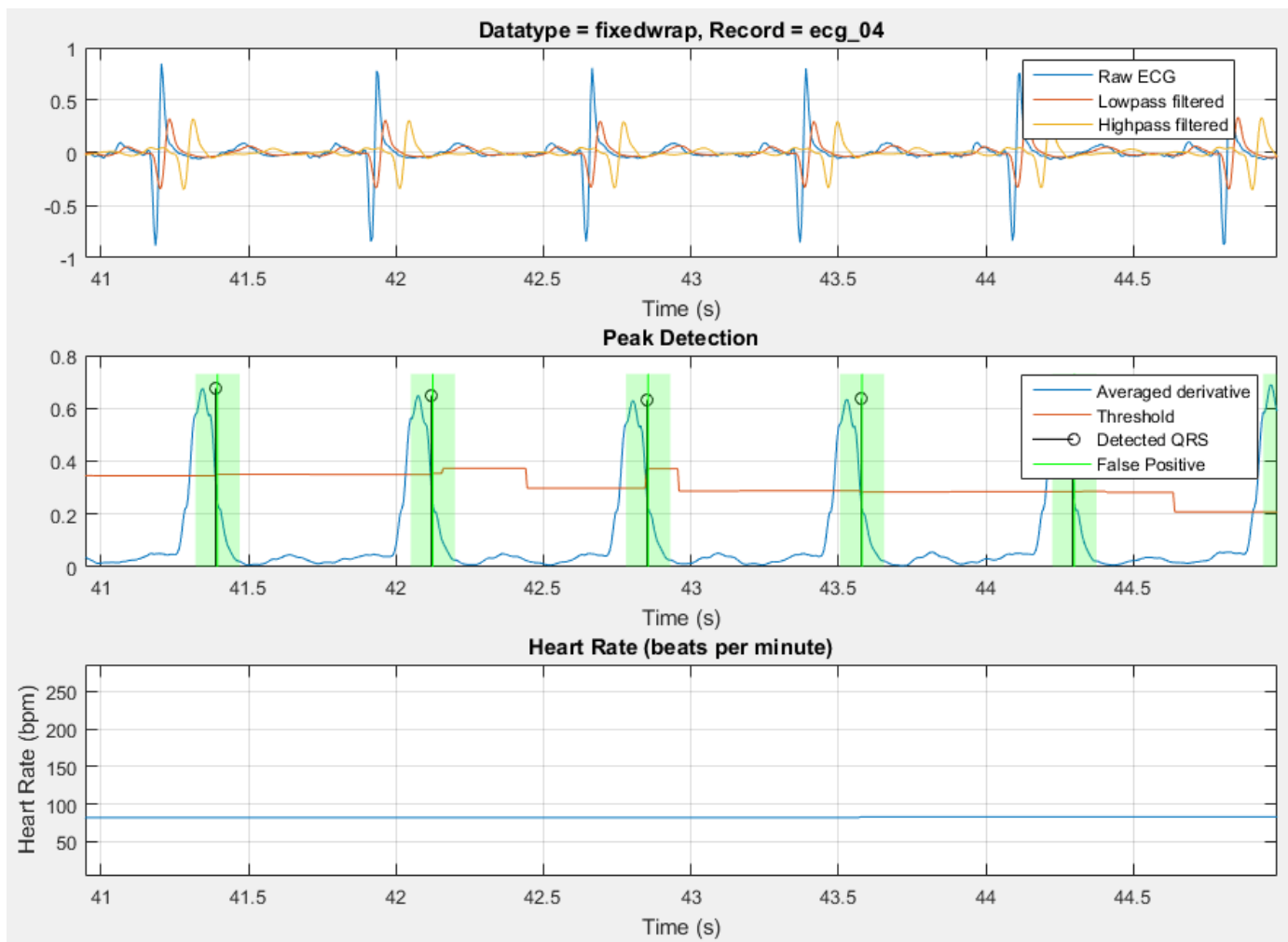
```

Because loop iterations must be completely independent of each other, you cannot call the `save` and `load` commands directly inside a `parfor`-loop. You can, however, call functions that call these commands. In this example, the functions `load_ecg_data` and `save_heart_rate_data` load and save the necessary data.

Description of System Under Test

The system under test in this example tests a simple QRS detector that measures the time difference between QRS detections to compute heart rate. The `test_heart_rate_detector_in_parallel` script passes ECG recordings to the detection algorithm.

The following plot is an example when the detector algorithm correctly identifies the QRS detections to compute the heartbeat.



The detection algorithm is simplified for this example. Examining the plots and results that are displayed when the example runs shows that the algorithm is not always very accurate.

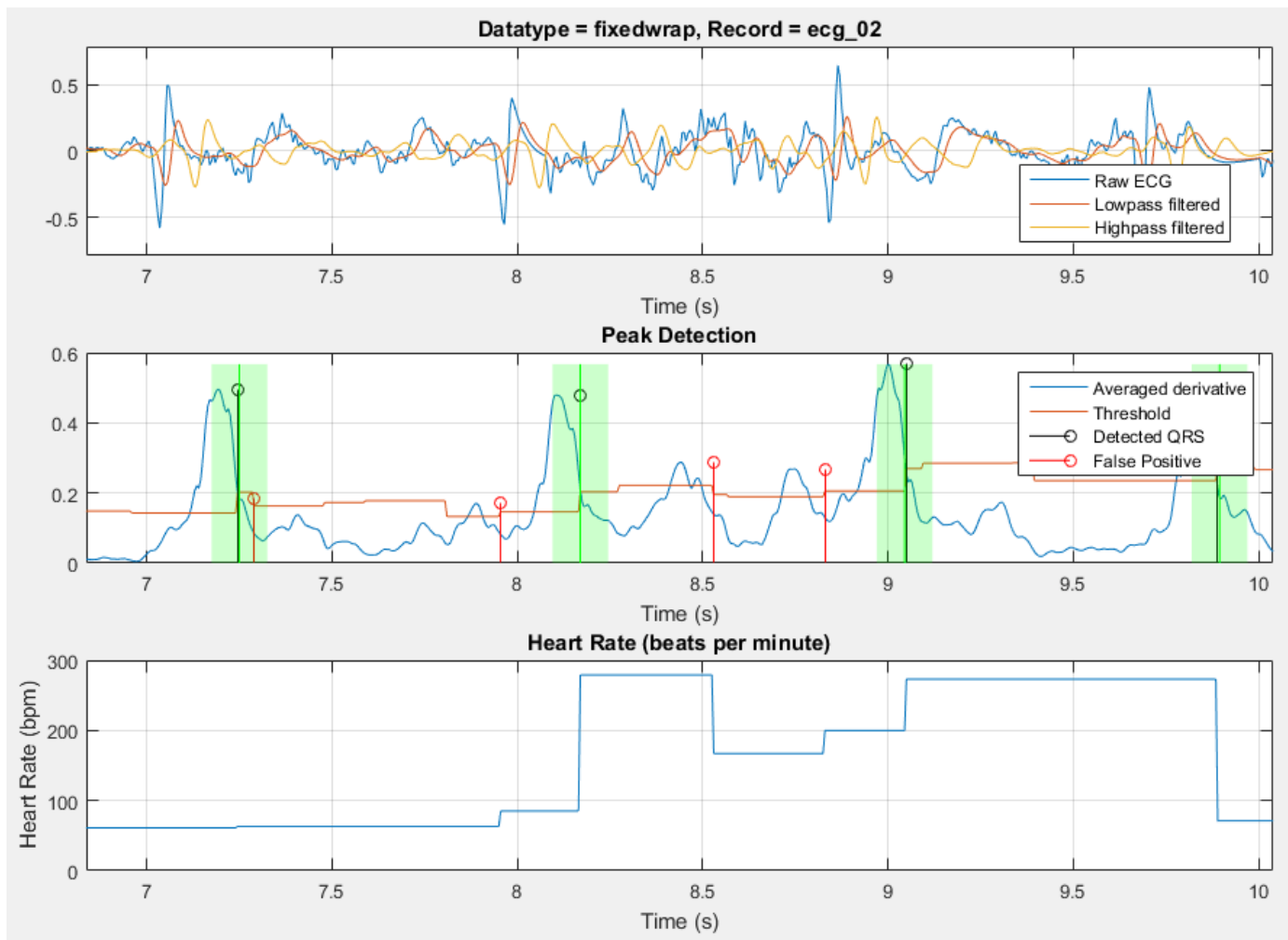
| Record | #QRS | TP | FP | FN | DER | Se | +P |
|--------|------|-----|----|----|-------|--------|-------|
| ecg_01 | 253 | 195 | 1 | 58 | 23.32 | 77.08 | 99.49 |
| ecg_02 | 133 | 133 | 18 | 0 | 13.53 | 100.00 | 88.08 |
| ecg_03 | 94 | 94 | 1 | 0 | 1.06 | 100.00 | 98.95 |

| | | | | | | | |
|--------|------|------|-----|----|-------|--------|--------|
| ecg_04 | 92 | 91 | 0 | 1 | 1.09 | 98.91 | 100.00 |
| ecg_05 | 93 | 91 | 1 | 2 | 3.23 | 97.85 | 98.91 |
| ecg_06 | 131 | 131 | 22 | 0 | 16.79 | 100.00 | 85.62 |
| ecg_07 | 174 | 173 | 2 | 0 | 1.15 | 100.00 | 98.86 |
| ecg_08 | 117 | 116 | 10 | 1 | 9.40 | 99.15 | 92.06 |
| ecg_09 | 137 | 137 | 1 | 0 | 0.73 | 100.00 | 99.28 |
| ecg_10 | 96 | 96 | 3 | 0 | 3.12 | 100.00 | 96.97 |
| ecg_11 | 73 | 73 | 1 | 0 | 1.37 | 100.00 | 98.65 |
| ecg_12 | 146 | 145 | 71 | 0 | 48.63 | 100.00 | 67.13 |
| ecg_13 | 144 | 144 | 5 | 0 | 3.47 | 100.00 | 96.64 |
| Totals | 1683 | 1619 | 136 | 62 | 11.76 | 96.31 | 92.25 |

Legend:

#QRS: Total number of QRS Complexes
 TP: Number of true positive
 FP: Number of false positive
 FN: Number of false negative
 DER: Detection error rate in percent
 Se: Sensitivity in percent
 +P: Positive prediction in percent

The following plot is an example when the detector algorithm identifies false positives due to noise in the recording.



All ECG recordings used in this example were measured on hobbyist equipment. You can use the PhysioNet database of recorded physiological signals to do a similar analysis on your own. The annotations on these recordings were not verified by doctors.

Run the Example

Run the `test_heart_rate_detector_in_parallel` example script.

```
test_heart_rate_detector_in_parallel
```

References

- [1] Patrick S. Hamilton, Open Source ECG Analysis Software (OSEA), E.P. Limited, Somerville, MA, <http://www.eplimited.com>, 2002.
- [2] Gari D Clifford, Francisco Azuaje, and Patrick E. McSharry. Advanced Methods and Tools for ECG Data Analysis, Artech House, 2006.
- [3] American National Standard ANSI/AAMI EC38:2007 Medical electrical equipment — Part 2-47: Particular requirements for the safety, including essential performance, of ambulatory electrocardiographic systems, Association for the Advancement of Medical Instrumentation, 2008.

[4] George B. Moody, "Evaluating ECG Analyzers", WaveForm DataBase Applications Guide, Harvard-MIT Division of Health Sciences and Technology, Cambridge, MA, WFDB10.5.23, 13 March 2014.

[5] Ida Laila binti Ahmad, Masnani binti Mohamed, Norul Ain binti Ab Ghani, "Development of a Concept Demonstrator for QRS Complex Detection using Combined Algorithms", 2012 IEEE EMBS International Conference on Biomedical Engineering and Sciences, Langkawi, 17th-19th December 2012.

[6] R. Harikumar, S.N. Shivappriya, "Analysis of QRS Detection Algorithm for Cardiac Abnormalities—A Review", International Journal of Soft Computing and Engineering (IJSCE), ISSN: 2231-2307, Volume-1, Issue-5, November 2011.

See Also

`for` | `parfor` | "Parallel for-Loops (`parfor`)" (Parallel Computing Toolbox)

Real-Time Image Acquisition, Image Processing, and Fixed-Point Blob Analysis for Target Practice Analysis

This example shows how to acquire real-time images from a GigE Vision® camera or webcam, process the images using fixed-point blob analysis, and determine world coordinates to score a laser pistol target.

The technology featured in this example is used in a wide range of applications, such as estimating distances to objects in front of a car [1], medical image analysis of cells [2], and detecting asteroids [3].

Key features of this example include:

- Fixed-point blob analysis for collecting measurements
- Real-time image acquisition
- Camera calibration to determine world coordinates of image points
- Correct images for lens distortion to ensure accuracy of collected measurements in world units
- Determine world coordinates of image points by mapping pixel locations to locations in real-world units

Required Products

This example uses these products for the algorithm:

- MATLAB®
- Fixed-Point Designer™
- Computer Vision Toolbox™
- Image Acquisition Toolbox™
- Image Processing Toolbox™

If you run the example in simulation mode, then you do not need a camera. Simulation mode loads recorded images and runs them through the algorithm as if a camera was attached.

If you use a GigE Vision camera, then you need this support package:

- Image Acquisition Toolbox™ Support Package for GigE Vision® Hardware

If you use a webcam, then you need this support package:

- MATLAB® Support Package for USB Webcams

Hardware Setup

Cameras

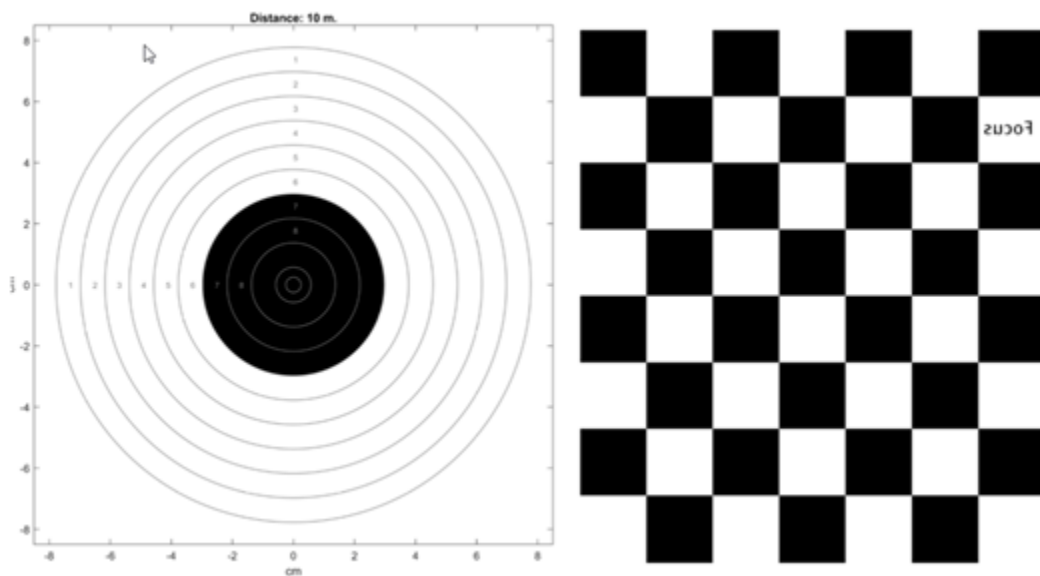
Image Acquisition Toolbox™ enables you to acquire images and video from cameras and frame grabbers directly into MATLAB® and Simulink®. Use the Image Acquisition Toolbox™ Support Package for GigE Vision® Hardware or the MATLAB® Support Package for USB Webcams to set up a camera to acquire the real-time images to perform the analysis.

For more information on setting up the camera, see “Setting Up Image Acquisition Hardware” (Image Acquisition Toolbox).

Target

Use these commands to create a target to print for use in this example. This code generates a postscript file that can be opened and printed double-sided, with the target on one side and the checkerboard for camera calibration on the other side.

```
distance_m = 10;  
offset_mm = 0;  
print_target = true;  
LaserTargetExample.make_target_airpistol10m(distance, offset_mm, print_target)
```



You can find an example target, `airpistoltarget_10m.pdf`, in the `+LaserTargetExample/targets_for_printing` folder.

Setup

Set up the camera so that it faces the checkerboard side of the target. The shooter faces the target. You can keep the target and camera in fixed positions by mounting them on a board.



Algorithm

Calibrate the Image

Camera calibration is the process of estimating the parameters of the lens and the image sensor. These parameters measure objects captured by the camera. Use the Camera Calibrator (Computer Vision Toolbox) app to detect the checkerboard pattern on the back of the target and remove any distortion. Determine the threshold of ambient light on the target. You may need to adjust the camera settings or the lighting so that the image is not saturated. Use the `pointsToWorld` (Computer Vision Toolbox) function to determine world coordinates of the image points.

For more information, see “What Is Camera Calibration?” (Computer Vision Toolbox).

Find and Score the Shot

The algorithm scores the shots by detecting the bright light of the laser pistol. While shooting, get a frame and detect if there is a bright spot. If there is a bright spot over the specified threshold, process that frame.

Use blob analysis to find the center of the bright spot and translate the location from pixel coordinates to world coordinates. The blob analysis is done in fixed point because the image is stored as an 8-bit signed integer. After finding the center of the bright spot in world coordinates, calculate its distance from the bullseye at the origin and assign a point value to the shot.

Run the Example

To start the simulation, execute the run script stored in the +LaserTargetExample folder.

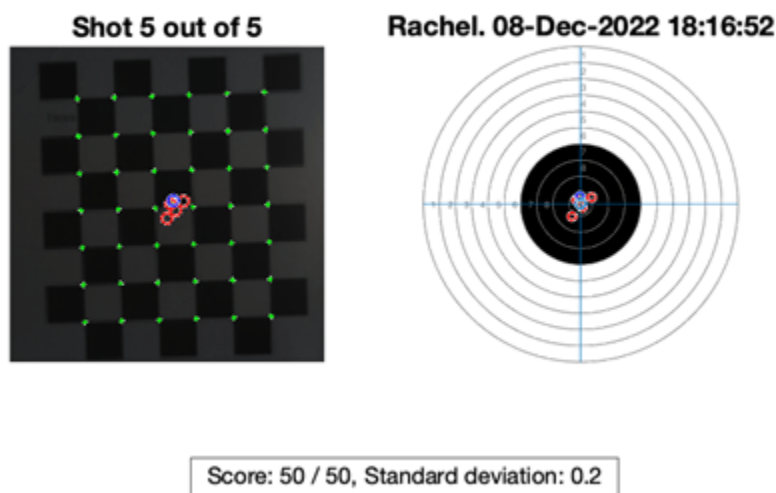
```
LaserTargetExample.run
```

This prompt appears:

```
(1) gigecam
(2) webcam
(3) simulation
Enter the number of the source type: 3
```

The script prompts you to select the source to use for the simulation. Enter 3 to watch a simulation of a previously recorded session. There is one previously recorded simulation available in the example files. A simulation recording is saved each time you run the example using a live camera and will be added to the simulation list. Enter a number to begin the simulation.

```
(1) saved_shots_20221208T181652
Enter number of file from list: 1
```

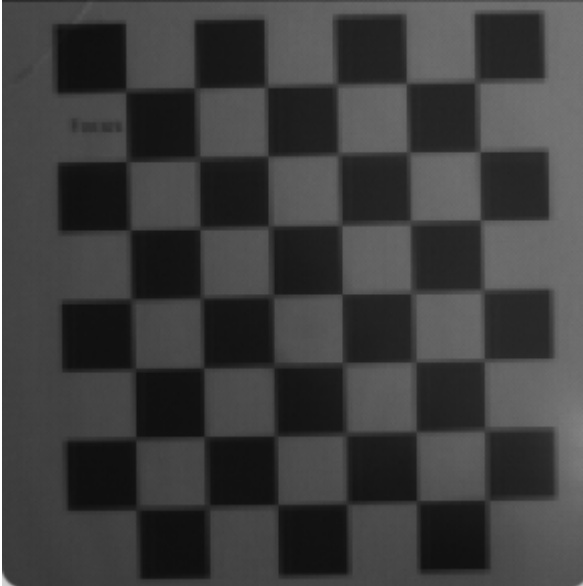


Entering a value of 1 or 2 prompts you to set up a GigE Vision camera or a webcam. It may take a couple of minutes for the camera to be found. When the camera is found, a **Deployable Video Player** window appears.

Adjust the camera so that:

- All of the squares on the back of the target are in view.
- The target is in focus.
- There are no bright spots that can be confused with a laser hit.

It does not matter if the image is not straight. The algorithm corrects for this. Close the **Deployable Video Player** to continue.



The example then prompts you to enter the distance from the shooter to the target in meters and the name of the shooter. After shooting five shots at the target with a laser pistol, you will be prompted to enter the name of another shooter or empty return to quit.

Use a Different Camera

To set up the example using your own camera, use the Camera Calibrator (Computer Vision Toolbox) app to detect the checkerboard on the back of the target and remove distortion. Save the calibration variables in a MAT-file. The calibration variables for the GigE Vision camera and a webcam are saved in these MAT-files:

- +LaserTargetExample/gigecam_240x240_calibration_parameters.mat
- +LaserTargetExample/webcam_LifeCam_480x480_camera_parameters.mat

Edit one of these files by substituting the settings with appropriate values for your camera:

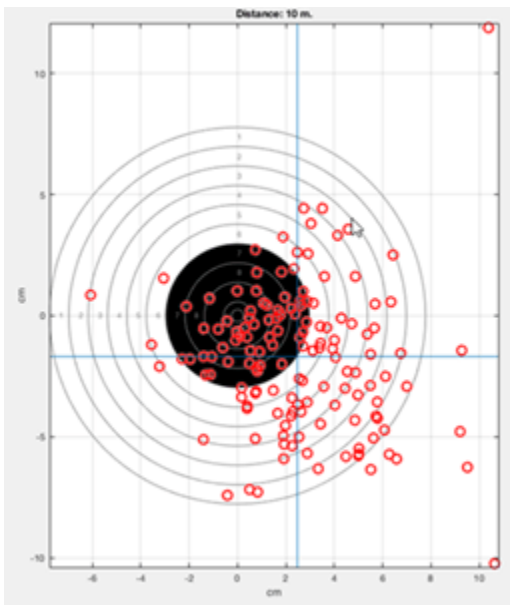
- +LaserTargetExample/gigecam_setup.m
- +LaserTargetExample/webcam_setup.m

Explore Data

Shot Database

Each time you shoot, the hits are recorded in a file named `ShotDatabase.csv`. You can use the `readtable` function to load the data into a table object to visualize it. For example, after shooting, which populates the `ShotDatabase.csv` file, this code plots the center of a group of many shots:

```
T = readtable('ShotDatabase.csv');
LaserTargetExample.make_target_airpistol10m;
LaserTargetExample.plot_shot_points(T.X, T.Y);
ax = gca;
line(mean(T.X)*[1,1], ax.YLim);
line(ax.XLim, mean(T.Y)*[1,1]);
grid on;
```



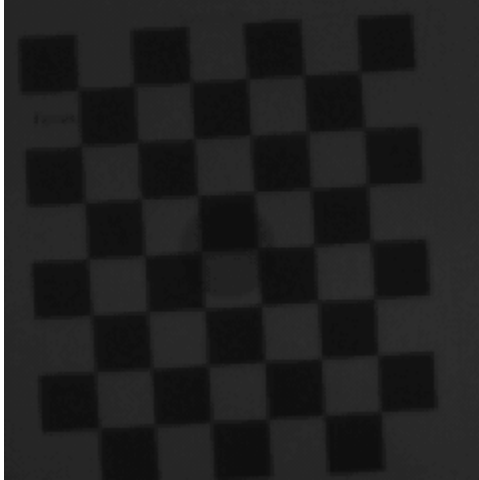
Simulation Recordings

Each time you shoot, the video frames in which shots were detected are stored in files in a folder named `simulation_recordings`. You can load these files and explore the raw data from the shots. You can also edit the algorithm.

The variable `frames` contains the first frame which was used for calibration, plus ten frames for each detected shot. The first frame in each run of ten is where a shot was detected. You can see your hand movement in the subsequent frames. You can make a short animation of the data using this code:

```
d = dir(fullfile('simulation_recordings', '*.mat'));
record = load(fullfile(d(1).folder, d(1).name));
t = LaserTargetExample.SerialDateNumber_to_seconds(...
    record.times);
t = t-t(1);
figure
for k = 1:size(record.frames, 3)
    imshow(record.frames(:,:,k), ...
        'InitialMagnification','fit');
```

```
        title(sprintf('Time since beginning of round: %.3f seconds',...  
                    t(k)))  
        drawnow  
end
```



References

[1] “Tracking Cars Using Foreground Detection” (Computer Vision Toolbox)

[2] “Cell Counting” (Computer Vision Toolbox)

[3] Rizza, Antonio, Felice Piccolo, Mattia Pugliatti, Paolo Panicucci, and Francesco Toppoto. "Hardware-in-the-Loop Simulation Framework for CubeSats Proximity Operations: Application to the Milani Mission." *73rd International Astronautical Congress, Paris, France, September 18-22, 2022*.

See Also

`undistortImage` | `pointsToWorld` | `vision.BlobAnalysis` | `detectCheckerboardPoints`

More About

- “Acquire Images from GigE Vision Cameras” (Image Acquisition Toolbox)
- “Install the MATLAB Support Package for USB Webcams” (Image Acquisition Toolbox)

Code Acceleration and Code Generation from MATLAB for Fixed-Point Algorithms

- “Code Acceleration and Code Generation from MATLAB” on page 12-2
- “Requirements for Generating Compiled C Code Files” on page 12-3
- “Functions Supported for Code Acceleration or C Code Generation” on page 12-4
- “Workflow for Fixed-Point Code Acceleration and Generation” on page 12-5
- “Accelerate Code Using fiaccel” on page 12-6
- “File Infrastructure and Paths Setup” on page 12-11
- “Detect and Debug Code Generation Errors” on page 12-14
- “Set Up C Compiler and Compilation Options” on page 12-16
- “MEX Configuration Dialog Box Options” on page 12-18
- “Specify Configuration Parameters in Command-Line Workflow Interactively” on page 12-21
- “Best Practices for Accelerating Fixed-Point Code” on page 12-24
- “Code Generation Reports” on page 12-27
- “Generate C Code from Code Containing Global Data” on page 12-31
- “Define Input Properties Programmatically in MATLAB File” on page 12-35
- “Specify Cell Array Inputs at the Command Line” on page 12-42
- “Specify Global Cell Arrays at the Command Line” on page 12-47
- “Control Run-Time Checks” on page 12-48
- “Fix Run-Time Stack Overflows” on page 12-50
- “Code Generation with MATLAB Coder” on page 12-51
- “Fixed-Point FIR Code Example Parameter Values” on page 12-52
- “Accelerate Code for Variable-Size Data” on page 12-54
- “Code Generation Readiness Tool” on page 12-61
- “Check Code Using the Code Generation Readiness Tool” on page 12-64
- “Check Code Using the MATLAB Code Analyzer” on page 12-65
- “Fix Errors Detected at Code Generation Time” on page 12-66
- “Avoid Multiword Operations in Generated Code” on page 12-67
- “Find Potential Data Type Issues in Generated Code” on page 12-69

Code Acceleration and Code Generation from MATLAB

In many cases, you may want your code to run faster and more efficiently. Code acceleration provides optimizations for accelerating fixed-point algorithms through MEX file building. In Fixed-Point Designer the `fiaccel` function converts your MATLAB code to a MEX function and can greatly accelerate the execution speed of your fixed-point algorithms.

Code generation creates efficient, production-quality C/C++ code for desktop and embedded applications. There are several ways to use Fixed-Point Designer software to generate C/C++ code.

| Use... | To... | Requires... | See... |
|---------------------------------|---|---|--|
| MATLAB Coder (codegen) function | Automatically convert MATLAB code to C/C++ code | MATLAB Coder code generation software license | "Generate C Code at the Command Line" (MATLAB Coder) |
| MATLAB Function | Use MATLAB code in your Simulink models that generate embeddable C/C++ code | Simulink license | "Implement MATLAB Functions in Simulink with MATLAB Function Blocks" |

MATLAB code generation supports variable-size arrays and matrices with known upper bounds. To learn more about using variable-size signals, see "Code Generation for Variable-Size Arrays" on page 29-2.

Requirements for Generating Compiled C Code Files

You use the `fiaccl` function to generate MEX code from a MATLAB algorithm. The algorithm must meet these requirements:

- Must be a MATLAB function, not a script
- Must meet the requirements listed on the `fiaccl` reference page
- Does not call custom C code using any of the following MATLAB Coder constructs:
 - `coder.ceval`
 - `coder.ref`
 - `coder.rref`
 - `coder.wref`

Functions Supported for Code Acceleration or C Code Generation

The following general limitations apply to the use of Fixed-Point Designer functions in generated code, with `fiaccel`:

- `fipref` and `quantizer` objects are not supported.
- Word lengths greater than 128 bits are not supported.
- You cannot change the `fi` `fimath` or `numericType` of a given `fi` variable after that variable has been created.
- The `boolean` value of the `DataTypeMode` and `DataType` properties are not supported.
- For all `SumMode` property settings other than `FullPrecision`, the `CastBeforeSum` property must be set to `true`.
- You can use parallel for (`parfor`) loops in code compiled with `fiaccel`, but those loops are treated like regular `for` loops.
- When you compile code containing `fi` objects with nontrivial slope and bias scaling, you may see different results in generated code than you achieve by running the same code in MATLAB.

To view a list of the Fixed-Point Designer functions that are supported for code acceleration or C/C++ code generation, refer to the **Fixed-Point Designer** category of these tables:

- Functions and Objects Supported for C/C++ Code Generation (Category List)
- Functions and Objects Supported for C/C++ Code Generation (Alphabetical List)

In these tables, a **i** icon before the name of a function indicates that there are specific usage notes and limitations related to code acceleration or code generation for that function. To view these usage notes and limitations, in the corresponding reference page, scroll down to the **Extended Capabilities** section at the bottom and expand the **C/C++ Code Generation** section.

Workflow for Fixed-Point Code Acceleration and Generation

| Step | Action | Details |
|------|---|--|
| 1 | Set up your C compiler. | See “Set Up C Compiler” on page 12-16. |
| 2 | Set up your file infrastructure. | See “File Infrastructure and Paths Setup” on page 12-11. |
| 3 | Make your MATLAB algorithm suitable for code generation | See “Best Practices for Accelerating Fixed-Point Code” on page 12-24. |
| 4 | Set compilation options. | See “Set Up C Compiler and Compilation Options” on page 12-16. |
| 5 | Specify properties of primary function inputs. | See “Specify Properties of Entry-Point Function Inputs” on page 31-2. |
| 6 | Run <code>fiaccel</code> with the appropriate command-line options. | See “Recommended Compilation Options for <code>fiaccel</code> ” on page 12-24. |

Accelerate Code Using `fiaccl`

In this section...

“Speeding Up Fixed-Point Execution with `fiaccl`” on page 12-6

“Running `fiaccl`” on page 12-6

“Generated Files and Locations” on page 12-6

“Data Type Override Using `fiaccl`” on page 12-9

“Specifying Default `fimath` Values for MEX Functions” on page 12-9

Speeding Up Fixed-Point Execution with `fiaccl`

You can convert fixed-point MATLAB code to MEX functions using `fiaccl`. The generated MEX functions contain optimizations to automatically accelerate fixed-point algorithms to compiled C/C++ code speed in MATLAB. The `fiaccl` function can greatly increase the execution speed of your algorithms.

Running `fiaccl`

The basic command is:

```
fiaccl M_fcn
```

By default, `fiaccl` performs the following actions:

- Searches for the function `M_fcn` stored in the file `M_fcn.m` as specified in “Compile Path Search Order” on page 12-11.
- Compiles `M_fcn` to MEX code.
- If there are no errors or warnings, generates a platform-specific MEX file in the current folder, using the naming conventions described in “File Naming Conventions” on page 12-26.
- If there are errors, does not generate a MEX file, but produces an error report in a default output folder, as described in “Generated Files and Locations” on page 12-6.
- If there are warnings, but no errors, generates a platform-specific MEX file in the current folder, but does report the warnings.

You can modify this default behavior by specifying one or more compiler options with `fiaccl`, separated by spaces on the command line.

Generated Files and Locations

`fiaccl` generates files in the following locations:

| Generates: | In: |
|-----------------------------|----------------|
| Platform-specific MEX files | Current folder |

| Generates: | In: |
|---|---|
| code generation reports (if errors or warnings occur during compilation) | Default output folder: fiaccel/mex/M_fcn_name/html |

You can change the name and location of generated files by using the options `-o` and `-d` when you run `fiaccel`.

In this example, you will use the `fiaccel` function to compile different parts of a simple algorithm. By comparing the run times of the two cases, you will see the benefits and best use of the `fiaccel` function.

Comparing Run Times When Accelerating Different Algorithm Parts

The algorithm used throughout this example replicates the functionality of the MATLAB `sum` function, which sums the columns of a matrix. To see the algorithm, type `open fi_matrix_column_sum.m` at the MATLAB command line.

```
function B = fi_matrix_column_sum(A)
% Sum the columns of matrix A.
%#codegen
[m,n] = size(A);
w = get(A,'WordLength') + ceil(log2(m));
f = get(A,'FractionLength');
B = fi(zeros(1,n),true,w,f);
for j = 1:n
    for i = 1:m
        B(j) = B(j) + A(i,j);
    end
end
end
```

Trial 1: Best Performance

The best way to speed up the execution of the algorithm is to compile the entire algorithm using the `fiaccel` function. To evaluate the performance improvement provided by the `fiaccel` function when the entire algorithm is compiled, run the following code.

The first portion of code executes the algorithm using only MATLAB functions. The second portion of the code compiles the entire algorithm using the `fiaccel` function. The MATLAB `tic` and `toc` functions keep track of the run times for each method of execution.

```
% MATLAB
fipref('NumericTypeDisplay','short');
A = fi(randn(1000,10));
tic
B = fi_matrix_column_sum(A)
t_matrix_column_sum_m = toc

% fiaccel
fiaccel fi_matrix_column_sum -args {A} ...
-I [matlabroot '/toolbox/fixedpoint/fidemos']
tic
B = fi_matrix_column_sum_mex(A);
t_matrix_column_sum_mex = toc
```

Trial 2: Worst Performance

Compiling only the smallest unit of computation using the `fiaccel` function leads to much slower execution. In some cases, the overhead that results from calling the `mex` function inside a nested loop can cause even slower execution than using MATLAB functions alone. To evaluate the performance of the `mex` function when only the smallest unit of computation is compiled, run the following code.

The first portion of code executes the algorithm using only MATLAB functions. The second portion of the code compiles the smallest unit of computation with the `fiaccel` function, leaving the rest of the computations to MATLAB.

```
% MATLAB
tic
[m,n] = size(A);
w = get(A,'WordLength') + ceil(log2(m));
f = get(A,'FractionLength');
B = fi(zeros(1,n),true,w,f);
for j = 1:n
    for i = 1:m
        B(j) = fi_scalar_sum(B(j),A(i,j));
        % B(j) = B(j) + A(i,j);
    end
end
t_scalar_sum_m = toc

% fiaccel
fiaccel fi_scalar_sum -args {B(1),A(1,1)} ...
-I [matlabroot '/toolbox/fixedpoint/fidemos']
tic
[m,n] = size(A);
w = get(A,'WordLength') + ceil(log2(m));
f = get(A,'FractionLength');
B = fi(zeros(1,n),true,w,f);
for j = 1:n
    for i = 1:m
        B(j) = fi_scalar_sum_mex(B(j),A(i,j));
        % B(j) = B(j) + A(i,j);
    end
end
t_scalar_sum_mex = toc
```

Ratio of Times

A comparison of Trial 1 and Trial 2 appears in the following table. Your computer may record different times than the ones the table shows, but the ratios should be approximately the same. There is an extreme difference in ratios between the trial where the entire algorithm was compiled using `fiaccel` (`t_matrix_column_sum_mex.m`) and where only the scalar sum was compiled (`t_scalar_sum_mex.m`). Even the file with no `fiaccel` compilation (`t_matrix_column_sum_m`) did better than when only the smallest unit of computation was compiled using `fiaccel` (`t_scalar_sum_mex`).

| X (Overall Performance Rank) | Time | X/Best | X _m /X _{mex} |
|----------------------------------|---------|---------|----------------------------------|
| Trial 1: Best Performance | | | |
| t_matrix_column_sum_m (2) | 1.99759 | 84.4917 | 84.4917 |

| X (Overall Performance Rank) | Time | X/Best | X_m/X_mex |
|-----------------------------------|-----------|---------|-----------|
| t_matrix_column_sum_mex (1) | 0.0236424 | 1 | |
| Trial 2: Worst Performance | | | |
| t_scalar_sum_m (4) | 10.2067 | 431.71 | 2.08017 |
| t_scalar_sum_mex (3) | 4.90664 | 207.536 | |

Data Type Override Using `fiaccel`

Fixed-Point Designer software ships with an example of how to generate a MEX function from MATLAB code. The code in the example takes the weighted average of a signal to create a lowpass filter. To run the example in the Help browser select **MATLAB Examples** under Fixed-Point Designer, and then select Fixed-Point Lowpass Filtering Using MATLAB for Code Generation.

You can specify data type override in this example by typing an extra command at the MATLAB prompt in the “Define Fixed-Point Parameters” section of the example. To turn data type override on, type the following command at the MATLAB prompt after running the `reset(fipref)` command in that section:

```
fipref('DataTypeOverride','TrueDoubles')
```

This command tells Fixed-Point Designer software to create all `fi` objects with type `fi double`. When you compile the code using the `fiaccel` command in the “Compile the M-File into a MEX File” section of the example, the resulting MEX-function uses floating-point data.

Specifying Default `fimath` Values for MEX Functions

MEX functions generated with `fiaccel` use the MATLAB default global `fimath`. The MATLAB factory default global `fimath` has the following properties:

```
RoundingMethod: Nearest
OverflowAction: Saturate
ProductMode: FullPrecision
SumMode: FullPrecision
```

When running MEX functions that depend on the MATLAB default `fimath` value, do not change this value during your MATLAB session. Otherwise, MATLAB generates a warning, alerting you to a mismatch between the compile-time and run-time `fimath` values. For example, create the following MATLAB function:

```
function y = test %#codegen
y = fi(0);
```

The function `test` constructs a `fi` object without explicitly specifying a `fimath` object. Therefore, `test` relies on the default `fimath` object in effect at compile time.

Generate the MEX function `test_mex` to use the factory setting of the MATLAB default `fimath`.

```
resetglobalfimath;
fiaccel test
```

`fiaccel` generates a MEX function, `test_mex`, in the current folder.

Run `test_mex`.

```
test_mex
ans =
    0
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 15
```

Modify the MATLAB default `fimath` value so it no longer matches the setting used at compile time.

```
F = fimath('RoundingMethod','Floor');
globalfimath(F);
```

Clear the MEX function from memory and rerun it.

```
clear test_mex
test_mex
```

The mismatch is detected and MATLAB generates a warning.

```
testglobalfimath_mex
Warning: This function was generated with a
different default fimath than the current default.
ans =
    0
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 15
```

To avoid this issue, separate the `fimath` properties from your algorithm by using types tables. For more information, see “Separate Data Type Definitions from Algorithm” on page 11-6.

File Infrastructure and Paths Setup

In this section...

“Compile Path Search Order” on page 12-11

“Naming Conventions” on page 12-11

Compile Path Search Order

`fiaccl` resolves function calls by searching first on the code generation path and then on the MATLAB path. By default, `fiaccl` tries to compile and generate code for functions it finds on the path unless you explicitly declare the function to be extrinsic. An extrinsic function is a function on the MATLAB path that is dispatched to MATLAB software for execution. `fiaccl` does not compile extrinsic functions, but rather dispatches them to MATLAB for execution.

Naming Conventions

MATLAB enforces naming conventions for functions and generated files.

- “Reserved Prefixes” on page 12-11
- “Reserved Keywords” on page 12-11
- “Conventions for Naming Generated files” on page 12-13

Reserved Prefixes

MATLAB reserves the prefix `eml` for global C functions and variables in generated code. For example, run-time library function names all begin with the prefix `emlrt`, such as `emlrtCallMATLAB`. To avoid naming conflicts, do not name C functions or primary MATLAB functions with the prefix `eml`.

Reserved Keywords

- “C Reserved Keywords” on page 12-11
- “C++ Reserved Keywords” on page 12-12
- “Reserved Keywords for Code Generation” on page 12-13

MATLAB Coder software reserves certain words for its own use as keywords of the generated code language. MATLAB Coder keywords on page 12-13 are reserved for use internal to MATLAB Coder software and should not be used in MATLAB code as identifiers or function names. C reserved keywords on page 12-11 should also not be used in MATLAB code as identifiers or function names. If your MATLAB code contains reserved keywords that the code generator cannot rename, the code generation build does not complete and an error message is displayed. To address this error, modify your code to use identifiers or names that are not reserved.

If you are generating C++ code using the MATLAB Coder software, in addition, your MATLAB code must not contain the “C++ Reserved Keywords” on page 12-12.

C Reserved Keywords

| | | | |
|---------------------|---------------------|---------------------|---------------------|
| <code>assert</code> | <code>extern</code> | <code>setjmp</code> | <code>string</code> |
| <code>auto</code> | <code>fenv</code> | <code>short</code> | <code>struct</code> |

| | | | |
|----------|----------|-------------|----------|
| break | float | signal | switch |
| case | for | signed | tgmath |
| char | goto | sizeof | threads |
| const | if | static | time |
| complex | int | stdalign | typedef |
| continue | inttypes | stdarg | uchar |
| ctype | iso646 | stdatomic | union |
| default | limits | stdbool | unsigned |
| do | locale | stddef | void |
| double | long | stdint | volatile |
| else | math | stdio | wchar |
| enum | register | stdlib | wctype |
| errno | return | stdnoreturn | while |

C++ Reserved Keywords

| | | | |
|--------------------|--------------|-----------------|---------------|
| algorithm | cstdint | iostream | sstream |
| any | cstdint | istream | stack |
| array | cstdint | iterator | static_cast |
| atomic | cstdint | limits | stdexcept |
| bitset | cstring | list | streambuf |
| cassert | ctgmth | locale | string_view |
| catch | ctime | map | stringstream |
| ccomplex | cuchar | memory | system_error |
| cctype | cwchar | memory_resource | template |
| cerrno | cwctype | mutable | this |
| cfenv | delete | mutex | thread |
| cfloat | deque | namespace | throw |
| chrono | dynamic_cast | new | try |
| cinttypes | exception | numeric | tuple |
| ciso646 | execution | operator | typeid |
| class | explicit | optional | type_traits |
| climits | export | ostream | typeid |
| locale | filesystem | private | typeinfo |
| cmath | forward_list | protected | typename |
| codecvt | friend | public | unordered_map |
| complex | fstream | queue | unordered_set |
| condition_variable | functional | random | using |
| const_cast | future | ratio | utility |

| | | | |
|-----------|------------------|------------------|----------|
| csetjmp | initializer_list | regex | valarray |
| csignal | inline | reinterpret_cast | vector |
| cstdalign | iomanip | scoped_allocator | virtual |
| cstdarg | ios | set | wchar_t |
| cstdbool | iosfwd | shared_mutex | |

Reserved Keywords for Code Generation

| | | | |
|-----------|-----------------------|-----------------------|-------------------|
| abs | fortran | localZCE | rtNaN |
| asm | HAVESTDIO | localZCSV | SeedFileBuffer |
| bool | id_t | matrix | SeedFileBufferLen |
| boolean_T | int_T | MODEL | single |
| byte_T | int8_T | MT | TID01EQ |
| char_T | int16_T | NCSTATES | time_T |
| cint8_T | int32_T | NULL | true |
| cint16_T | int64_T | NUMST | TRUE |
| cint32_T | INTEGER_CODE | pointer_T | uint_T |
| creal_T | LINK_DATA_BUFFER_SIZE | PROFILING_ENABLED | uint8_T |
| creal32_T | LINK_DATA_STREAM | PROFILING_NUM_SAMPLES | uint16_T |
| creal64_T | localB | real_T | uint32_T |
| cuint8_T | localC | real32_T | uint64_T |
| cuint16_T | localDWork | real64_T | UNUSED_PARAMETER |
| cuint32_T | localP | RT | USE_RTMODEL |
| ERT | localX | RT_MALLOC | VCAST_FLUSH_DATA |
| false | localXdis | rtInf | vector |
| FALSE | localXdot | rtMinusInf | |

Conventions for Naming Generated files

MATLAB provides platform-specific extensions for MEX files.

| Platform | MEX File Extension |
|-------------------|--------------------|
| Linux® x86-64 | .mexa64 |
| Windows® (32-bit) | .mexw32 |
| Windows x64 | .mexw64 |

Detect and Debug Code Generation Errors

In this section...

“Debugging Strategies” on page 12-14

“Error Detection at Design Time” on page 12-14

“Error Detection at Compile Time” on page 12-15

Debugging Strategies

To prepare your algorithms for code generation, MathWorks recommends that you choose a debugging strategy for detecting and correcting violations in your MATLAB applications, especially if they consist of a large number of MATLAB files that call each other's functions. Here are two best practices:

| Debugging Strategy | What to Do | Pros | Cons |
|------------------------|--|--|---|
| Bottom-up verification | <ol style="list-style-type: none"> 1 Verify that your lowest-level (leaf) functions are suitable for code generation. 2 Work your way up the function hierarchy incrementally to compile and verify each function, ending with the top-level function. | <ul style="list-style-type: none"> • Efficient • Safe • Easy to isolate syntax violations | Requires application tests that work from the bottom up |
| Top-down verification | <ol style="list-style-type: none"> 1 Declare all functions called by the top-level function to be extrinsic so <code>fiaccl</code> does not compile them. 2 Verify that your top-level function is suitable for code generation. 3 Work downward in the function hierarchy to: <ol style="list-style-type: none"> a. Remove extrinsic declarations one by one b. Compile and verify each function, ending with the leaf functions. | Lets you retain your top-level tests | Introduces extraneous code that you must remove after code verification, including: <ul style="list-style-type: none"> • Extrinsic declarations • Additional assignment statements as necessary to convert opaque values returned by extrinsic functions to nonopaque values. |

Error Detection at Design Time

To detect potential issues for MEX file building as you write your MATLAB algorithm, add the `%#codegen` directive to the code that you want `fiaccl` to compile. Adding this directive indicates that you intend to generate code from the algorithm and turns on detailed diagnostics during MATLAB code analysis.

Error Detection at Compile Time

Before you can successfully generate code from a MATLAB algorithm, you must verify that the algorithm does not contain syntax and semantics violations that would cause compile-time errors, as described in “Detect and Debug Code Generation Errors” on page 12-14.

`fiaccl` checks for all potential syntax violations at compile time. When `fiaccl` detects errors or warnings, it automatically produces a code generation report that describes the issues and provides links to the offending code. See “Code Generation Reports” on page 12-27.

If your MATLAB code calls functions on the MATLAB path, `fiaccl` attempts to compile these functions unless you declare them to be extrinsic.

Set Up C Compiler and Compilation Options

In this section...

“Set Up C Compiler” on page 12-16

“C Code Compiler Configuration Object” on page 12-16

“Compilation Options Modification at the Command Line Using Dot Notation” on page 12-16

“How fiaccel Resolves Conflicting Options” on page 12-17

Set Up C Compiler

Fixed-Point Designer automatically locates and uses a supported installed compiler. For the current list of supported compilers, see Supported and Compatible Compilers.

You can use `mex -setup` to change the default compiler. See “Change Default Compiler”.

C Code Compiler Configuration Object

For C code generation to a MEX file, MATLAB provides a configuration object `coder.mexconfig` for fine-tuning the compilation. To set MEX compilation options:

- 1 Define the compiler configuration object in the MATLAB workspace by issuing a constructor command:

```
comp_cfg = coder.mexconfig
```

MATLAB displays the list of compiler options and their current values in the command window.

- 2 Modify the compilation options as necessary. See “Compilation Options Modification at the Command Line Using Dot Notation” on page 12-16
- 3 Invoke `fiaccel` with the `-config` option and specify the configuration object as its argument:

```
fiaccel -config comp_cfg myMfile
```

The `-config` option instructs `fiaccel` to convert `myFile.m` to a MEX function, based on the compilation settings in `comp_cfg`.

Compilation Options Modification at the Command Line Using Dot Notation

Use dot notation to modify the value of compilation options, using this syntax:

```
configuration_object.property = value
```

Dot notation uses assignment statements to modify configuration object properties. For example, to change the maximum size function to inline and the stack size limit for inlined functions during MEX generation, enter this code at the command line:

```
co_cfg = coder.mexconfig
co_cfg.InlineThreshold = 25;
co_cfg.InlineStackLimit = 4096;
fiaccel -config co_cfg myFun
```


How fiaccel Resolves Conflicting Options

`fiaccel` takes the union of all options, including those specified using configuration objects, so that you can specify options in any order.

MEX Configuration Dialog Box Options

MEX Configuration Dialog Box Options

The following table describes parameters for fine-tuning the behavior of `fiaccl` for converting MATLAB files to MEX:

| Parameter | Equivalent Command-Line Property and Values (default in bold) | Description |
|-------------------------------------|--|---|
| Report | | |
| Create code generation report | GenerateReport true, false | Document generated code in a report. |
| Launch report automatically | LaunchReport true, false | Specify whether to automatically open report after code generation completes. |
| | | Note Requires that you enable Create code generation report |
| Debugging | | |
| Echo expressions without semicolons | EchoExpressions true , false | Specify whether or not actions that do not terminate with a semicolon appear in the MATLAB Command Window. |
| Enable debug build | EnableDebugging true, false | Compile the generated code in debug mode. |
| Language and Semantics | | |
| Constant Folding Timeout | ConstantFoldingTimeout <i>integer</i> , 10000 | Specify the maximum number of instructions to be executed by the constant folder. |
| Dynamic memory allocation | DynamicMemoryAllocation 'off' , 'AllVariableSizeArrays' | Enable dynamic memory allocation for variable-size data. By default, dynamic memory allocation is disabled and <code>fiaccl</code> allocates memory statically on the stack. When you select dynamic memory allocation, <code>fiaccl</code> allocates memory for all variable-size data dynamically on the heap. You <i>must</i> use dynamic memory allocation for all unbounded variable-size data. |
| Enable variable sizing | EnableVariableSizing true , false | Enable support for variable-size arrays. |

| Parameter | Equivalent Command-Line Property and Values (default in bold) | Description |
|----------------------------------|---|---|
| Extrinsic calls | ExtrinsicCalls true , false | <p>Allow calls to extrinsic functions.</p> <p>When enabled (true), the compiler generates code for the call to a MATLAB function, but does not generate the function's internal code.</p> <p>When disabled (false), the compiler ignores the extrinsic function. Does not generate code for the call to the MATLAB function—as long as the extrinsic function does not affect the output of the caller function. Otherwise, the compiler issues a compiler error.</p> |
| Global Data Synchronization Mode | GlobalDataSyncMethod <i>string</i> , SyncAlways , SyncAtEntryAndExits, NoSync | <p>Controls when global data is synchronized with the MATLAB global workspace. By default, (SyncAlways), synchronizes global data at MEX function entry and exit and for all extrinsic calls. This synchronization ensures maximum consistency between MATLAB and generated code. If the extrinsic calls do not affect global data, use this option with the <code>coder.extrinsic -sync:off</code> option to turn off synchronization for these calls.</p> <p>SyncAtEntryAndExits synchronizes global data at MEX function entry and exit only. If only a few extrinsic calls affect global data, use this option with the <code>coder.extrinsic -sync:on</code> option to turn on synchronization for these calls.</p> <p>NoSync disables synchronization. Ensure that your generated code does not interact with MATLAB before disabling synchronization. Otherwise, inconsistencies might occur.</p> |
| Saturate on integer overflow | SaturateOnIntegerOverflow true , false | Add checks in the generated code to detect integer overflow or underflow. |

| Parameter | Equivalent Command-Line Property and Values (default in bold) | Description |
|---|--|--|
| Safety (disable for faster MEX) | | |
| Ensure memory integrity | IntegrityChecks true , false | Detects violations of memory integrity while building MATLAB Function blocks and stops simulation with a diagnostic message. Setting IntegrityChecks to false also disables the run-time stack. |
| Ensure responsiveness | ResponsivenessChecks true , false | Enables responsiveness checks in code generated from MATLAB algorithms. |
| Function Inlining and Stack Allocation | | |
| Inline Stack Limit | InlineStackLimit <i>integer</i> , 4000 | Specify the stack size limit on inlined functions. |
| Inline Threshold | InlineThreshold <i>integer</i> , 10 | Specify the maximum size of functions to be inlined. |
| Inline Threshold Max | InlineThresholdMax <i>integer</i> , 200 | Specify the maximum size of functions after inlining. |
| Stack Usage Max | StackUsageMax <i>integer</i> , 200000 | Specify the maximum stack usage per application in bytes. Set a limit that is lower than the available stack size. Otherwise, a runtime stack overflow might occur. Overflows are detected and reported by the C compiler, not by <code>fiaccel</code> . |
| Optimizations | | |
| Use BLAS library if possible | EnableBLAS true , false | Speed up low-level matrix operations during simulation by calling the Basic Linear Algebra Subprograms (BLAS) library. |

See Also

- “Control Run-Time Checks” on page 12-48
- “Code Generation for Variable-Size Arrays” on page 29-2
- “Generate C Code from Code Containing Global Data” on page 12-31

Specify Configuration Parameters in Command-Line Workflow Interactively

After you have created a code generation configuration object at the command line, you can modify the properties of the object interactively by using the Configuration Parameters dialog box.

For more information on configuring the code generation process by using configuration objects, see “Configure Build Settings” (MATLAB Coder).

Create and Modify Configuration Objects by Using the Dialog Box

- 1 Create a configuration object as described in “Creating Configuration Objects” (MATLAB Coder).

For example, to create a `coder.MexCodeConfig` configuration object for MEX code generation:

```
mexcfg = coder.config('mex');
```

- 2 Open the property dialog box by using one of these methods:

- In the MATLAB workspace, double-click the configuration object variable.
- At the MATLAB command prompt, issue the `open` command, passing it the configuration object variable:

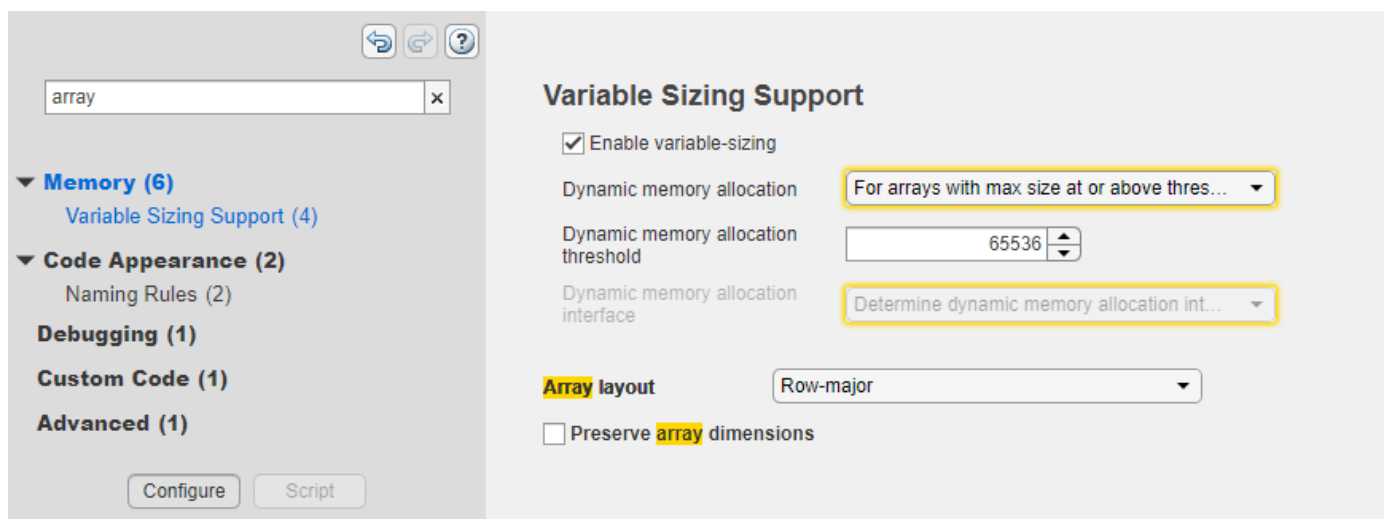
```
open mexcfg
```

- 3 In the dialog box, modify configuration parameters as required.

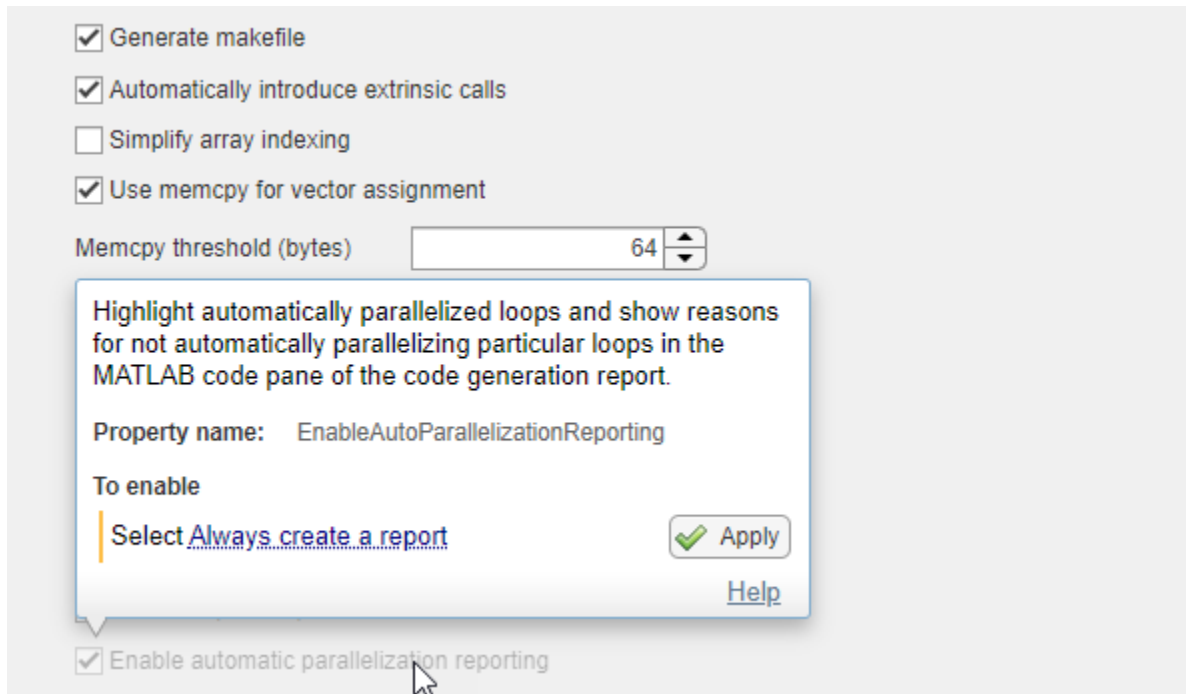
Additional Functionalities in the Dialog Box

To enable you to easily modify the configuration parameters in an interactive fashion, the Configuration Parameters dialog box provides these functionalities:

- *Search*: When you search for a string, you see the filtered results across all the settings categories. The search string might be present in a setting name, the name of an option for a setting, or in a tooltip.



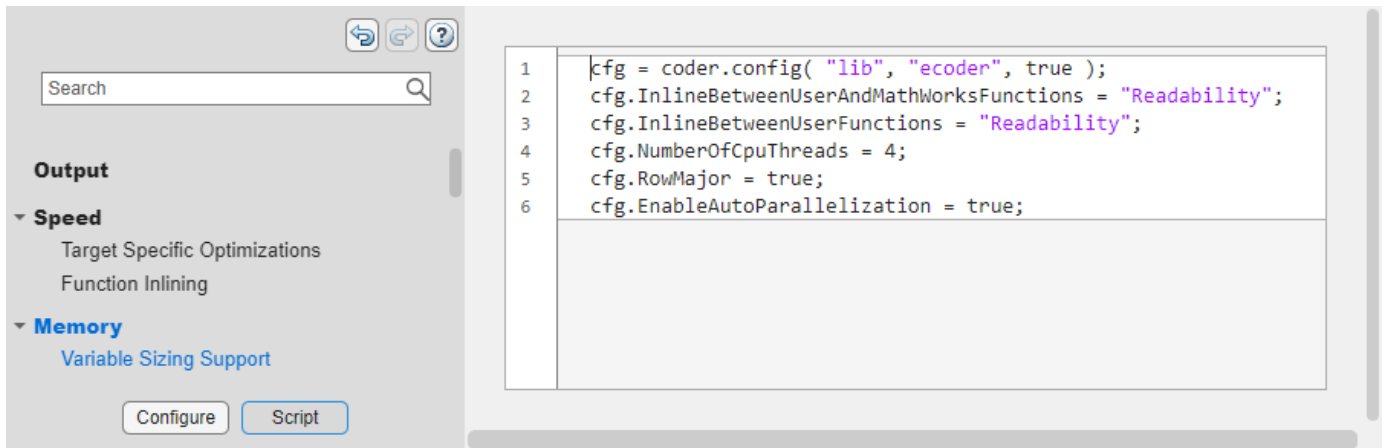
- *Informative tooltips*: The tooltip for each individual setting contains the corresponding configuration object property name, a **Help** link for that property, and the name of any additional product that using that property requires. If the property is disabled, the tooltip also contains links to other properties that you must set to enable this property. You can make that change in the tooltip itself.



- *Settings with nondefault values*: The dialog box shows settings that have nondefault values in bold font. To reset such a setting to its default values, click the **Reset** button in the tooltip.
- *MISRA Compliance pane*: If you have Embedded Coder®, the **MISRA Compliance** pane displays the settings that might impact MISRA™ compliance of the generated code. To set all of these settings to the recommended values, click **Set to Recommended Values**.

See “Generate C/C++ Code with Improved MISRA and AUTOSAR Compliance” (Embedded Coder).

- *Generate equivalent script*: You can view the command-line script that produces your current settings by clicking the **Script** button located at the bottom of the list of categories. You can switch from the script mode back to the interactive mode by clicking the **Configure** button.



See Also

[coder.MexCodeConfig](#) | [coder.CodeConfig](#) | [coder.EmbeddedCodeConfig](#) | [coder.DeepLearningCodeConfig](#)

More About

- “Configure Build Settings” (MATLAB Coder)

Best Practices for Accelerating Fixed-Point Code

In this section...

“Recommended Compilation Options for `fiaccl`” on page 12-24

“Build Scripts” on page 12-24

“Check Code Interactively Using MATLAB Code Analyzer” on page 12-25

“Separating Your Test Bench from Your Function Code” on page 12-25

“Preserving Your Code” on page 12-25

“File Naming Conventions” on page 12-26

Recommended Compilation Options for `fiaccl`

- `-args` - Specify input parameters by example

Use the `-args` option to specify the properties of primary function inputs as a cell array of example values at the same time as you generate code for the MATLAB file with `fiaccl`. The cell array can be a variable or literal array of constant values. The cell array should provide the same number and order of inputs as the primary function.

When you use the `-args` option you are specifying the data types and array dimensions of these parameters, not the values of the variables. For more information, see “Define Input Properties by Example at the Command Line” (MATLAB Coder).

Note Alternatively, you can use the `assert` function to define properties of primary function inputs directly in your MATLAB file. For more information, see “Define Input Properties Programmatically in MATLAB File” on page 12-35.

- `-report` - Generate code generation report

Use the `-report` option to generate a report in HTML format at code generation time to help you debug your MATLAB code and verify that it is suitable for code generation. If you do not specify the `-report` option, `fiaccl` generates a report only if build errors or warnings occur.

The code generation report contains the following information:

- Summary of code generation results, including type of target and number of warnings or errors
- Build log that records build and linking activities
- Links to generated files
- Error and warning messages (if any)

For more information, see `fiaccl`.

Build Scripts

Use build scripts to call `fiaccl` to generate MEX functions from your MATLAB function.

A build script automates a series of MATLAB commands that you want to perform repeatedly from the command line, saving you time and eliminating input errors. For instance, you can use a build script to clear your workspace before each build and to specify code generation options.

This example shows a build script to run `fiaccl` to process `lms_02.m`:

```
close all;
clear all;
clc;

N = 73113;

fiaccl -report lms_02.m ...
      -args { zeros(N,1) zeros(N,1) }
```

In this example, the following actions occur:

- `close all` deletes all figures whose handles are not hidden. See `close` in the MATLAB Graphics function reference for more information.
- `clear all` removes all variables, functions, and MEX-files from memory, leaving the workspace empty. This command also clears all breakpoints.

Note Remove the `clear all` command from the build scripts if you want to preserve breakpoints for debugging.

- `clc` clears all input and output from the Command Window display, giving you a “clean screen.”
- `N = 73113` sets the value of the variable `N`, which represents the number of samples in each of the two input parameters for the function `lms_02`
- `fiaccl -report lms_02.m -args { zeros(N,1) zeros(N,1) }` calls `fiaccl` to accelerate simulation of the file `lms_02.m` using the following options:
 - `-report` generates a code generation report
 - `-args { zeros(N,1) zeros(N,1) }` specifies the properties of the function inputs as a cell array of example values. In this case, the input parameters are `N`-by-1 vectors of real doubles.

Check Code Interactively Using MATLAB Code Analyzer

The code analyzer checks your code for problems and recommends modifications to maximize performance and maintainability. You can use the code analyzer to check your code continuously in the MATLAB Editor while you work.

To ensure that continuous code checking is enabled:

- 1 On the MATLAB **Home** tab, click **Preferences**. Select **Code Analyzer** to view the list of code analyzer preferences.
- 2 Select the **Enable integrated warning and error messages** check box.

Separating Your Test Bench from Your Function Code

Separate your core algorithm from your test bench. Create a separate test script to do all the pre- and post-processing such as loading inputs, setting up input values, calling the function under test, and outputting test results. See the example on the `fiaccl` reference page.

Preserving Your Code

Preserve your code before making further modifications. This practice provides a fallback in case of error and a baseline for testing and validation. Use a consistent file naming convention, as described

in “File Naming Conventions” on page 12-26. For example, add a 2-digit suffix to the file name for each file in a sequence. Alternatively, use a version control system.

File Naming Conventions

Use a consistent file naming convention to identify different types and versions of your MATLAB files. This approach keeps your files organized and minimizes the risk of overwriting existing files or creating two files with the same name in different folders.

For example, the file naming convention in the Generating MEX Functions getting started tutorial is:

- The suffix `_build` identifies a build script.
- The suffix `_test` identifies a test script.
- A numerical suffix, for example, `_01` identifies the version of a file. These numbers are typically two-digit sequential integers, beginning with 01, 02, 03, and so on.

For example:

- The file `build_01.m` is the first version of the build script for this tutorial.
- The file `test_03.m` is the third version of the test script for this tutorial.

Code Generation Reports

In this section...

“Report Generation” on page 12-27
“Report Location” on page 12-27
“Errors and Warnings” on page 12-27
“Files and Functions” on page 12-27
“MATLAB Source” on page 12-28
“MATLAB Variables” on page 12-29
“Code Insights” on page 12-30
“Report Limitations” on page 12-30

When you enable report generation or an error occurs, `fiaccl` generates a code generation report. Use the report to debug your MATLAB functions and verify that they are suitable for code generation. The report provides type information for the variables and expressions in your functions. This information helps you to find sources of error messages and to understand type propagation rules.

Report Generation

To control generation and opening of the report, use `fiaccl` options:

- To generate a report, use the `-report` option.
- To generate and open a report, use the `-launchreport` option.

Alternatively, use configuration object properties:

- To generate a report, set `GenerateReport` to `true`.
- If you want `fiaccl` to open the report for you, set `LaunchReport` to `true`.

Report Location

The code generation report is named `report.mldatx`. It is located in the `html` subfolder of the code generation output folder. If you have MATLAB R2018a or later, you can open the `report.mldatx` file by double-clicking it.

Errors and Warnings

View code generation error, warning, and information messages on the **All Messages** tab. To highlight the source code for an error or warning, click the message. It is a best practice to address the first message because subsequent errors and warnings can be related to the first message.

Files and Functions

In the **MATLAB Source** pane, the **Function List** view organizes functions according to the containing file. To visualize functions according to the call structure, use the **Call Tree** view.

To view a function in the code pane of the report, click the function in the list. Clicking a function opens the file that contains the function. To edit the selected file in the MATLAB Editor, click **Edit in MATLAB** or click a line number in the code pane.

Specialized Functions or Classes

When a function is called with different types of inputs or a class uses different types for its properties, the code generator produces specializations. In the **MATLAB Source** pane, numbered functions (or classes) indicate specializations. For example:

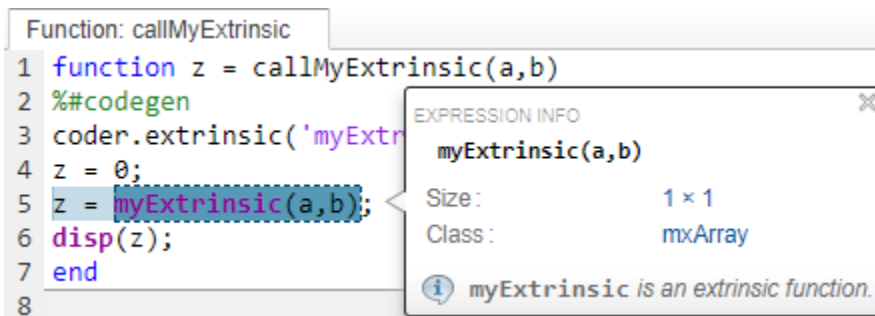
```
fx fcn > 1
fx fcn > 2
```

MATLAB Source

To view a MATLAB function in the code pane, click the name of the function in the **MATLAB Source** pane. In the code pane, when you pause on a variable or expression, a tooltip displays information about its size, type, and complexity. Additionally, syntax highlighting helps you to identify MATLAB syntax elements and certain code generation attributes, such as whether a function is extrinsic or whether an argument is constant.

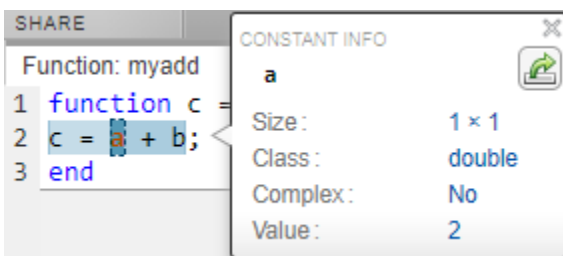
Extrinsic Functions

The report identifies an extrinsic function with purple text. The tooltip indicates that the function is extrinsic.




Constant Arguments

Orange text indicates a compile-time constant argument to an entry-point function or a specialized function. The tooltip includes the constant value.



Knowing the value of a constant argument helps you to understand the generated function signatures. It also helps you to see when code generation creates function specializations for different constant argument values.

To export the value to a variable in the workspace, click the Export icon .

MATLAB Variables

The **Variables** tab provides information about the variables for the selected MATLAB function. To select a function, click the function in the **MATLAB Source** pane.

The variables table shows:


- Class, size, and complexity
- Properties of fixed-point types

This information helps you to understand type propagation and identify type mismatch errors.

Visual Indicators on the Variables Tab

This table describes the symbols, badges, and other indicators in the variables table.

| Column in the Variables Table | Indicator | Description |
|-------------------------------|-----------|---|
| Name | expander | Variable has elements or properties that you can see by clicking the expander. |
| Name | {:} | Homogenous cell array (all elements have the same properties). |
| Name | {n} | nth element of a heterogeneous cell array. |
| Class | v > n | v is reused with a different class, size, and complexity. The number n identifies a reuse with a unique set of properties. When you pause on a renamed variable, the report highlights only the instances of this variable that share the class, size, and complexity. See “Reuse the Same Variable with Different Properties” on page 18-10. |
| Size | :n | Variable-size array with an upper bound of n. |
| Size | :? | Variable-size array with no upper bound. |

| Column in the Variables Table | Indicator | Description |
|-------------------------------|---|--|
| Size | italics | Variable-size array whose dimensions do not change size during execution. |
| Class | sparse prefix | Sparse array. |
| Class | complex prefix | Complex number. |
| Class |  | Fixed-point type. To see the fixed-point properties, click the badge. |

Code Insights

If you enable potential differences reporting, you can view the messages on the **Code Insights** tab. The report includes potential differences messages only if you enabled potential differences reporting. See “Potential Differences Reporting” on page 19-23.

Report Limitations

- The entry-point summary shows the individual elements of `varargin` and `varargout`, but the variables table does not show them.
- The report does not show full information for unrolled loops. It displays data types of one arbitrary iteration.
- The report does not show information about dead code.

See Also

`fiaccel`

More About

- “Accelerate Code Using `fiaccel`” on page 12-6
- “Reuse the Same Variable with Different Properties” on page 18-10

Generate C Code from Code Containing Global Data

In this section...

“Workflow Overview” on page 12-31
 “Declaring Global Variables” on page 12-31
 “Defining Global Data” on page 12-31
 “Synchronizing Global Data with MATLAB” on page 12-32
 “Limitations of Using Global Data” on page 12-34

Workflow Overview

To generate MEX functions from MATLAB code that uses global data:

- 1 Declare the variables as global in your code.
- 2 Define and initialize the global data before using it.

For more information, see “Defining Global Data” on page 12-31.

- 3 Compile your code using `fiaccl`.

If you use global data, you must also specify whether you want to synchronize this data between MATLAB and the generated code. If there is no interaction between MATLAB and the generated code, it is safe to disable synchronization. Otherwise, you should enable synchronization. For more information, see “Synchronizing Global Data with MATLAB” on page 12-32.

Declaring Global Variables

For code generation, you must declare global variables before using them in your MATLAB code. Consider the `use_globals` function that uses two global variables AR and B.

```
function y = use_globals()
%#codegen
% Turn off inlining to make
% generated code easier to read
coder.inline('never');
% Declare AR and B as global variables
global AR;
global B;
AR(1) = B(1);
y = AR * 2;
```

Defining Global Data

You can define global data either in the MATLAB global workspace or at the command line. If you do not initialize global data at the command line, `fiaccl` looks for the variable in the MATLAB global workspace. If the variable does not exist, `fiaccl` generates an error.

Defining Global Data in the MATLAB Global Workspace

To compile the `use_globals` function described in “Declaring Global Variables” on page 12-31 using `fiaccl`:

- 1 Define the global data in the MATLAB workspace. At the MATLAB prompt, enter:

```
global AR B;
AR = fi(ones(4),1,16,14);
B = fi([1 2 3],1,16,13);
```

- 2 Compile the function to generate a MEX file named `use_globalsx`.

```
fiaccel -o use_globalsx use_globals
```

Defining Global Data at the Command Line

To define global data at the command line, use the `fiaccel -global` option. For example, to compile the `use_globals` function described in “Declaring Global Variables” on page 12-31, specify two global inputs `AR` and `B` at the command line.

```
fiaccel -o use_globalsx ...
    -global {'AR',fi(ones(4)),'B',fi([1 2 3])} use_globals
```

Alternatively, specify the type and initial value with the `-globals` flag using the format `-globals {'g', {type, initial_value}}`.

Defining Variable-Sized Global Data

To provide initial values for variable-sized global data, specify the type and initial value with the `-globals` flag using the format `-globals {'g', {type, initial_value}}`. For example, to specify a global variable `g1` that has an initial value `[1 1]` and upper bound `[2 2]`, enter:

```
fiaccel foo -globals {'g1',{coder.typeof(0,[2 2],1),[1 1]}}
```

For a detailed explanation of `coder.typeof` syntax, see `coder.typeof`.

Synchronizing Global Data with MATLAB

Why Synchronize Global Data?

The generated code and MATLAB each have their own copies of global data. To ensure consistency, you must synchronize their global data whenever the two interact. If you do not synchronize the data, their global variables might differ. The level of interaction determines when to synchronize global data.

When to Synchronize Global Data

By default, synchronization between global data in MATLAB and generated code occurs at MEX function entry and exit and for all extrinsic calls, which are calls to MATLAB functions on the MATLAB path that `fiaccel` dispatches to MATLAB for execution. This behavior ensures maximum consistency between generated code and MATLAB.

To improve performance, you can:

- Select to synchronize only at MEX function entry and exit points.
- Disable synchronization when the global data does not interact.
- Choose whether to synchronize before and after each extrinsic call.

The following table summarizes which global data synchronization options to use. To learn how to set these options, see “How to Synchronize Global Data” on page 12-33.

Global Data Synchronization Options

| If you want to... | Set the global data synchronization mode to: | Synchronize before and after extrinsic calls? |
|---|---|--|
| Ensure maximum consistency when all extrinsic calls modify global data. | At MEX-function entry, exit and extrinsic calls (default) | Yes. Default behavior. |
| Ensure maximum consistency when most extrinsic calls modify global data, but a few do not. | At MEX-function entry, exit and extrinsic calls (default) | Yes. Use the <code>coder.extrinsic -sync:off</code> option to turn off synchronization for the extrinsic calls that do not affect global data. |
| Ensure maximum consistency when most extrinsic calls do not modify global data, but a few do. | At MEX-function entry and exit | Yes. Use the <code>coder.extrinsic -sync:on</code> option to synchronize only the calls that modify global data. |
| Maximize performance when synchronizing global data, and none of your extrinsic calls modify global data. | At MEX-function entry and exit | No. |
| Communicate between generated code files only. No interaction between global data in MATLAB and generated code. | Disabled | No. |

How to Synchronize Global Data

To control global data synchronization, set the global data synchronization mode and select whether to synchronize extrinsic functions. For guidelines on which options to use, see “When to Synchronize Global Data” on page 12-32.

You control the synchronization of global data with extrinsic functions using the `coder.extrinsic -sync:on` and `-sync:off` options.

Controlling the Global Data Synchronization Mode from the Command Line

- 1 Define the compiler options object in the MATLAB workspace by issuing a constructor command:

```
comp_cfg = coder.mexconfig
```

- 2 From the command line, set the `GlobalDataSyncMethod` property to `Always`, `SyncAtEntryAndExits` or `NoSync`, as applicable. For example:

```
comp_cfg.GlobalDataSyncMethod = 'SyncAtEntryAndExits';
```

- 3 Use the `comp_cfg` configuration object when compiling your code by specifying it using the `-config` compilation option. For example,

```
fiaccel -config comp_cfg myFile
```

Controlling Synchronization for Extrinsic Function Calls

You can control whether synchronization between global data in MATLAB and generated code occurs before and after you call an extrinsic function. To do so, use the `coder.extrinsic -sync:on` and `-sync:off` options.

By default, global data is:

- Synchronized before and after each extrinsic call if the global data synchronization mode is `At MEX-function entry, exit and extrinsic calls`. If you are sure that certain extrinsic calls do not affect global data, turn off synchronization for these calls using the `-sync:off` option. Turning off synchronization improves performance. For example, if functions `foo1` and `foo2` *do not* affect global data, turn off synchronization for these functions:

```
coder.extrinsic('-sync:off', 'foo1', 'foo2');
```

- Not synchronized if the global data synchronization mode is `At MEX-function entry and exit`. If the code has a few extrinsic calls that affect global data, turn on synchronization for these calls using the `-sync:on` option. For example, if functions `foo1` and `foo2` *do* affect global data, turn on synchronization for these functions:

```
coder.extrinsic('-sync:on', 'foo1', 'foo2');
```

- Not synchronized if the global data synchronization mode is `Disabled`. When synchronization is disabled, you cannot control the synchronization for specific extrinsic calls. The `-sync:on` option has no effect.

Clear Global Data

Because MEX functions and MATLAB each have their own copies of global data, you must `clear` both copies to ensure that consecutive MEX runs produce the same results. The `clear global` command removes only the copy of the global data in the MATLAB workspace. To remove both copies of the data, use the `clear global` and `clear mex` commands together. The `clear all` command also removes both copies.

Limitations of Using Global Data

You cannot use global data with the `coder.varsize` function. Instead, use a `coder.typeof` object to define variable-sized global data as described in “Defining Variable-Sized Global Data” on page 12-32.

See Also

More About

- “Specify Global Cell Arrays at the Command Line” on page 12-47
- “Convert Code Containing Global Data to Fixed Point” on page 7-56

Define Input Properties Programmatically in MATLAB File

In this section...

“How to Use `assert`” on page 12-35
 “Rules for Using `assert` Function” on page 12-38
 “Specifying Properties of Primary Fixed-Point Inputs” on page 12-38
 “Specifying Properties of Cell Arrays” on page 12-39
 “Specifying Class and Size of Scalar Structure” on page 12-40
 “Specifying Class and Size of Structure Array” on page 12-41

How to Use `assert`

You can use the MATLAB `assert` function to define properties of primary function inputs directly in your MATLAB file.

Use the `assert` function to invoke standard MATLAB functions for specifying the class, size, and complexity of primary function inputs.

Specify Any Class

```
assert ( isa ( param, 'class_name' ) )
```

Sets the input parameter *param* to the MATLAB class *class_name*. For example, to set the class of input *U* to a 32-bit signed integer, call:

```
...
assert(isa(U, 'embedded.fi'));
...
```

Note If you set the class of an input parameter to `fi`, you must also set its `numerictype`, see “Specify `numerictype` of Fixed-Point Input” on page 12-37. You can also set its `fimath` properties, see “Specify `fimath` of Fixed-Point Input” on page 12-38.

If you set the class of an input parameter to `struct`, you must specify the properties of each field in the structure in the order in which you define the fields in the structure definition.

Specify `fi` Class

```
assert ( isfi ( param ) )
assert ( isa ( param, 'embedded.fi' ) )
```

Sets the input parameter *param* to the MATLAB class `fi` (fixed-point numeric object). For example, to set the class of input *U* to `fi`, call:

```
...
assert(isfi(U));
...
```

or

```
...  
assert(isa(U, 'embedded.fi'));  
...
```

Note If you set the class of an input parameter to `fi`, you must also set its `numericity`, see “Specify `numericity` of Fixed-Point Input” on page 12-37. You can also set its `fimath` properties, see “Specify `fimath` of Fixed-Point Input” on page 12-38.

Specify Structure Class

```
assert ( isstruct ( param ) )
```

Sets the input parameter *param* to the MATLAB class `struct` (structure). For example, to set the class of input *U* to a `struct`, call:

```
...  
assert(isstruct(U));  
...
```

or

```
...  
assert(isa(U, 'struct'));  
...
```

Note If you set the class of an input parameter to `struct`, you must specify the properties of each field in the structure in the order in which you define the fields in the structure definition.

Specify Cell Array Class

```
assert(iscell( param ))  
assert(isa(param, 'cell'))
```

Sets the input parameter *param* to the MATLAB class `cell` (cell array). For example, to set the class of input *C* to a `cell`, call:

```
...  
assert(iscell(C));  
...
```

or

```
...  
assert(isa(C, 'cell'));  
...
```

To specify the properties of cell array elements, see “Specifying Properties of Cell Arrays” on page 12-39.

Specify Any Size

```
assert ( all ( size ( param ) == [ dims ] ) )
```

Sets the input parameter *param* to the size specified by dimensions *dims*. For example, to set the size of input *U* to a 3-by-2 matrix, call:

```
...
assert(all(size(U)== [3 2]));
...
```

Specify Scalar Size

```
assert ( isscalar (param ) )
assert ( all ( size (param) == [ 1 ] ) )
```

Sets the size of input parameter *param* to scalar. For example, to set the size of input *U* to scalar, call:

```
...
assert(isscalar(U));
...
```

or

```
...
assert(all(size(U)== [1]));
...
```

Specify Real Input

```
assert ( isreal (param ) )
```

Specifies that the input parameter *param* is real. For example, to specify that input *U* is real, call:

```
...
assert(isreal(U));
...
```

Specify Complex Input

```
assert ( ~isreal (param ) )
```

Specifies that the input parameter *param* is complex. For example, to specify that input *U* is complex, call:

```
...
assert(~isreal(U));
...
```

Specify numerictype of Fixed-Point Input

```
assert ( isequal ( numerictype ( fiparam ), T ) )
```

Sets the `numerictype` properties of `fi` input parameter *fiparam* to the `numerictype` object *T*. For example, to specify the `numerictype` property of fixed-point input *U* as a signed `numerictype` object *T* with 32-bit word length and 30-bit fraction length, use the following code:

```
...
% Define the numerictype object.
T = numerictype(1, 32, 30);

% Set the numerictype property of input U to T.
assert(isequal(numerictype(U),T));
...
```

Specify fimath of Fixed-Point Input

```
assert ( isequal ( fimath ( fiparam ), F ) )
```

Sets the `fimath` properties of `fi` input parameter `fiparam` to the `fimath` object `F`. For example, to specify the `fimath` property of fixed-point input `U` so that it saturates on integer overflow, use the following code:

```
...  
% Define the fimath object.  
F = fimath('OverflowAction','Saturate');  
  
% Set the fimath property of input U to F.  
assert(isequal(fimath(U),F));  
...
```

Specify Multiple Properties of Input

```
assert ( function1 ( params ) && function2 ( params ) && function3 ( params ) && ... )
```

Specifies the class, size, and complexity of one or more inputs using a single `assert` function call. For example, the following code specifies that input `U` is a double, complex, 3-by-3 matrix, and input `V` is a 16-bit unsigned integer:

```
...  
assert(isa(U,'double') && ~isreal(U) && all(size(U) == [3 3]) && isa(V,'uint16'));  
...
```

Rules for Using assert Function

Follow these rules when using the `assert` function to specify the properties of primary function inputs:

- Call `assert` functions at the beginning of the primary function, before any flow-control operations such as `if` statements or subroutine calls.
- Do not call `assert` functions inside conditional constructs, such as `if`, `for`, `while`, and `switch` statements.
- If you set the class of an input parameter to `fi`:
 - You must also set its `numericType`, see “Specify `numericType` of Fixed-Point Input” on page 12-37.
 - You can also set its `fimath` properties, see “Specify `fimath` of Fixed-Point Input” on page 12-38.
- If you set the class of an input parameter to `struct`, you must specify the class, size, and complexity of each field in the structure in the order in which you define the fields in the structure definition.

Specifying Properties of Primary Fixed-Point Inputs

In the following example, the primary MATLAB function `emcsqrtfi` takes one fixed-point input: `x`. The code specifies the following properties for this input:

| Property | Value |
|------------|---|
| class | fi |
| numericity | numericity object T, as specified in the primary function |
| fimath | fimath object F, as specified in the primary function |
| size | scalar (by default) |
| complexity | real (by default) |

```
function y = emcsqrtfi(x)
T = numericity('WordLength',32,'FractionLength',23,...
    'Signed',true);
F = fimath('SumMode','SpecifyPrecision',...
    'SumWordLength',32,'SumFractionLength',23,...
    'ProductMode','SpecifyPrecision',...
    'ProductWordLength',32,'ProductFractionLength',23);
assert(isfi(x));
assert(isequal(numericity(x),T));
assert(isequal(fimath(x),F));

y = sqrt(x);
```

Specifying Properties of Cell Arrays

To specify the MATLAB class `cell` (cell array), use one of the following syntaxes:

```
assert(iscell(param))
assert(isa(param, 'cell'))
```

For example, to set the class of input `C` to `cell`, use:

```
...
assert(iscell(C));
...
```

or

```
...
assert(isa(C, 'cell'));
...
```

You can also specify the size of the cell array and the properties of the cell array elements. The number of elements that you specify determines whether the cell array is homogeneous or heterogeneous. See “Code Generation for Cell Arrays” (MATLAB Coder).

If you specify the properties of the first element only, the cell array is homogeneous. For example, the following code specifies that `C` is a 1x3 homogeneous cell array whose elements are 1x1 double.

```
...
assert(isa(C, 'cell'));
assert(all(size(C) == [1 3]));
assert(isa(C{1}, 'double'));
...
```

If you specify the properties of each element, the cell array is heterogeneous. For example, the following code specifies a 1x2 heterogeneous cell array whose first element is 1x1 char and whose second element is 1x3 double.

```
...
assert(isa(C, 'cell'));
assert(all(size(C) == [1 2]));
assert(isa(C{1}, 'char'));
assert(all(size(C{2}) == [1 3]));
assert(isa(C{2}, 'double'));
...
```

If you specify more than one element, you cannot specify that the cell array is variable size, even if all elements have the same properties. For example, the following code specifies a variable-size cell array. Because the code specifies the properties of the first and second elements, code generation fails.

```
...
assert(isa(C, 'cell'));
assert(all(size(C) <= [1 2]));
assert(isa(C{1}, 'double'));
assert(isa(C{2}, 'double'));
...
```

In the previous example, if you specify the first element only, you can specify that the cell array is variable-size. For example:

```
...
assert(isa(C, 'cell'));
assert(all(size(C) <= [1 2]));
assert(isa(C{1}, 'double'));
...
```

Specifying Class and Size of Scalar Structure

Assume you have defined `S` as the following scalar MATLAB structure:

```
S = struct('r',double(1),'i',fi(4,true,8,0));
```

This code specifies the class and size of `S` and its fields when passed as an input to your MATLAB function:

```
function y = fcn(S)

% Specify the class of the input as struct.
assert(isstruct(S));

% Specify the size of the fields r and i
% in the order in which you defined them.
T = numerictype('Wordlength', 8,'FractionLength', ...
    0,'signed',true);
assert(isa(S.r,'double'));
assert(isfi(S.i) && isequal(numerictype(S.i),T));

y = S;
```

Note The only way to name a field in a structure is to set at least one of its properties. Therefore in the preceding example, an `assert` function specifies that field `S.r` is of type `double`, even though `double` is the default.

Specifying Class and Size of Structure Array

For structure arrays, you must choose a representative element of the array for specifying the properties of each field. For example, assume you have defined `S` as the following 1-by-2 array of MATLAB structures:

```
S = struct('r',{double(1), double(2)},'i',...
          {fi(4,1,8,0), fi(5,1,8,0)});
```

The following code specifies the class and size of each field of structure input `S` using the first element of the array:

```
function y = fcn(S)

% Specify the class of the input S as struct.
assert(isstruct(S));
T = numerictype('Wordlength', 8,'FractionLength', ...
              0,'signed',true);

% Specify the size of the fields r and i
% based on the first element of the array.
assert(all(size(S) == [1 2]));
assert(isa(S(1).r,'double'));
assert(isfi(S(1).i) && isequal(numerictype(S(1).i),T));

y = S;
```

Note The only way to name a field in a structure is to set at least one of its properties. Therefore in the example above, an `assert` function specifies that field `S(1).r` is of type `double`, even though `double` is the default.

Specify Cell Array Inputs at the Command Line

To specify cell array inputs at the command line, use the same methods that you use for other types of inputs. You can:

- Provide an example cell array input to the `-args` option of the `fiaccel` command.
- Provide a `coder.CellType` object to the `-args` option of the `fiaccel` command. To create a `coder.CellType` object, use `coder.typeof`.
- Use `coder.Constant` to specify a constant cell array input.

For code generation, cell arrays are classified as homogeneous or heterogeneous. See “Code Generation for Cell Arrays” on page 30-2. When you provide an example cell array to `fiaccel` or `coder.typeof`, the function determines whether the cell array type is homogeneous or heterogeneous. If the cell array elements have the same class and size, `coder.typeof` returns a homogeneous cell array type. If the elements have different classes, `coder.typeof` returns a heterogeneous cell array type. For some cell arrays, the classification as homogeneous or heterogeneous is ambiguous. For example, the type for `{1 [2 3]}` can be a 1x2 heterogeneous type. The first element is double and the second element is 1x2 double. The type can also be a 1x3 homogeneous type in which the elements have class double and size 1x:2. For these ambiguous cases, `coder.typeof` uses heuristics to classify the type as homogeneous or heterogeneous. If you want a different classification, use the `coder.CellType` `makeHomogeneous` or `makeHeterogeneous` methods. The `makeHomogeneous` method makes a homogeneous copy of a type. The `makeHeterogeneous` method makes a heterogeneous copy of a type.

The `makeHomogeneous` and `makeHeterogeneous` methods permanently assign the classification as homogeneous and heterogeneous, respectively. You cannot later use one of these methods to create a copy that has a different classification.

Specify Cell Array Inputs by Example

To specify a cell array input by example, provide an example cell array in the `-args` option of the `fiaccel` command.

For example:

- To specify a 1x3 cell array whose elements have class double:

```
fiaccel myfunction -args {[1 2 3]} -report
```

The input argument is a 1x3 homogeneous cell array whose elements are 1x1 double.

- To specify a 1x2 cell array whose first element has class char and whose second element has class double:

```
fiaccel myfunction -args {'a', 1} -report
```

The input argument is a 1x2 heterogeneous cell array whose first element is 1x1 char and whose second element is 1x1 double.

Specify the Type of the Cell Array Input

To specify the type of a cell array input, use `coder.typeof` to create a `coder.CellType` object. Pass the `coder.CellType` object to the `-args` option of the `fiaccel` command.

For example:

- To specify a 1x3 cell array whose elements have class double:

```
t = coder.typeof({1 2 3});
fiaccl myfunction -args {t} -report
```

The input argument is a 1x3 homogeneous cell array whose elements are 1x1 double.

- To specify a 1x2 cell array whose first element has class char and whose second element has class double:

```
t = coder.typeof({'a', 1});
fiaccl myfunction -args {t}
```

The input argument is a 1x2 heterogeneous cell array whose first element is a 1x1 char and whose second element is a 1x1 double.

You can also use the advanced function `coder.newtype` to create a `coder.CellType` object.

Make a Homogeneous Copy of a Type

If `coder.typeof` returns a heterogeneous cell array type, but you want a homogeneous type, use the `makeHomogeneous` method to make a homogeneous copy of the type.

The following code creates a heterogeneous type.

```
t = coder.typeof({1 [2 3]})
t =
coder.CellType
  1x2 heterogeneous cell
  f0: 1x1 double
  f1: 1x2 double
```

To make a homogeneous copy of the type, use:

```
t = makeHomogeneous(t)
t =
coder.CellType
  1x2 locked homogeneous cell
  base: 1x:2 double
```

Alternatively, use this notation:

```
t = makeHomogeneous(coder.typeof({1 [2 3]}))
t =
coder.CellType
  1x2 locked homogeneous cell
  base: 1x:2 double
```

The classification as homogeneous is locked (permanent). You cannot later use the `makeHeterogeneous` method to make a heterogeneous copy of the type.

If the elements of a type have different classes, such as char and double, you cannot use `makeHomogeneous` to make a homogeneous copy of the type.

Make a Heterogeneous Copy of a Type

If `coder.typeof` returns a homogeneous cell array type, but you want a heterogeneous type, use the `makeHeterogeneous` method to make a heterogeneous copy of the type.

The following code creates a homogeneous type.

```
t = coder.typeof({1 2 3})
t =
coder.CellType
  1x3 homogeneous cell
  base: 1x1 double
```

To make the type heterogeneous, use:

```
t = makeHeterogeneous(t)
t =
coder.CellType
  1x3 locked heterogeneous cell
  f1: 1x1 double
  f2: 1x1 double
  f3: 1x1 double
```

Alternatively, use this notation:

```
t = makeHeterogeneous(coder.typeof({1 2 3}))
t =
coder.CellType
  1x3 locked heterogeneous cell
  f1: 1x1 double
  f2: 1x1 double
  f3: 1x1 double
```

The classification as heterogeneous is locked (permanent). You cannot later use the `makeHomogeneous` method to make a homogeneous copy of the type.

If a type is variable size, you cannot use `makeHeterogeneous` to make a heterogeneous copy of it.

Specify Variable-Size Cell Array Inputs

You can specify variable-size cell array inputs in the following ways:

- In the `coder.typeof` call.

For example, to specify a variable-size cell array whose first dimension is fixed and whose second dimension has an upper bound of 5:

```
t = coder.typeof({1}, [1 5], [0 1])
```

```
t =
coder.CellType
  1x5 homogeneous cell
  base: 1x1 double
```

For elements with the same classes, but different sizes, you can use `coder.typeof` size and variable dimensions arguments to create a variable-size homogeneous cell array type. For example, the following code does not use the size and variable dimensions arguments. This code creates a type for a heterogeneous cell array.

```
t = coder.typeof({1 [2 3]})
t =
coder.CellType
  1x2 heterogeneous cell
  f0: 1x1 double
  f1: 1x2 double
```

The following code, that uses the size and dimensions arguments, creates a type for a variable-size homogeneous type cell array:

```
t = coder.typeof({1 [2 3]}, [1 5], [0 1])
t =
coder.CellType
  1x5 locked homogeneous cell
  base: 1x2 double
```

- Use `coder.resize`.

For example, to specify a variable-size cell array whose first dimension is fixed and whose second dimension has an upper bound of 5:

```
t = coder.typeof({1});
t = coder.resize(t, [1 5], [0,1])
t =
coder.CellType
  1x5 homogeneous cell
  base: 1x1 double
```

You cannot use `coder.resize` with a heterogeneous cell array type.

Specify Constant Cell Array Inputs

To specify that a cell array input is constant, use the `coder.Constant` function with the `-args` option of the `fiaccel` command. For example:

```
fiaccel myfunction -args {coder.Constant({'red',1,'green',2,'blue',3})} -report
```

The input is a 1x6 heterogeneous cell array. The sizes and classes of the elements are:

- 1x3 char

- 1x1 double
- 1x5 char
- 1x1 double
- 1x4 char
- 1x1 double

See Also

`coder.CellType` | `coder.typeof` | `coder.resize` | `coder.newtype`

Related Examples

- “Define Input Properties by Example at the Command Line” on page 31-4
- “Specify Constant Inputs at the Command Line” on page 31-6

More About

- “Code Generation for Cell Arrays” on page 30-2

Specify Global Cell Arrays at the Command Line

To specify global cell array inputs, use the `-globals` option of the `fiaccl` command with this syntax:

```
fiaccl myfunction -globals {global_var, {type_object, initial_value}}
```

For example:

- To specify that the global variable `g` is a 1x3 cell array whose elements have class `double` and whose initial value is `{1 2 3}`, use:

```
fiaccl myfunction -globals {'g', {coder.typeof({1 1 1}), {1 2 3}}}
```

Alternatively, use:

```
t = coder.typeof({1 1 1});
fiaccl myfunction -globals {'g', {t, {1 2 3}}}
```

The global variable `g` is a 1x3 homogeneous cell array whose elements are 1x1 `double`.

To make `g` heterogeneous, use:

```
t = makeHeterogeneous(coder.typeof({1 1 1}));
fiaccl myfunction -globals {'g', {t, {1 2 3}}}
```

- To specify that `g` is a cell array whose first element has type `char`, whose second element has type `double`, and whose initial value is `{'a', 1}`, use:

```
fiaccl myfunction -globals {'g', {coder.typeof({'a', 1}), {'a', 1}}}
```

The global variable `g` is a 1x2 heterogeneous cell array whose first element is 1x1 `char` and whose second element is 1x1 `double`.

- To specify that `g` is a cell array whose first element has type `double`, whose second element is a 1x2 `double` array, and whose initial value is `{1 [2 3]}`, use:

```
fiaccl myfunction -globals {'g', {coder.typeof({1 [2 3]}), {1 [2 3]}}}
```

Alternatively, use:

```
t = coder.typeof({1 [2 3]});
fiaccl myfunction -globals {'g', {t, {1 [2 3]}}}
```

The global variable `g` is a 1x2 heterogeneous cell array whose first element is 1x1 `double` and whose second element is 1x2 `double`.

Global variables that are cell arrays cannot have variable size.

See Also

`fiaccl` | `coder.typeof`

Related Examples

- “Generate C Code from Code Containing Global Data” on page 12-31

Control Run-Time Checks

In this section...

“Types of Run-Time Checks” on page 12-48

“When to Disable Run-Time Checks” on page 12-48

“How to Disable Run-Time Checks” on page 12-49

Types of Run-Time Checks

In simulation, the code generated for your MATLAB functions includes the following run-time checks and external function calls.

- Memory integrity checks

These checks detect violations of memory integrity in code generated for MATLAB functions and stop execution with a diagnostic message.

Caution For safety, these checks are enabled by default. Without memory integrity checks, violations will result in unpredictable behavior.

- Responsiveness checks in code generated for MATLAB functions

These checks enable periodic checks for Ctrl+C breaks in code generated for MATLAB functions. Enabling responsiveness checks also enables graphics refreshing.

Caution For safety, these checks are enabled by default. Without these checks the only way to end a long-running execution might be to terminate MATLAB.

- Extrinsic calls to MATLAB functions

Extrinsic calls to MATLAB functions, for example to display results, are enabled by default for debugging purposes. For more information about extrinsic functions, see “Use the coder.extrinsic Construct” on page 14-7.

When to Disable Run-Time Checks

Generally, generating code with run-time checks enabled results in more generated code and slower simulation than generating code with the checks disabled. Similarly, extrinsic calls are time consuming and have an adverse effect on performance. Disabling run-time checks and extrinsic calls usually results in streamlined generated code and faster simulation, with these caveats:

| Consider disabling... | Only if... |
|-------------------------|--|
| Memory integrity checks | You are sure that your code is safe and that all array bounds and dimension checking is unnecessary. |
| Responsiveness checks | You are sure that you will not need to stop execution of your application using Ctrl+C . |

| Consider disabling... | Only if... |
|-----------------------|---|
| Extrinsic calls | You are only using extrinsic calls to functions that do not affect application results. |

How to Disable Run-Time Checks

To disable run-time checks:

- 1 Define the compiler options object in the MATLAB workspace by issuing a constructor command:

```
comp_cfg = coder.MEXConfig
```

- 2 From the command line set the `IntegrityChecks`, `ExtrinsicCalls`, or `ResponsivenessChecks` properties false, as applicable:

```
comp_cfg.IntegrityChecks = false;  
comp_cfg.ExtrinsicCalls = false;  
comp_cfg.ResponsivenessChecks = false;
```

Fix Run-Time Stack Overflows

If your C compiler reports a run-time stack overflow, set the value of the maximum stack usage parameter to be less than the available stack size. Create a command-line configuration object, `coder.mexconfig` and then set the `StackUsageMax` parameter.

Code Generation with MATLAB Coder

MATLAB Coder `codegen` automatically converts MATLAB code directly to C code. It generates standalone C code that is bit-true to fixed-point MATLAB code. Using Fixed-Point Designer and MATLAB Coder software you can generate C code with algorithms containing integer math only (i.e., without any floating-point math).

Fixed-Point FIR Code Example Parameter Values

| Block | Parameter | Value |
|--------------|------------------------------------|------------------|
| Constant | Constant value | b |
| | Interpret vector parameters as 1-D | Not selected |
| | Sampling mode | Sample based |
| | Sample time | inf |
| | Mode | Fixed point |
| | Signedness | Signed |
| | Scaling | Slope and bias |
| | Word length | 12 |
| | Slope | 2 ⁻¹² |
| | Bias | 0 |
| Constant1 | Constant value | x+noise |
| | Interpret vector parameters as 1-D | Unselected |
| | Sampling mode | Sample based |
| | Sample time | 1 |
| | Mode | Fixed point |
| | Signedness | Signed |
| | Scaling | Slope and bias |
| | Word length | 12 |
| | Slope | 2 ⁻⁸ |
| | Bias | 0 |
| Constant2 | Constant value | zi |
| | Interpret vector parameters as 1-D | Unselected |
| | Sampling mode | Sample based |
| | Sample time | inf |
| | Mode | Fixed point |
| | Signedness | Signed |
| | Scaling | Slope and bias |
| | Word length | 12 |
| | Slope | 2 ⁻⁸ |
| | Bias | 0 |
| To Workspace | Variable name | yout |
| | Limit data points to last | inf |
| | Decimation | 1 |

| Block | Parameter | Value |
|----------------------|--|--------------|
| | Sample time | -1 |
| | Save format | Array |
| | Log fixed-point data as a fi object | Selected |
| To Workspace1 | Variable name | zf |
| | Limit data points to last | inf |
| | Decimation | 1 |
| | Sample time | -1 |
| | Save format | Array |
| | Log fixed-point data as a fi object | Selected |
| To Workspace2 | Variable name | noisyx |
| | Limit data points to last | inf |
| | Decimation | 1 |
| | Sample time | -1 |
| | Save format | Array |
| | Log fixed-point data as a fi object | Selected |

Accelerate Code for Variable-Size Data

In this section...

“Disable Support for Variable-Size Data” on page 12-54

“Control Dynamic Memory Allocation” on page 12-54

“Accelerate Code for MATLAB Functions with Variable-Size Data” on page 12-55

“Accelerate Code for a MATLAB Function That Expands a Vector in a Loop” on page 12-56

Variable-size data is data whose size might change at run time. MATLAB supports bounded and unbounded variable-size data for code generation. Bounded variable-size data has fixed upper bounds. This data can be allocated statically on the stack or dynamically on the heap. Unbounded variable-size data does not have fixed upper bounds. This data must be allocated on the heap. By default, for MEX and C/C++ code generation, support for variable-size data is enabled and dynamic memory allocation is enabled for variable-size arrays whose size exceeds a configurable threshold.

Disable Support for Variable-Size Data

By default, for MEX and C/C++ code acceleration, support for variable-size data is enabled. You modify variable sizing settings at the command line.

- 1 Create a configuration object for code generation.

```
cfg = coder.mexconfig;
```

- 2 Set the `EnableVariableSizing` option:

```
cfg.EnableVariableSizing = false;
```

- 3 Using the `-config` option, pass the configuration object to `fiaccel` :

```
fiaccel -config cfg foo
```

Control Dynamic Memory Allocation

By default, dynamic memory allocation is enabled for variable-size arrays whose size exceeds a configurable threshold. If you disable support for variable-size data, you also disable dynamic memory allocation. You can modify dynamic memory allocation settings at the command line.

- 1 Create a configuration object for code acceleration. For example, for a MEX function:

```
mexcfg = coder.mexconfig;
```

- 2 Set the `EnableDynamicMemoryAllocation` option:

| Setting | Action |
|--|---|
| <code>mexcfg.EnableDynamicMemoryAllocation=false;</code> | Dynamic memory allocation is disabled. All variable-size data is allocated statically on the stack. |

| Setting | Action |
|---|--|
| <code>mexcfg.EnableDynamicMemoryAllocation=true;</code> | Dynamic memory allocation is enabled for all variable-size arrays whose size (in bytes) is greater than or equal to the value specified using the Dynamic memory allocation threshold parameter. Variable-size arrays whose size is less than this threshold are allocated on the stack. |

- Optionally, if you set `Enable dynamic memory allocation` to `true`, configure `Dynamic memory allocation threshold` to fine tune memory allocation.
- Using the `-config` option, pass the configuration object to `fiaccel`:

```
fiaccel -config mexcfg foo
```

Accelerate Code for MATLAB Functions with Variable-Size Data

Here is a basic workflow that generates MEX code.

- In the MATLAB Editor, add the compilation directive `%#codegen` at the top of your function.

This directive:

- Indicates that you intend to generate code for the MATLAB algorithm
- Turns on checking in the MATLAB Code Analyzer to detect potential errors during code generation

- Address issues detected by the Code Analyzer.

In some cases, the MATLAB Code Analyzer warns you when your code assigns data a fixed size but later grows the data, such as by assignment or concatenation in a loop. If that data is supposed to vary in size at run time, you can ignore these warnings.

- Generate a MEX function using `fiaccel`. Use the following command-line options:

- `-args {coder.typeof...}` if you have variable-size inputs
- `-report` to generate a code generation report

For example:

```
fiaccel -report foo -args {coder.typeof(0,[2 4],1)}
```

This command uses `coder.typeof` to specify one variable-size input for function `foo`. The first argument, `0`, indicates the input data type (`double`) and complexity (`real`). The second argument, `[2 4]`, indicates the size, a matrix with two dimensions. The third argument, `1`, indicates that the input is variable sized. The upper bound is 2 for the first dimension and 4 for the second dimension.

Note During compilation, `fiaccel` detects variables and structure fields that change size after you define them, and reports these occurrences as errors. In addition, `fiaccel` performs a runtime check to generate errors when data exceeds upper bounds.

4 Fix size mismatch errors:

| Cause: | How To Fix: | For More Information: |
|--|--|---|
| You try to change the size of data after its size has been locked. | Declare the data to be variable sized. | See “Diagnosing and Fixing Size Mismatch Errors” on page 29-15. |

5 Fix upper bounds errors

| Cause: | How To Fix: | For More Information: |
|---|---|---|
| MATLAB cannot determine or compute the upper bound | Specify an upper bound. | See “Specify Upper Bounds for Variable-Size Arrays” and “Diagnosing and Fixing Size Mismatch Errors” on page 29-15. |
| MATLAB attempts to compute an upper bound for unbounded variable-size data. | If the data is unbounded, enable dynamic memory allocation. | See “Control Dynamic Memory Allocation” on page 12-54 |

6 Generate C/C++ code using the `fiaccl` function.

Accelerate Code for a MATLAB Function That Expands a Vector in a Loop

- “About the MATLAB Function `uniquetol`” on page 12-56
- “Step 1: Add Compilation Directive for Code Generation” on page 12-56
- “Step 2: Address Issues Detected by the Code Analyzer” on page 12-57
- “Step 3: Generate MEX Code” on page 12-57
- “Step 4: Fix the Size Mismatch Error” on page 12-58
- “Step 5: Compare Execution Speed of MEX Function to Original Code” on page 12-59

About the MATLAB Function `uniquetol`

This example uses the function `uniquetol`. This function returns in vector `B` a version of input vector `A`, where the elements are unique to within tolerance `tol` of each other. In vector `B`, $\text{abs}(B(i) - B(j)) > \text{tol}$ for all `i` and `j`. Initially, assume input vector `A` can store up to 100 elements.

```
function B = uniquetol(A, tol)
A = sort(A);
B = A(1);
k = 1;
for i = 2:length(A)
    if abs(A(k) - A(i)) > tol
        B = [B A(i)];
        k = i;
    end
end
```

Step 1: Add Compilation Directive for Code Generation

Add the `%#codegen` compilation directive at the top of the function:


```
function B = uniquetol(A, tol) %#codegen
A = sort(A);
B = A(1);
k = 1;
for i = 2:length(A)
    if abs(A(k) - A(i)) > tol
        B = [B A(i)];
        k = i;
    end
end
```

Step 2: Address Issues Detected by the Code Analyzer

The Code Analyzer detects that variable `B` might change size in the `for`-loop. It issues this warning:

The variable 'B' appears to change size on every loop iteration.
Consider preallocating for speed.

In this function, vector `B` should expand in size as it adds values from vector `A`. Therefore, you can ignore this warning.

Step 3: Generate MEX Code

To generate MEX code, use the `fiaccel` function.

- 1 Generate a MEX function for `uniquetol`:

```
T = numerictype(1, 16, 15);
fiaccel -report uniquetol -args {coder.typeof(fi(0,T),[1 100],1),coder.typeof(fi(0,T))}
```

What do these command-line options mean?

`T = numerictype(1, 16, 15)` creates a signed `numerictype` object with a 16-bit word length and 15-bit fraction length that you use to specify the data type of the input arguments for the function `uniquetol`.

The `fiaccel` function `-args` option specifies the class, complexity, and size of each input to function `uniquetol`:

- The first argument, `coder.typeof`, defines a variable-size input. The expression `coder.typeof(fi(0,T),[1 100],1)` defines input `A` as a vector of real, signed embedded.fi objects that have a 16-bit word length and 15-bit fraction length. The vector has a fixed upper bound; its first dimension is fixed at 1 and its second dimension can vary in size up to 100 elements.

For more information, see “Specify Variable-Size Inputs at the Command Line” (MATLAB Coder).

- The second argument, `coder.typeof(fi(0,T))`, defines input `tol` as a real, signed embedded.fi object with a 16-bit word length and 15-bit fraction length.

The `-report` option instructs `fiaccel` to generate a code generation report, even if no errors or warnings occur.

For more information, see the `fiaccel` reference page.

Executing this command generates a compiler error:

??? Size mismatch (size [1 x 1] ~= size [1 x 2]).
The size to the left is the size
of the left-hand side of the assignment.

- 2 Open the error report and select the **Variables** tab.

Function: `uniquetol`

```

1 function B = uniquetol(A, tol) %#codegen
2 A = sort(A);
3 coder.varsize('B');
4 B = A(1);
5 k = 1;
6 for i = 2:length(A)
7     if abs(A(k) - A(i)) > tol
8         B = [B A(i)];
9         k = i;
10    end
11 end

```

| Order | Variable | Type | Size | Class | Complex | Signedness | WL | FL |
|-------|----------|--------|---------|-------------|---------|------------|----|----|
| 1 | B | Output | 1 x 1 | embedded.fi | No | Signed | 16 | 15 |
| 2 | A > 1 | Input | 1 x 100 | embedded.fi | No | Signed | 16 | 15 |
| 3 | A > 2 | Local | 1 x ? | embedded.fi | No | Signed | 16 | 15 |
| 4 | tol | Input | 1 x 1 | embedded.fi | No | Signed | 16 | 15 |
| 5 | k | Local | 1 x 1 | double | No | - | - | - |
| 6 | i | Local | 1 x 1 | double | No | - | - | - |

The error indicates a size mismatch between the left-hand side and right-hand side of the assignment statement `B = [B A(i)]`; . The assignment `B = A(1)` establishes the size of B as a fixed-size scalar (1 x 1). Therefore, the concatenation of `[B A(i)]` creates a 1 x 2 vector.

Step 4: Fix the Size Mismatch Error

To fix this error, declare B to be a variable-size vector.

- 1 Add this statement to the `uniquetol` function:

```
coder.varsize('B');
```

It should appear before B is used (read). For example:

```
function B = uniquetol(A, tol) %#codegen
A = sort(A);
```

```
coder.varsize('B');
```

```
B = A(1);
k = 1;
for i = 2:length(A)
```

```

    if abs(A(k) - A(i)) > tol
        B = [B A(i)];
        k = i;
    end
end

```

The function `coder. varsize` declares every instance of `B` in `uniquetol` to be variable sized.

- 2 Generate code again using the same command:

```
fiaccl -report uniquetol -args {coder.typeof(fi(0,T),[1 100],1),coder.typeof(fi(0,T))}
```

In the current folder, `fiaccl` generates a MEX function for `uniquetol` named `uniquetol_mex` and provides a link to the code generation report.

- 3 Click the *View report* link.
- 4 In the code generation report, select the **Variables** tab.

Function: `uniquetol`

```

1 function B = uniquetol(A, tol) %#codegen
2 A = sort(A);
3 coder. varsize('B');
4 B = A(1);
5 k = 1;
6 for i = 2:length(A)
7     if abs(A(k) - A(i)) > tol
8         B = [B A(i)];
9         k = i;
10    end
11 end

```

| Summary | All Messages (0) | Variables | | | | | | |
|---------|------------------|-----------|---------|-------------|---------|------------|----|----|
| Order | Variable | Type | Size | Class | Complex | Signedness | WL | FL |
| 1 | B | Output | 1 x ? | embedded.fi | No | Signed | 16 | 15 |
| 2 | A | Input | 1 x 100 | embedded.fi | No | Signed | 16 | 15 |
| 3 | tol | Input | 1 x 1 | embedded.fi | No | Signed | 16 | 15 |
| 4 | k | Local | 1 x 1 | double | No | - | - | - |
| 5 | i | Local | 1 x 1 | double | No | - | - | - |

The size of variable `B` is `1 x ?`, indicating that it is variable size with no upper bounds.

Step 5: Compare Execution Speed of MEX Function to Original Code

Run the original MATLAB algorithm and MEX function with the same inputs for the same number of loop iterations and compare their execution speeds.

- 1 Create inputs of the correct class, complexity, and size to pass to the `uniquetol` MATLAB and MEX functions.

```
x = fi(rand(1,90), T);  
tol = fi(0, T);
```

- 2 Run the original `uniquetol` function in a loop and time how long it takes to execute 10 iterations of the loop.

```
tic; for k=1:10, b = uniquetol(x,tol); end; tSim=toc
```

- 3 Run the generated MEX function with the same inputs for the same number of loop iterations.

```
tic; for k=1:10, b = uniquetol_mex(x,tol); end; tSim_mex=toc
```

- 4 Compare the execution times.

```
r = tSim/tSim_mex
```

This example shows that generating a MEX function using `fiaccel` greatly accelerates the execution of the fixed-point algorithm.

Code Generation Readiness Tool

The code generation readiness tool screens MATLAB code for features and functions that code generation does not support. The tool provides a report that lists the source files that contain unsupported features and functions. It is possible that the tool does not detect all code generation issues. Under certain circumstances, it is possible that the tool can report false errors. Therefore, before you generate code, verify that your code is suitable for code generation by generating a MEX function.

The code generation readiness tool does not report functions that the code generator automatically treats as extrinsic. Examples of such functions are `plot`, `disp`, and `figure`.

Issues Tab

2 Code generation readiness issues - Code might require changes Language C/C++ (MATLAB Coder)
 2 Unsupported functions Refresh Edit
 2 Files analyzed

Issues Files Group by: Issue

- Unsupported function: `hascycles` (1)
- Unsupported function: `isdag` (1)

Unsupported function: `hascycles`

```
foo2.m
1 function [tf1,tf2] = foo2(source,target)
2 G = digraph(source,target);
3 tf1 = hascycles(G);
4 tf2 = isdag(G);
5 end
6
```

On the **Issues** tab, the tool displays information about:

- MATLAB syntax issues. These issues are reported in the MATLAB editor. To learn more about the issues and how to fix them, use the Code Analyzer.
- Unsupported MATLAB function calls, language features, and data types.

You can also:

- View your MATLAB code inside the Code Generation Readiness Tool. When you select an issue, the part of your MATLAB code that caused this issue gets highlighted.
- Group the readiness results either by issue or by file.
- Select the language that the code generation readiness analysis uses.
- Refresh the code generation readiness analysis if you updated your MATLAB code.
- Export the analysis report either as plain text file or as a `coder.ScreenerInfo` object in the base workspace.

Files Tab

The screenshot shows the Code Generation Readiness Tool interface. At the top left, there is a warning icon and the text "2 Code generation readiness issues - Code might require changes". Below this, it says "2 Unsupported functions" and "2 Files analyzed". On the right, there is a dropdown menu for "Language C/C++ (MATLAB Coder)" with "Refresh" and "Edit" buttons. Below the header, there are tabs for "Issues" and "Files". The "Files" tab is active, showing a tree view with "foo1" expanded to show "foo2". Below the tree, there is a message "Unsupported function: hascycles". At the bottom, the MATLAB code for "foo2.m" is displayed, with "hascycles(G)" highlighted in red and a red exclamation mark icon in the margin.

If the code that you are checking calls functions in other MATLAB code files, the **Files** tab shows the call dependency between these files. If you select **Show MathWorks Functions**, the report also lists the MathWorks functions that your function calls.

See Also

`coder.screener` | `coder.ScreenerInfo` Properties

Related Examples

- “MATLAB Language Features Supported for C/C++ Code Generation” (MATLAB Coder)

- “Functions and Objects Supported for C/C++ Code Generation” (MATLAB Coder)

Check Code Using the Code Generation Readiness Tool

Run Code Generation Readiness Tool at the Command Line

- 1 Navigate to the folder that contains the file that you want to check for code generation readiness.
- 2 At the MATLAB command prompt, enter:

```
coder.screener('filename')
```

The **Code Generation Readiness** tool opens for the file named `filename`, provides a code generation readiness score, and lists issues that must be fixed prior to code generation.

Run the Code Generation Readiness Tool From the Current Folder Browser

- 1 In the current folder browser, right-click the file that you want to check for code generation readiness.
- 2 From the context menu, select **Check Code Generation Readiness**.

The **Code Generation Readiness** tool opens for the selected file and provides a code generation readiness score and lists issues that must be fixed prior to code generation.

See Also

- “Code Generation Readiness Tool” on page 12-61

Check Code Using the MATLAB Code Analyzer

The code analyzer checks your code for problems and recommends modifications. You can use the code analyzer to check your code interactively in the MATLAB Editor while you work.

To verify that continuous code checking is enabled:

- 1** In MATLAB, select the **Home** tab and then click **Preferences**.
- 2** In the **Preferences** dialog box, select **Code Analyzer**.
- 3** In the **Code Analyzer Preferences** pane, verify that **Enable integrated warning and error messages** is selected.

Fix Errors Detected at Code Generation Time

When the code generator detects errors or warnings, it automatically generates an error report. The error report describes the issues and provides links to the MATLAB code with errors.

To fix the errors, modify your MATLAB code to use only those MATLAB features that are supported for code generation. For more information, see “Algorithm Design Basics”. Choose a debugging strategy for detecting and correcting code generation errors in your MATLAB code. For more information, see “Debugging Strategies” on page 12-14.

When code generation is complete, the software generates a MEX function that you can use to test your implementation in MATLAB.

If your MATLAB code calls functions on the MATLAB path, unless the code generator determines that these functions should be extrinsic or you declare them to be extrinsic, it attempts to compile these functions. See “Resolution of Function Calls for Code Generation” on page 14-2. To get detailed diagnostics, add the `%#codegen` directive to each external function that you want codegen to compile.

See Also

- “Code Generation Reports” on page 12-27
- “Why Test MEX Functions in MATLAB?” (MATLAB Coder)
- “When to Generate Code from MATLAB Algorithms” on page 19-2
- “Debugging Strategies” on page 12-14
- “Use the `coder.extrinsic` Construct” on page 14-7

Avoid Multiword Operations in Generated Code

This example shows how to avoid multiword operations in generated code by using the `accumpos` function instead of simple addition in your MATLAB algorithm. Similarly, you can use `accumneg` for subtraction.

This example requires a MATLAB Coder license.

Write a simple MATLAB algorithm that adds two numbers and returns the result.

```
function y = my_add1(a,b)
    y = a+b;
end
```

Write a second MATLAB algorithm that adds two numbers using `accumpos` and returns the result.

```
function y = my_add2(a,b)
    y = accumpos(a,b); %floor, wrap
end
```

`accumpos` adds `a` and `b` using the data type of `a`. `b` is cast into the data type of `a`. If `a` is a `fi` object, by default, `accumpos` sets the rounding mode to 'Floor' and the overflow action to 'Wrap'. It ignores the `fimath` properties of `a` and `b`.

Compare the outputs of the two functions in MATLAB.

```
a = fi(1.25,1,32,5);
b = fi(0.125,0,32);
```

```
y1 = my_add1(a,b)
y2 = my_add2(a,b)
```

```
y1 =
```

```
1.3750
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 63
    FractionLength: 34
```

```
y2 =
```

```
1.3750
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 32
    FractionLength: 5
```

For the simple addition, the word length grows but using `accumpos`, the word length of the result is the same as that of `a`.

Generate C code for the function `my_add1`. First, disable use of the `long long` data type because it is not usually supported by the target hardware.

```
hw = coder.HardwareImplementation;
hw.ProdHWDeviceType = 'Generic->32-bit Embedded Processor';
```

```
hw.ProdLongLongMode = false;
hw.ProdBitPerLong = 32;
cfg = coder.config('lib');
cfg.HardwareImplementation = hw;
codegen my_add1 -args {a,b} -report -config cfg
```

MATLAB Coder generates a C static library and provides a link to the code generation report.

View the generated code for the simple addition. Click the [View report](#) link to open the code generation report and then scroll to the code for the `my_add1` function.

```
/* Function Declarations */
static void MultiWordAdd(const unsigned long u1[], const unsigned long u2[],
    unsigned long y[], int n);
static void MultiWordSignedWrap(const unsigned long u1[], int n1, unsigned int
    n2, unsigned long y[]);
static void sLong2MultiWord(long u, unsigned long y[], int n);
static void sMultiWord2MultiWord(const unsigned long u1[], int n1, unsigned long
    y[], int n);
static void sMultiWord2sMultiWordSat(const unsigned long u1[], int n1, unsigned
    long y[], int n);
static void sMultiWordShl(const unsigned long u1[], int n1, unsigned int n2,
    unsigned long y[], int n);
static void sMultiWordShr(const unsigned long u1[], int n1, unsigned int n2,
    unsigned long y[], int n);
static void uLong2MultiWord(unsigned long u, unsigned long y[], int n);
```

The generated C code contains multiple multiword operations.

Generate C code for the function `my_add2`.

```
codegen my_add2 -args {a,b} -report -config cfg
```

View the generated code for the addition using `accumpos`. Click the [View report](#) link to open the code generation report and then scroll to the code for the `my_add2` function.

```
int my_add2(int a, unsigned int b)
{
    int y;
    y = a + (int)(b >> 29);
    /* floor, wrap */
    return y;
}
```

For this function, the generated code contains no multiword operations.

Find Potential Data Type Issues in Generated Code

In this section...

- “Data Type Issues Overview” on page 12-69
- “Enable Highlighting of Potential Data Type Issues” on page 12-69
- “Find and Address Cumbersome Operations” on page 12-69
- “Find and Address Expensive Rounding” on page 12-70
- “Find and Address Expensive Comparison Operations” on page 12-71
- “Find and Address Multiword Operations” on page 12-71


Data Type Issues Overview

When you convert MATLAB code to fixed point, you can highlight potential data type issues in the generated report. The report highlights MATLAB code that requires single-precision, double-precision, or expensive fixed-point operations.

- The double-precision check highlights expressions that result in a double-precision operation. When trying to achieve a strict-single or fixed-point design, manual inspection of code can be time-consuming and error prone.
- The single-precision check highlights expressions that result in a single operation.
- The expensive fixed-point operations check identifies optimization opportunities for fixed-point code. It highlights expressions in the MATLAB code that require cumbersome multiplication or division, expensive rounding, expensive comparison, or multiword operations. For more information on optimizing generated fixed-point code, see “Tips for Making Generated Code More Efficient” on page 49-9.

Enable Highlighting of Potential Data Type Issues

Enable the highlight option using the Fixed-Point Converter app

- 1 On the **Convert to Fixed Point** page, click the **Settings** arrow .
- 2 Under **Plotting and Reporting**, set **Highlight potential data type issues** to Yes.

When conversion is complete, open the fixed-point conversion report to view the highlighting. Click **View report** in the **Type Validation Output** tab.

Enable the highlight option using the command-line interface

- 1 Create a fixed-point code configuration object:


```
fixptcfg = coder.config('fixpt');
```
- 2 Set the `HighlightPotentialDataTypeIssues` property of the configuration object to `true`.


```
fixptcfg.HighlightPotentialDataTypeIssues = true;
```

Find and Address Cumbersome Operations

Cumbersome operations usually occur due to an insufficient range of output. Avoid inputs to a multiply or divide operation that have word lengths larger than the base integer type of your

processor. Software can process operations with larger word lengths, but this approach requires more code and runs slower.

This example requires Embedded Coder and Fixed-Point Designer. The target word length for the processor in this example is 64.

- 1 Create the function `myMul`.

```
function out = myMul(in1, in2)
    out = fi(in1*in2, 1, 64, 0);
end
```

- 2 Generate code for `myMul`.

```
cfg = coder.config('lib');
cfg.GenerateReport = true;
cfg.HighlightPotentialDataTypeIssues = true;
fm = fimath('ProductMode', 'SpecifyPrecision', 'ProductWordLength', 64);
codegen -config cfg myMul -args {fi(1, 1, 64, 4, fm), fi(1, 1, 64, 4, fm)}
```

- 3 Click **View report**.

- 4 In the code generation report, click the **Code Insights** tab.

- 5 Expand the **Potential data type issues** section. Then, expand the **Expensive fixed-point operations** section.



The report flags the expression `in1 * in2`. To resolve the issue, modify the data types of `in1` and `in2` so that the word length of the product does not exceed the target word length of 64.

Find and Address Expensive Rounding

Traditional handwritten code, especially for control applications, almost always uses "no effort" rounding. For example, for unsigned integers and two's complement signed integers, shifting right and dropping the bits is equivalent to rounding to floor. To get results comparable to, or better than, what you expect from traditional handwritten code, use the `floor` rounding method.

This example requires Embedded Coder and Fixed-Point Designer.

- 1 Create the function `myRounding`.

```
function [quot] = myRounding(in1, in2)
    quot = in1 / in2;
end
```

- 2 Generate code for `myRounding`.

```
cfg = coder.config('lib');
cfg.GenerateReport = true;
cfg.HighlightPotentialDataTypeIssues = true;
codegen -config cfg myRounding -args {fi(1, 1, 16, 2), fi(1, 1, 16, 4)}
```

- 3 Click **View report**.

- 4 In the code generation report, click the **Code Insights** tab.

- 5 Expand the **Potential data type issues** section. Then, expand the **Expensive fixed-point operations** section.



The division operation `in1/in2` uses the default rounding method, nearest. Changing the rounding method to `Floor` provides a more efficient implementation.

Find and Address Expensive Comparison Operations

Comparison operations generate extra code when a casting operation is required to do the comparison. For example, before comparing an unsigned integer to a signed integer, one of the inputs must be cast to the signedness of the other. Consider optimizing the data types of the input arguments so that a cast is not required in the generated code.

This example requires Embedded Coder and Fixed-Point Designer.

- 1 Create the function `myRelop`.

```
function out = myRelop(in1, in2)
    out = in1 > in2;
end
```

- 2 Generate code for `myRelop`.

```
cfg = coder.config('lib');
cfg.GenerateReport = true;
cfg.HighlightPotentialDataTypesIssues = true;
codegen -config cfg myRelop -args {fi(1, 1, 14, 3, 1), fi(1, 0, 14, 3, 1)}
```

- 3 Click **View report**.
- 4 In the code generation report, click the **Code Insights** tab.
- 5 Expand the **Potential data type issues** section. Then, expand the **Expensive fixed-point operations** section.



The first input argument, `in1`, is signed, while `in2` is unsigned. Extra code is generated because a cast must occur before the two inputs can be compared.

Change the signedness and scaling of one of the inputs to generate more efficient code.

Find and Address Multiword Operations

Multiword operations can be inefficient on hardware. When an operation has an input or output data type larger than the largest word size of your processor, the generated code contains multiword operations. You can avoid multiword operations in the generated code by specifying local `fimath` properties for variables. You can also manually specify input and output word lengths of operations that generate multiword code.

This example requires Embedded Coder and Fixed-Point Designer. In this example, the target word length is 64.

- 1 Create the function `myMul`.

```
function out = myMul(in1, in2)
    out = in1 * in2;
end
```

- 2 Generate code for `myMul`.

```
cfg = coder.config('lib');
cfg.GenerateReport = true;
cfg.HighlightPotentialDataTypeIssues = true;
codegen -config cfg myMul -args {fi(1, 1, 33, 4), fi(1, 1, 32, 4)}
```

- 3 Click **View report**.
- 4 In the code generation report, click the **Code Insights** tab.
- 5 Expand the **Potential data type issues** section. Then, expand the **Expensive fixed-point operations** section.



- 6 The report flags the `in1 * in2` operation in line 2 of `myMul`. In the code pane, pause over `in1`, `in2`, and the expression `in1 * in2`. You see that:
 - The word length of `in1` is 33 bits and the word length of `in2` is 32 bits.
 - The word length of the expression `in1 * in2` is 65 bits.

The software detects a multiword operation because the word length 65 is larger than the target word length of 64.

- 7 To resolve this issue, modify the data types of `in1` and `in2` so that the word length of the product does not exceed the target word length. Alternatively, specify the `ProductMode` property of the local `fimath` object.

See Also

More About

- “Highlight Potential Data Type Issues in a Report” (Embedded Coder)
- “Code Generation Reports” (MATLAB Coder)

Interoperability with Other Products

- “fi Objects with Simulink” on page 13-2
- “fi Objects with DSP System Toolbox” on page 13-7
- “Ways to Generate Code” on page 13-10

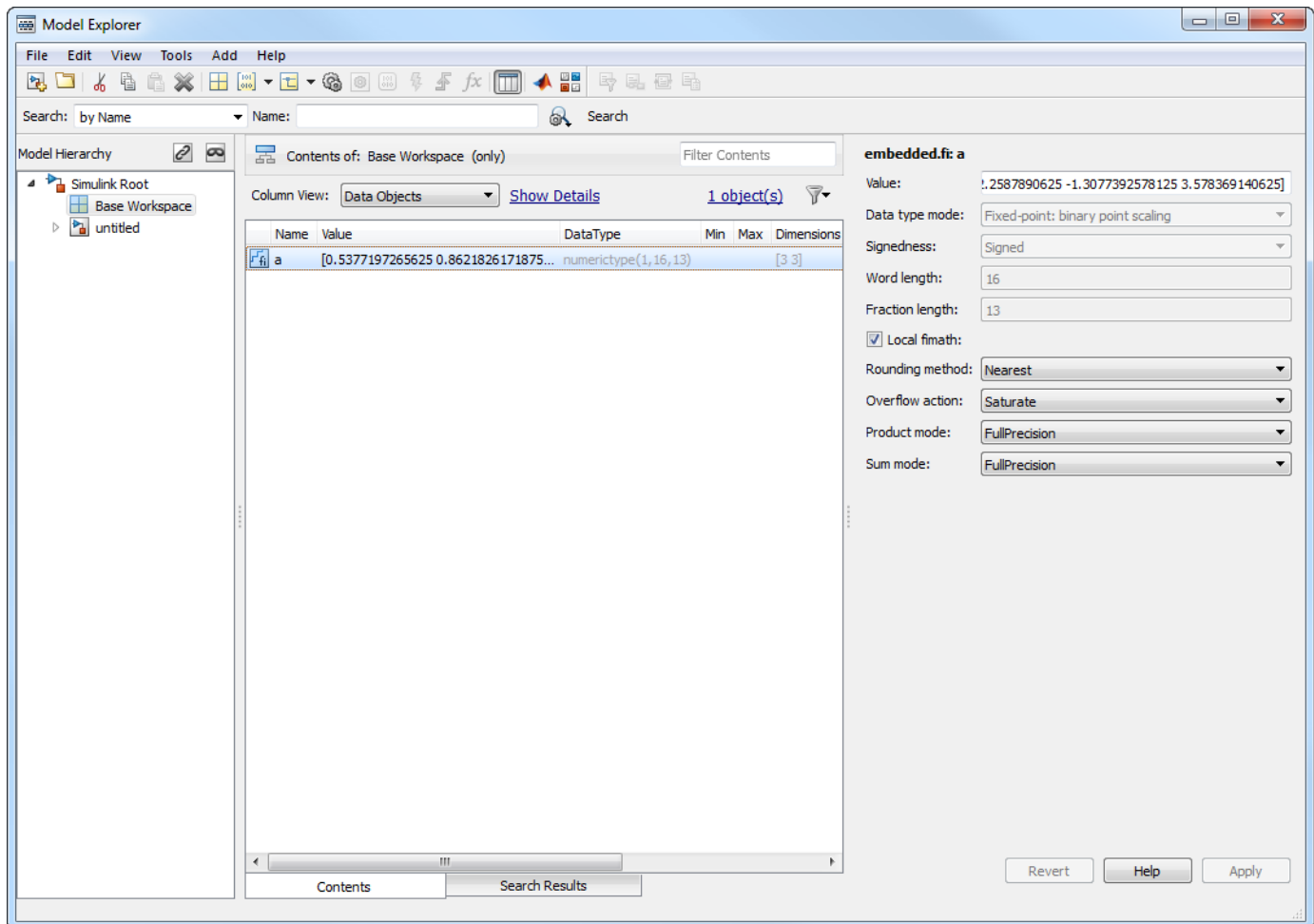
fi Objects with Simulink

In this section...

- “View and Edit fi objects in Model Explorer” on page 13-2
- “Reading Fixed-Point Data from the Workspace” on page 13-3
- “Writing Fixed-Point Data to the Workspace” on page 13-3
- “Setting the Value and Data Type of Block Parameters” on page 13-6
- “Logging Fixed-Point Signals” on page 13-6
- “Accessing Fixed-Point Block Data During Simulation” on page 13-6

View and Edit fi objects in Model Explorer

You can view and edit `fi` objects and their local `fimath` properties using Model Explorer in Simulink. You can change the writable properties of `fi` objects from the Model Explorer, but you cannot change the numeric type properties of `fi` objects after creation.



Reading Fixed-Point Data from the Workspace

You can read fixed-point data from the MATLAB workspace into a Simulink model via the From Workspace block. To do so, the data must be in a structure format with a `fi` object in the `values` field. In array format, the From Workspace block only accepts real, double-precision data.

To read in `fi` data, the **Interpolate data** parameter of the From Workspace block must not be selected, and the **Form output after final data value by** parameter must be set to anything other than Extrapolation.

Writing Fixed-Point Data to the Workspace

You can write fixed-point output from a model to the MATLAB workspace via the To Workspace block in either array or structure format. Fixed-point data written by a To Workspace block to the workspace in structure format can be read back into a Simulink model in structure format by a From Workspace block.

Note To write fixed-point data to the MATLAB workspace as a `fi` object, select the **Log fixed-point data as a fi object** check box on the To Workspace block dialog. Otherwise, fixed-point data is converted to `double` and written to the workspace as `double`.

For example, you can use the following code to create a structure in the MATLAB workspace with a `fi` object in the `values` field. You can then use the From Workspace block to bring the data into a Simulink model.

```
a = fi([sin(0:10)' sin(10:-1:0)'])
```

```
a =
```

```

      0    -0.5440
  0.8415    0.4121
  0.9093    0.9893
  0.1411    0.6570
 -0.7568   -0.2794
 -0.9589   -0.9589
 -0.2794   -0.7568
  0.6570    0.1411
  0.9893    0.9093
  0.4121    0.8415
 -0.5440         0
```

```

      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 16
      FractionLength: 15
```

```
s.signals.values = a
```

```
s =
```

```
signals: [1x1 struct]
```

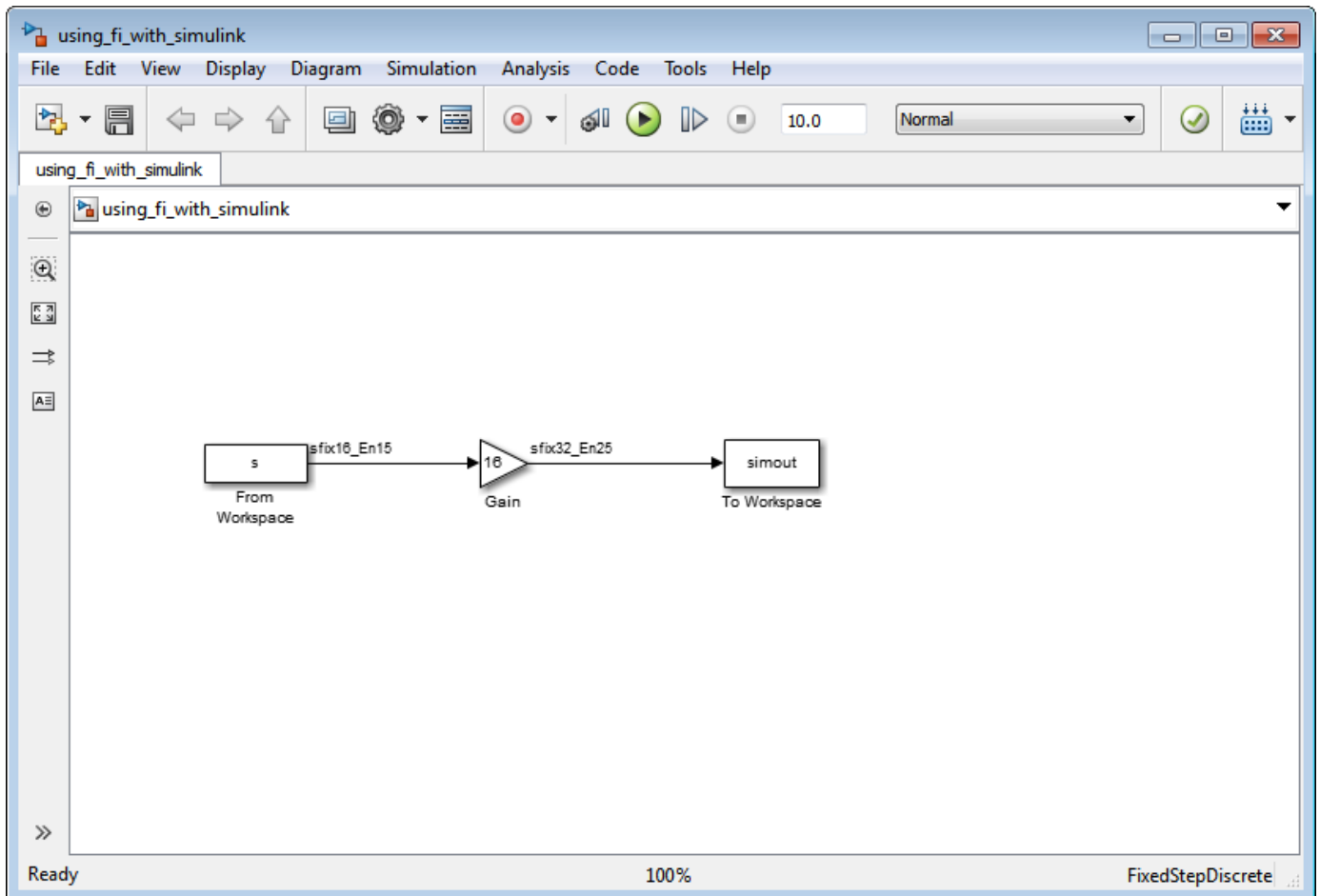
```
s.signals.dimensions = 2  
  
s =  
    signals: [1x1 struct]  
  
s.time = [0:10]'  
  
s =  
    signals: [1x1 struct]  
    time: [11x1 double]
```

The From Workspace block in the following model has the `fi` structure `s` in the **Data** parameter.

Remember, to write fixed-point data to the MATLAB workspace as a `fi` object, select the **Log fixed-point data as a fi object** check box on the To Workspace block dialog. Otherwise, fixed-point data is converted to double and written to the workspace as double.

In the model, the following parameters in the **Solver** pane of the **Model Configuration Parameters** dialog have the indicated settings:

- **Start time** — 0.0
- **Stop time** — 10.0
- **Type** — Fixed-step
- **Solver** — Discrete (no continuous states)
- **Fixed step size (fundamental sample time)** — 1.0



The To Workspace block writes the result of the simulation to the MATLAB workspace as a `fi` structure.

```
simout.signals.values
```

```
ans =
```

```

      0   -8.7041
 13.4634   6.5938
 14.5488  15.8296
   2.2578  10.5117
 -12.1089  -4.4707
 -15.3428 -15.3428
  -4.4707 -12.1089
 10.5117   2.2578
 15.8296  14.5488
   6.5938  13.4634
 -8.7041    0

```

```

DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 32

```

FractionLength: 25

Setting the Value and Data Type of Block Parameters

You can use Fixed-Point Designer expressions to specify the value and data type of block parameters in Simulink. For more information, see “Specify Fixed-Point Data Types”.

Logging Fixed-Point Signals

When fixed-point signals are logged to the MATLAB workspace via signal logging, they are always logged as `fi` objects.

To enable signal logging for a signal:

- 1 Select the signal.
- 2 Open the **Record** dropdown.
- 3 Select **Log/Unlog Selected Signals**.

For more information, refer to “Save Signal Data Using Signal Logging”.

When you log signals from a referenced model or Stateflow® chart in your model, the word lengths of `fi` objects may be larger than you expect. The word lengths of fixed-point signals in referenced models and Stateflow charts are logged as the next largest data storage container size.

Accessing Fixed-Point Block Data During Simulation

Simulink provides an application program interface (API) that enables programmatic access to block data, such as block inputs and outputs, parameters, states, and work vectors, while a simulation is running. You can use this interface to develop MATLAB programs capable of accessing block data while a simulation is running or to access the data from the MATLAB command line. Fixed-point signal information is returned to you via this API as `fi` objects. For more information on the API, refer to “Accessing Block Data During Simulation” in the Simulink documentation.

fi Objects with DSP System Toolbox

In this section...

“Reading Fixed-Point Signals from the Workspace” on page 13-7

“Writing Fixed-Point Signals to the Workspace” on page 13-7

Reading Fixed-Point Signals from the Workspace

You can read fixed-point data from the MATLAB workspace into a Simulink model using the Signal From Workspace and Triggered Signal From Workspace blocks from DSP System Toolbox software. Enter the name of the defined `fi` variable in the **Signal** parameter of the Signal From Workspace or Triggered Signal From Workspace block.

Writing Fixed-Point Signals to the Workspace

Fixed-point output from a model can be written to the MATLAB workspace via the To Workspace or Triggered To Workspace block from the blockset. The fixed-point data is always written as a 2-D or 3-D array.

Note To write fixed-point data to the MATLAB workspace as a `fi` object, select the **Log fixed-point data as a fi object** check box on the Signal To Workspace or Triggered To Workspace block dialog. Otherwise, fixed-point data is converted to `double` and written to the workspace as `double`.

For example, you can use the following code to create a `fi` object in the MATLAB workspace. You can then use the Signal From Workspace block to bring the data into a Simulink model.

```
a = fi([sin(0:10)' sin(10:-1:0)'])
```

```
a =
```

```

      0   -0.5440
  0.8415   0.4121
  0.9093   0.9893
  0.1411   0.6570
 -0.7568  -0.2794
 -0.9589  -0.9589
 -0.2794  -0.7568
  0.6570   0.1411
  0.9893   0.9093
  0.4121   0.8415
 -0.5440      0

```

```

      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 16
      FractionLength: 15

```

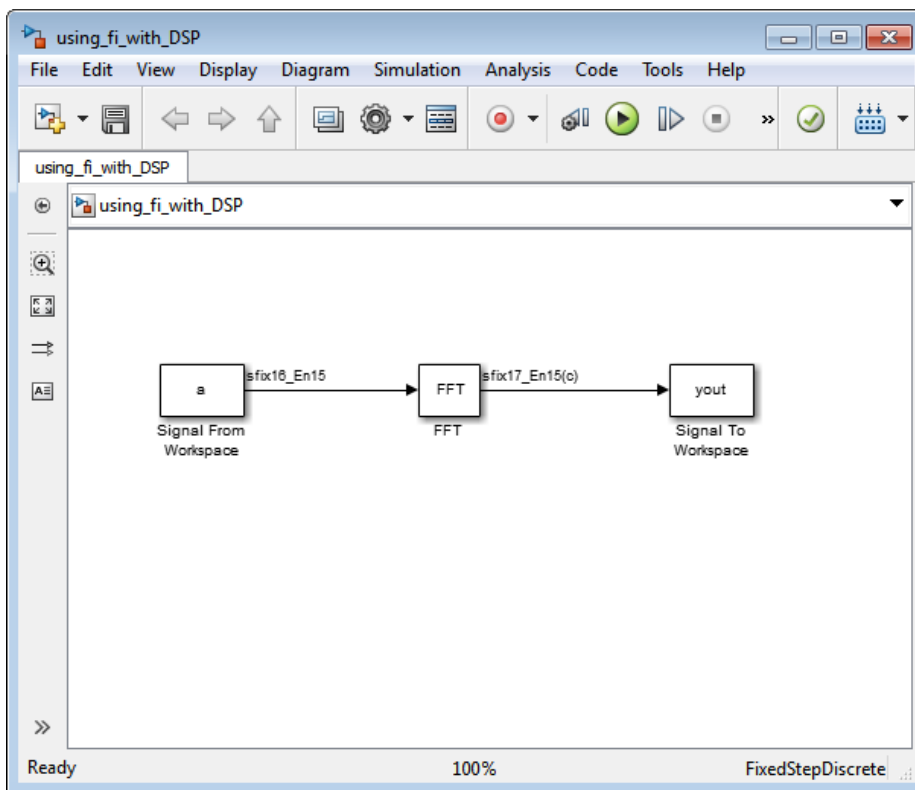
The Signal From Workspace block in the following model has these settings:

- **Signal** — a
- **Sample time** — 1
- **Samples per frame** — 2
- **Form output after final data value by** — Setting to zero

The following parameters in the **Solver** pane of the **Model Configuration Parameters** dialog have these settings:

- **Start time** — 0.0
- **Stop time** — 10.0
- **Type** — Fixed-step
- **Solver** — Discrete (no continuous states)
- **Fixed step size (fundamental sample time)** — 1.0

Remember, to write fixed-point data to the MATLAB workspace as a **fi** object, select the **Log fixed-point data as a fi object** check box on the Signal To Workspace block dialog. Otherwise, fixed-point data is converted to **double** and written to the workspace as **double**.



The Signal To Workspace block writes the result of the simulation to the MATLAB workspace as a **fi** object.

yout =

(:,:,1) =


```
0.8415 -0.1319
-0.8415 -0.9561
```

```
(:,:,2) =
```

```
1.0504 1.6463
0.7682 0.3324
```

```
(:,:,3) =
```

```
-1.7157 -1.2383
0.2021 0.6795
```

```
(:,:,4) =
```

```
0.3776 -0.6157
-0.9364 -0.8979
```

```
(:,:,5) =
```

```
1.4015 1.7508
0.5772 0.0678
```

```
(:,:,6) =
```

```
-0.5440 0
-0.5440 0
```

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 17
FractionLength: 15
```

Ways to Generate Code

There are several ways to use Fixed-Point Designer software to generate code:

- The Fixed-Point Designer `fiaccel` function converts your fixed-point MATLAB code to a MEX function and can greatly accelerate the execution speed of your fixed-point algorithms.
- The MATLAB Coder `codegen` function automatically converts MATLAB code to C/C++ code. Using the MATLAB Coder software allows you to accelerate your MATLAB code that uses Fixed-Point Designer software. To use the `codegen` function with Fixed-Point Designer software, you also need to have a MATLAB Coder license. For more information, see “Generate C Code at the Command Line” (MATLAB Coder).
- The MATLAB Function block allows you to use MATLAB code in your Simulink models that generate embeddable C/C++ code. To use the MATLAB Function block with Fixed-Point Designer software, you also need a Simulink license. For more information on the MATLAB Function block, see the Simulink documentation.

Calling Functions for Code Generation

- “Resolution of Function Calls for Code Generation” on page 14-2
- “Resolution of File Types on Code Generation Path” on page 14-4
- “Compilation Directive %#codegen” on page 14-5
- “Use MATLAB Engine to Execute a Function Call in Generated Code” on page 14-6
- “Code Generation for Recursive Functions” on page 14-12
- “Force Code Generator to Use Run-Time Recursion” on page 14-14
- “Avoid Duplicate Functions in Generated Code” on page 14-17

Resolution of Function Calls for Code Generation

From a MATLAB function, you can call local functions, supported toolbox functions, and other MATLAB functions. MATLAB resolves function names for code generation as follows:

Key Points About Resolving Function Calls

The diagram illustrates key points about how MATLAB resolves function calls for code generation:

- Searches two paths, the code generation path and the MATLAB path
See “Compile Path Search Order” on page 14-2.
- Attempts to compile functions unless the code generator determines that it should not compile them or you explicitly declare them to be extrinsic.

If a MATLAB function is not supported for code generation, you can declare it to be extrinsic by using the construct `coder.extrinsic`, as described in “Use the `coder.extrinsic` Construct” on page 14-7. During simulation, the code generator produces code for the call to an extrinsic function, but does not generate the internal code for the function. Therefore, simulation can run only on platforms where MATLAB software is installed. During standalone code generation, the code generator attempts to determine whether the extrinsic function affects the output of the function in which it is called — for example by returning `mxArrays` to an output variable. If the output does not change, code generation proceeds, but the extrinsic function is excluded from the generated code. Otherwise, compilation errors occur.

The code generator detects calls to many common visualization functions, such as `plot`, `disp`, and `figure`. The software treats these functions like extrinsic functions but you do not have to declare them extrinsic using the `coder.extrinsic` function.

- Resolves file type based on precedence rules described in “Resolution of File Types on Code Generation Path” on page 14-4

Compile Path Search Order

During code generation, function calls are resolved on two paths:

1 Code generation path

MATLAB searches this path first during code generation. The code generation path contains the toolbox functions supported for code generation.

2 MATLAB path

If the function is not on the code generation path, MATLAB searches this path.

MATLAB applies the same dispatcher rules when searching each path (see “Function Precedence Order”).

When to Use the Code Generation Path

Use the code generation path to override a MATLAB function with a customized version. A file on the code generation path shadows a file of the same name on the MATLAB path.

For more information on how to add additional folders to the code generation path, see “Paths and File Infrastructure Setup” (MATLAB Coder).

Resolution of File Types on Code Generation Path

MATLAB uses the following precedence rules for code generation:

Compilation Directive %#codegen

Add the %#codegen directive (or pragma) to your function after the function signature to indicate that you intend to generate code for the MATLAB algorithm. Adding this directive instructs the MATLAB Code Analyzer to help you diagnose and fix violations that would result in errors during code generation.

```
function y = my_fcn(x) %#codegen
```

```
.....
```

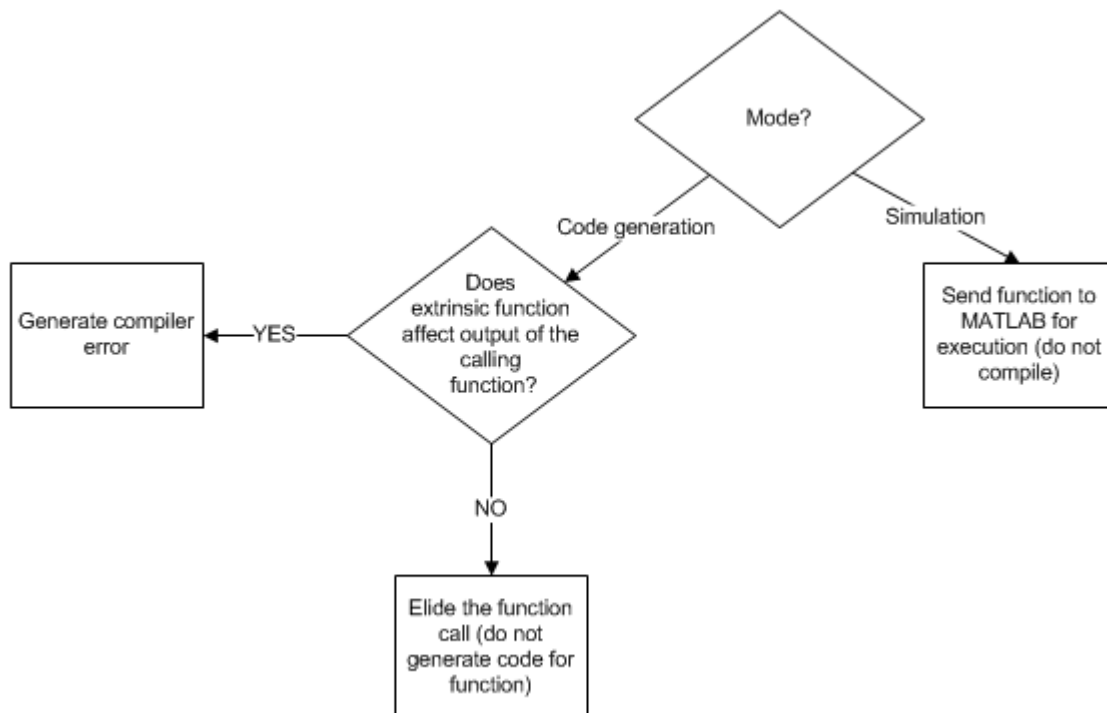
Note The %#codegen directive is not necessary for MATLAB Function blocks. Code inside a MATLAB Function block is always intended for code generation. The %#codegen directive, or the absence of it, does not change the error checking behavior.

Use MATLAB Engine to Execute a Function Call in Generated Code

When processing a call to a function `foo` in your MATLAB code, the code generator finds the definition of `foo` and generates code for its body. In some cases, you might want to bypass code generation and instead use the MATLAB engine to execute the call. Use `coder.extrinsic('foo')` to declare that calls to `foo` do not generate code and instead use the MATLAB engine for execution. In this context, `foo` is referred to as an extrinsic function. This functionality is available only when the MATLAB engine is available during execution. Examples of such situations include execution of MEX functions, Simulink simulations, or function calls at the time of code generation (also known as compile time).

If you generate standalone code for a function that calls `foo` and includes `coder.extrinsic('foo')`, the code generator attempts to determine whether `foo` affects the output. If `foo` does not affect the output, the code generator proceeds with code generation, but excludes `foo` from the generated code. Otherwise, the code generator produces a compilation error.

Including the `coder.extrinsic('foo')` directive inside a certain MATLAB function declares all calls to `foo` inside that MATLAB function as extrinsic. Alternatively, you might want to narrow the scope of extrinsic declaration to just one call to `foo`. See “Call MATLAB Functions Using `feval`” (MATLAB Coder).



When To Declare a Function as Extrinsic

These are some common situations in which you might consider declaring a MATLAB function as extrinsic:

- The function performs display or logging actions. Such functions are useful primarily during simulation and are not used in embedded systems.
- In your MEX execution or Simulink simulation, you want to use a MATLAB function that is not supported for code generation. This workflow does not apply to non-simulation targets.
- You instruct the code generator to constant fold a function call by using `coder.const`. In such situations, the function is called only during code generation when the MATLAB engine is available for executing the call.

Use the `coder.extrinsic` Construct

To declare a function `foo` as extrinsic, include this statement in your MATLAB code.

```
coder.extrinsic('foo')
```

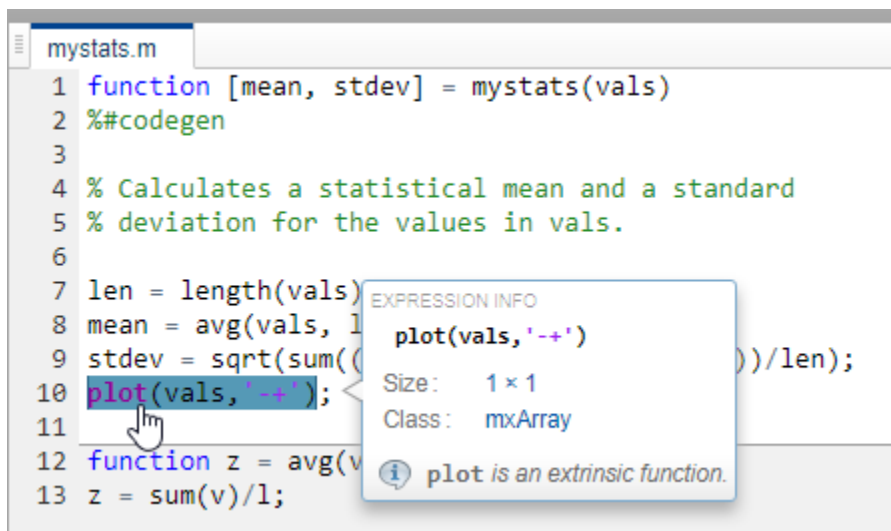
When declaring functions as extrinsic for code generation, adhere to these rules:

- Declare the function as extrinsic before you call it.
- Do not use the extrinsic declaration in conditional statements.
- Assign the return value of an extrinsic function to a known type. See “Working with mxArrays” (MATLAB Coder).

For additional information and examples, see `coder.extrinsic`.

The code generator automatically treats many common MATLAB visualization functions, such as `plot`, `disp`, and `figure`, as extrinsic. You do not have to explicitly declare them as extrinsic functions by using `coder.extrinsic`. For example, you might want to call `plot` to visualize your results in the MATLAB environment. If you generate a MEX function from a function that calls `plot`, and then run the generated MEX function, the code generator dispatches calls to the `plot` function to the MATLAB engine. If you generate a library or executable, the generated code does not contain calls to the `plot` function.

If you generate MEX or standalone C/C++ code by using MATLAB Coder, the code generation report highlights calls from your MATLAB code to extrinsic functions. By inspecting the report, you can determine which functions are supported only in the MATLAB environment.



```
mystats.m
1 function [mean, stdev] = mystats(vals)
2 %#codegen
3
4 % Calculates a statistical mean and a standard
5 % deviation for the values in vals.
6
7 len = length(vals)
8 mean = avg(vals, 1)
9 stdev = sqrt(sum((vals - mean).^2)/len);
10 plot(vals, '--');
11
12 function z = avg(v)
13 z = sum(v)/l;
```

EXPRESSION INFO

`plot(vals, '--')`

Size: 1 × 1

Class: mxArray

plot is an extrinsic function.

Scope of Extrinsic Function Declarations

The `coder.extrinsic` construct has function scope. For example, consider the following code:

```
function y = foo %#codegen
coder.extrinsic('rat','min');
[N D] = rat(pi);
y = 0;
y = min(N, D);
```

In this example, `rat` and `min` are treated as extrinsic every time they are called in the main function `foo`. There are two ways to narrow the scope of an extrinsic declaration inside the main function:

- Declare the MATLAB function extrinsic in a local function, as in this example:

```
function y = foo %#codegen
coder.extrinsic('rat');
[N D] = rat(pi);
y = 0;
y = mymin(N, D);

function y = mymin(a,b)
coder.extrinsic('min');
y = min(a,b);
```

Here, the function `rat` is extrinsic every time it is called inside the main function `foo`, but the function `min` is extrinsic only when called inside the local function `mymin`.

- Instead of using the `coder.extrinsic` construct, call the MATLAB function using `feval`. This approach is described in the next section.

Extrinsic Declaration for Nonstatic Methods

Suppose that you define a class `myClass` that has a nonstatic method `foo`, and then create an instance `obj` of this class. If you want to declare the method `obj.foo` as extrinsic in your MATLAB code that you intend for code generation, follow these rules:

- Write the call to `foo` as a function call. Do not write the call by using the dot notation.
- Declare `foo` to be extrinsic by using the syntax `coder.extrinsic('foo')`.

For example, define `myClass` as:

```
classdef myClass
    properties
        prop = 1
    end
    methods
        function y = foo(obj,x)
            y = obj.prop + x;
        end
    end
end
```

Here is an example MATLAB function that declares `foo` as extrinsic.

```
function y = myFunction(x) %#codegen
coder.extrinsic('foo');
```

```
obj = myClass;
y = foo(obj,x);
end
```

Nonstatic methods are also known as ordinary methods. See “Method Syntax”.

Additional Uses

Use the `coder.extrinsic` construct to:

- Call MATLAB functions that do not produce output during simulation without generating unnecessary code.
- Make your code self-documenting and easier to debug. You can scan the source code for `coder.extrinsic` statements to isolate calls to MATLAB functions, which can potentially create and propagate `mxArrays`. See “Working with `mxArrays`” (MATLAB Coder).

Call MATLAB Functions Using `feval`

To narrow the scope of extrinsic declaration to just one function call, use the function `feval`. `feval` is automatically interpreted as an extrinsic function during code generation. So, you can use `feval` to call functions that you want to execute in the MATLAB environment, rather than compile to generated code.

Consider this example:

```
function y = foo
coder.extrinsic('rat');
[N D] = rat(pi);
y = 0;
y = feval('min',N,D);
```

Because `feval` is extrinsic, the statement `feval('min',N,D)` is evaluated by MATLAB, not compiled, which has the same result as declaring the function `min` extrinsic for just this one call. By contrast, the function `rat` is extrinsic throughout the function `foo`.

The code generator does not support the use of `feval` to call local functions or functions that are located in a private folder.

Working with `mxArrays`

The run-time output of an extrinsic function is an `mxArray`, also known as a MATLAB array. The only valid operations for `mxArrays` are:

- Storing an `mxArray` in a variable.
- Passing an `mxArray` to an extrinsic function.
- Returning an `mxArray` from a function back to MATLAB.
- Converting an `mxArray` to a known type at run time. Assign the `mxArray` to a variable whose type is already defined by a prior assignment. See the following example.

To use an `mxArray` returned by an extrinsic function in other operations (for example, returning it from a MATLAB Function block to Simulink execution), you must first convert it to a known type.

If the input arguments of a function are `mxArrays`, the code generator automatically treats the function as extrinsic.

Convert `mxArrays` to Known Types

To convert an `mxArray` to a known type, assign the `mxArray` to a variable whose type is defined. At run time, the `mxArray` is converted to the type of the variable that it is assigned to. If the data in the `mxArray` is not consistent with the type of the variable, you get a run-time error.

For example, consider this code:

```
function y = foo %#codegen
coder.extrinsic('rat');
[N D] = rat(pi);
y = min(N,D);
```

Here, the top-level function `foo` calls the extrinsic MATLAB function `rat`, which returns two `mxArrays` representing the numerator `N` and denominator `D` of the rational fraction approximation of `pi`. You can pass these `mxArrays` to another MATLAB function, in this case, `min`. Because the inputs passed to `min` are `mxArrays`, the code generator automatically treats `min` as an extrinsic function. As a result, `min` returns an `mxArray`.

While generating a MEX function by using MATLAB Coder, you can directly assign this `mxArray` returned by `min` to the output `y` because the MEX function returns its output to MATLAB.

```
codegen foo
```

```
Code generation successful.
```

But if you put `foo` in a MATLAB Function block in a Simulink model and then update or run the model, you get this error:

```
Function output 'y' cannot be an mxArray in this context.
Consider preinitializing the output variable with a known type.
```

This error occurs because returning an `mxArray` back to Simulink is not supported. To fix this issue, define `y` to be the type and size of the value that you expect `min` to return, in this case, a scalar double:

```
function y = foo %#codegen
coder.extrinsic('rat');
[N D] = rat(pi);
y = 0; % Define y as a scalar of type double
y = min(N,D);
```

In this example, the output of the extrinsic function `min` affects the output `y` of the entry-point function `foo` for which you are generating code. If you attempt to generate standalone code (for example, a static library) for `foo`, the code generator is unable to ignore the extrinsic function call and produces a code generation error.

```
codegen -config:lib foo
```

```
??? The extrinsic function 'min' is not available for
standalone code generation. It must be eliminated for
stand-alone code to be generated. It could not be
eliminated because its outputs appear to influence the
calling function. Fix this error by not using 'min'
```

or by ensuring that its outputs are unused.

```
Error in ==> foo Line: 4 Column: 5  
Code generation failed: View Error Report
```

```
Error using codegen
```

Restrictions on Using Extrinsic Functions

The full MATLAB run-time environment is not supported during code generation. Therefore, the following restrictions apply when calling MATLAB functions extrinsically:

- MATLAB functions that inspect the caller, or read or write to the caller workspace, do not work during code generation. Such functions include:
 - `dbstack`
 - `evalin`
 - `assignin`
 - `save`
- Functions in generated code can produce unpredictable results if your extrinsic function performs these actions at run time:
 - Changes folders
 - Changes the MATLAB path
 - Deletes or adds MATLAB files
 - Changes warning states
 - Changes MATLAB preferences
 - Changes Simulink parameters
- The code generator does not support the use of `coder.extrinsic` to call functions that are located in a private folder.
- The code generator does not support the use of `coder.extrinsic` to call local functions.
- You can call extrinsic functions with up to 64 inputs and 64 outputs.

See Also

`coder.extrinsic` | `coder.const`

Code Generation for Recursive Functions

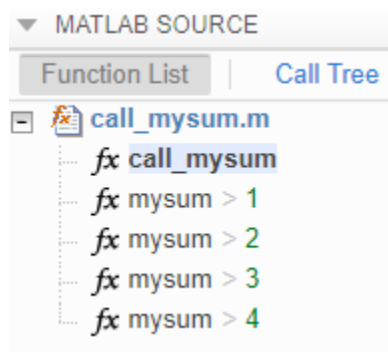
To generate code for recursive MATLAB functions, the code generator uses compile-time recursion on page 14-12 or run-time recursion on page 14-12. You can influence whether the code generator uses compile-time or run-time recursion by modifying your MATLAB code. See “Force Code Generator to Use Run-Time Recursion” on page 14-14.

You can disallow recursion on page 14-13 or disable run-time recursion on page 14-13 by modifying configuration parameters.

When you use recursive functions in MATLAB code that is intended for code generation, you must adhere to certain restrictions. See “Recursive Function Limitations for Code Generation” on page 14-13.

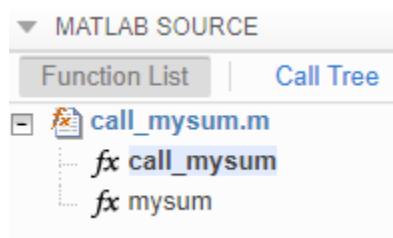
Compile-Time Recursion

With compile-time recursion, the code generator creates multiple versions of a recursive function in the generated code. The inputs to each version have values or sizes that are customized for that version. These versions are known as function specializations. You can see if the code generator used compile-time recursion by looking at the code generation report. Here is an example of compile-time recursion in the report.



Run-Time Recursion

With run-time recursion, the code generator produces a recursive function in the generated code. You can see if the code generator used run-time recursion by looking at the code generation report. Here is an example of run-time recursion in the report.



Disallow Recursion

In a code acceleration configuration object, set the value of the `CompileTimeRecursionLimit` configuration parameter to 0.

Disable Run-Time Recursion

Some coding standards, such as MISRA, do not allow recursion. To increase the likelihood of generating code that is compliant with MISRA C™, disable run-time recursion.

In a code acceleration configuration object, set `EnableRuntimeRecursion` to `false`.

If your code requires run-time recursion and run-time recursion is disabled, you must rewrite your code so that it uses compile-time recursion or does not use recursion.

Recursive Function Limitations for Code Generation

When you use recursion in MATLAB code that is intended for code generation, follow these restrictions:

- Assign all outputs of a run-time recursive function before the first recursive call in the function.
- Assign all elements of cell array outputs of a run-time recursive function.
- Inputs and outputs of run-time recursive functions cannot be classes.
- The `StackUsageMax` code acceleration configuration parameter is ignored for run-time recursion.

See Also

Related Examples

- “Force Code Generator to Use Run-Time Recursion” on page 14-14
- “Compile-Time Recursion Limit Reached” on page 49-33
- “Output Variable Must Be Assigned Before Run-Time Recursive Call” on page 49-36
- “Set Up C Compiler and Compilation Options” on page 12-16
- “Code Generation Reports” on page 12-27

Force Code Generator to Use Run-Time Recursion

When your MATLAB code includes recursive function calls, the code generator uses compile-time or run-time recursion. With compile-time recursion on page 14-12, the code generator creates multiple versions of the recursive function in the generated code. These versions are known as function specializations. With run-time recursion on page 14-12, the code generator produces a recursive function. If compile-time recursion results in too many function specializations or if you prefer run-time recursion, you can try to force the code generator to use run-time recursion. Try one of these approaches:

- “Treat the Input to the Recursive Function as a Nonconstant” on page 14-14
- “Make the Input to the Recursive Function Variable-Size” on page 14-15
- “Assign Output Variable Before the Recursive Call” on page 14-16

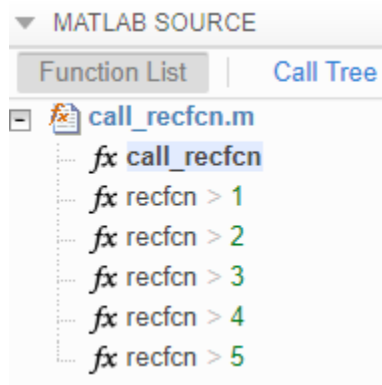
Treat the Input to the Recursive Function as a Nonconstant

Consider this function:

```
function y = call_recfcn(n)
A = ones(1,n);
x = 5;
y = recfcn(A,x);
end

function y = recfcn(A,x)
if size(A,2) == 1 || x == 1
    y = A(1);
else
    y = A(1)+recfcn(A(2:end),x-1);
end
end
```

`call_recfcn` calls `recfcn` with the value 5 for the second argument. `recfcn` calls itself recursively until `x` is 1. For each `recfcn` call, the input argument `x` has a different value. The code generator produces five specializations of `recfcn`, one for each call.



To force run-time recursion, in `call_recfcn`, in the call to `recfcn`, instruct the code generator to treat the value of the input argument `x` as a nonconstant value by using `coder.ignoreConst`.

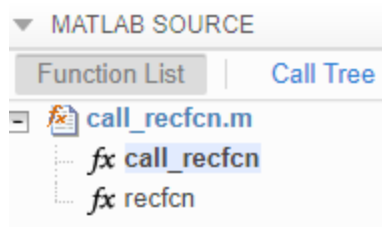

```

function y = call_recfcn(n)
A = ones(1,n);
x = coder.ignoreConst(5);
y = recfcn(A,x);
end

function y = recfcn(A,x)
if size(A,2) == 1 || x == 1
    y = A(1);
else
    y = A(1)+recfcn(A(2:end),x-1);
end
end

```

, you see only one specialization.



Make the Input to the Recursive Function Variable-Size

Consider this code:

```

function z = call_mysum(A)
%#codegen
z = mysum(A);
end

function y = mysum(A)
coder.inline('never');
if size(A,2) == 1
    y = A(1);
else
    y = A(1)+mysum(A(2:end));
end
end

```

If the input to mysum is fixed-size, the code generator uses compile-time recursion. To force the code generator to use run-time conversion, make the input to mysum variable-size by using `coder.varsize`.

```

function z = call_mysum(A)
%#codegen
B = A;
coder.varsize('B');
z = mysum(B);
end

function y = mysum(A)
coder.inline('never');
if size(A,2) == 1

```

```
        y = A(1);  
else  
    y = A(1)+ mysum(A(2:end));  
end  
end
```

Assign Output Variable Before the Recursive Call

The code generator uses compile-time recursion for this code:

```
function y = callrecursive(n)  
x = 10;  
y = myrecursive(x,n);  
end  
  
function y = myrecursive(x,n)  
coder.inline('never')  
if x > 1  
    y = n + myrecursive(x-1,n-1);  
  
else  
    y = n;  
end  
end
```

To force the code generator to use run-time recursion, modify `myrecursive` so that the output `y` is assigned before the recursive call. Place the assignment `y = n` in the `if` block and the recursive call in the `else` block.

```
function y = callrecursive(n)  
x = 10;  
y = myrecursive(x,n);  
end  
  
function y = myrecursive(x,n)  
coder.inline('never')  
if x == 1  
    y = n;  
else  
    y = n + myrecursive(x-1,n-1);  
end  
end
```

See Also

More About

- “Code Generation for Recursive Functions” on page 14-12
- “Output Variable Must Be Assigned Before Run-Time Recursive Call” on page 49-36
- “Compile-Time Recursion Limit Reached” on page 49-33

Avoid Duplicate Functions in Generated Code

Issue

You generate code and it contains multiple, duplicate copies of the same functions, with only slight differences, such as modifications to the function signature. For example, your generated code might contain functions called `foo` and `b_foo`. Duplicate functions can make the generated code more difficult to analyze and manage.

Cause

Duplicate functions in the generated code are the result of function specializations. The code generator specializes functions when it detects that they differ at different call sites by:

- Number of input or output variables.
- Type of input or output variables.
- Size of input or output variables.
- Values of input variables.

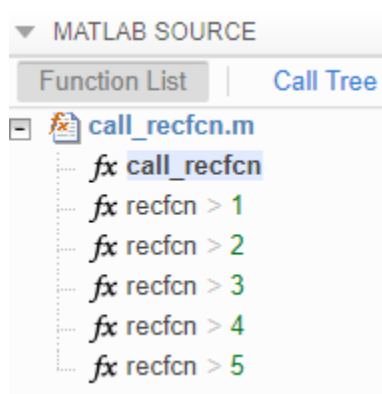
In some cases, these specializations are necessary for the generated C/C++ code because C/C++ functions do not have the same flexibility as MATLAB functions. In other cases, the code generator specializes functions to optimize the generated code or because of a lack of information.

Solution

In certain cases, you can alter your MATLAB code to avoid the generation of duplicate functions.

Identify Duplicate Functions by Using Code Generation Report

You can determine whether the code generator created duplicate functions by inspecting the code generation report or in Simulink, the MATLAB Function report. The report shows a list of the duplicate functions underneath the entry-point function. For example:



Duplicate Functions Generated for Multiple Input Sizes

If your MATLAB code calls a function multiple times and passes inputs of different sizes, the code generator can create specializations of the function for each size. To avoid this issue, use on the

function input. For example, this code uses `coder.ignoreSize` to avoid creating multiple copies of the function `indexOf`:

```
function [out1, out2] = test1(in)
    a = 1:10;
    b = 2:40;
    % Without coder.ignoreSize duplicate functions are generated
    out1 = indexOf(coder.ignoreSize(a), in);
    out2 = indexOf(coder.ignoreSize(b), in);
end

function index = indexOf(array, value)
    coder.inline('never');
    for i = 1:numel(array)
        if array(i) == value
            index = i;
            return
        end
    end
    index = -1;
    return
end
```

To generate code, enter:

```
codegen test1 -config:lib -report -args {1}
```

Duplicate Functions Generated for Different Input Values

If your MATLAB code calls a function and passes multiple different constant inputs, the code generator can create specializations of the function for each different constant. In this case, use to indicate to the code generator not to treat the value as an immutable constant. For example:

```
function [out3, out4] = test2(in)
    c = ['a', 'b', 'c'];
    if in > 0
        c(2)='d';
    end
    out3 = indexOf(c, coder.ignoreConst('a'));
    out4 = indexOf(c, coder.ignoreConst('b'));
end

function index = indexOf(array, value)
    coder.inline('never');
    for i = 1:numel(array)
        if array(i) == value
            index = i;
            return
        end
    end
    index = -1;
    return
end
```

To generate code, enter:

```
codegen test2 -config:lib -report -args {1}
```

Duplicate Functions Generated for Different Number of Outputs

If your MATLAB code calls a function and accepts a different number of outputs at different call sites, the code generator can produce specializations for each call. For example:

```
[a b] = foo();  
c = foo();
```

To make each call return the same number of outputs and avoid duplicate functions, use the ~ symbol:

```
[a b] = foo();  
[c, ~] = foo();
```


Code Generation for MATLAB Classes

- “MATLAB Classes Definition for Code Generation” on page 15-2
- “Classes That Support Code Generation” on page 15-7
- “Generate Code for MATLAB Value Classes” on page 15-8
- “Generate Code for MATLAB Handle Classes and System Objects” on page 15-12
- “Code Generation for Handle Class Destructors” on page 15-15
- “Class Does Not Have Property” on page 15-18
- “Handle Object Limitations for Code Generation” on page 15-19
- “System Objects in MATLAB Code Generation” on page 15-23
- “Specify Objects as Inputs” on page 15-26
- “Work Around Language Limitation: Code Generation Does Not Support Object Arrays” on page 15-29

MATLAB Classes Definition for Code Generation

To generate efficient standalone code for MATLAB classes, you must use classes differently than when running your code in the MATLAB environment.

Language Limitations

Although code generation support is provided for common features of classes such as properties and methods, there are a number of advanced features which are not supported, such as:

- Events
- Listeners
- Arrays of objects
- Recursive data structures
 - Linked lists
 - Trees
 - Graphs
- Nested functions in constructors
- Overloadable operators `subsref`, `subsassign`, and `subsindex`

In MATLAB, classes can define their own versions of the `subsref`, `subsassign`, and `subsindex` methods. Code generation does not support classes that have their own definitions of these methods.

- The empty method

In MATLAB, classes have a built-in static method, `empty`, which creates an empty array of the class. Code generation does not support this method.

- The following MATLAB handle class methods:
 - `addlistener`
 - `eq`
 - `findobj`
 - `findprop`
- The `AbortSet` property attribute

Code Generation Features Not Compatible with Classes

- You can generate code for entry-point MATLAB functions that use classes, but you cannot generate code directly for a MATLAB class.

For example, if `ClassNameA` is a class definition, you cannot generate code by executing:

```
codegen ClassNameA
```

- A handle class object cannot be an entry-point function input or output.
- A value class object can be an entry-point function input or output. However, if a value class object contains a handle class object, then the value class object cannot be an entry-point function input or output. A handle class object cannot be an entry-point function input or output.

- Code generation does not support global variables that are handle classes.
- Code generation does not support multiple outputs from constructors.
- Code generation does not support assigning an object of a value class into a nontunable property. For example, `obj.prop=v`; is invalid when `prop` is a nontunable property and `v` is an object based on a value class.
- You cannot use `coder.extrinsic` to declare a class or method as extrinsic.
- If an object has duplicate property names and the code generator tries to constant-fold the object, code generation can fail. The code generator constant-folds an object when it is used with `coder.Constant` or `coder.const`, or when it is an input to or output from a constant-folded extrinsic function.

Duplicate property names occur in an object of a subclass in these situations:

- The subclass has a property with the same name as a property of the superclass.
- The subclass derives from multiple superclasses that use the same name for a property.

Duplicate property names must be consistently constant or non-constant across multiple inheritance related classes. For example, code generation produces an error if an object with a constant property `aProp` inherits `aProp` from a superclass where `aProp` is defined as non-constant.

For information about when MATLAB allows duplicate property names, see “Subclassing Multiple Classes”.

Defining Class Properties for Code Generation

For code generation, you must define class properties differently than you do when running your code in the MATLAB environment:

- To test property validation, it is a best practice to run a MEX function over the full range of input values.
- After defining a property, do not assign it an incompatible type. Do not use a property before attempting to grow it.

When you define class properties for code generation, consider the same factors that you take into account when defining variables. In the MATLAB language, variables can change their class, size, or complexity dynamically at run time so you can use the same variable to hold a value of varying class, size, or complexity. C and C++ use static typing. Before using variables, to determine their type, the code generator requires a complete assignment to each variable. Similarly, before using properties, you must explicitly define their class, size, and complexity.

- Initial values:
 - If the property does not have an explicit initial value, the code generator assumes that it is undefined at the beginning of the constructor. The code generator does not assign an empty matrix as the default.
 - If the property does not have an initial value and the code generator cannot determine that the property is assigned prior to first use, the software generates a compilation error.
 - For System objects, if a nontunable property is a structure, you must completely assign the structure. You cannot do partial assignment using subscripting.

For example, for a nontunable property, you can use the following assignment:

```
mySystemObject.nonTunableProperty=struct('fieldA','a','fieldB','b');
```

You cannot use the following partial assignments:

```
mySystemObject.nonTunableProperty.fieldA = 'a';
mySystemObject.nonTunableProperty.fieldB = 'b';
```

- `coder. varsize` is not supported for class properties.
- If the initial value of a property is an object, then the property must be constant. To make a property constant, declare the `Constant` attribute in the property block. For example:

```
classdef MyClass
    properties (Constant)
        p1 = MyClass2;
    end
end
```

Code generation does not support a constant property that is assigned to an object that contains a System object.

- MATLAB computes class initial values at class loading time before code generation. If you use persistent variables in MATLAB class property initialization, the value of the persistent variable computed when the class loads belongs to MATLAB; it is not the value used at code generation time. If you use `coder.target` in MATLAB class property initialization, `coder.target('MATLAB')` returns `true (1)`.
- Variable-size properties:
 - Code generation supports upper-bounded and unbounded variable-size properties for both value and handle classes.
 - To generate unbounded variable-size class properties, enable dynamic memory allocation.
 - To make a variable-size class property, make two sequential assignments of a class property, one to a scalar and the next to an array.

```
classdef varSizeProp1 < handle
    properties
        prop
        varProp
    end
end

function extFunc(n)
    obj = varSizeProp1;
    % Assign a scalar value to the property.
    obj.prop = 1;
    obj.varProp = 1;
    % Assign an array to the same property to make it variable-sized.
    obj.prop = 1:98;
    obj.varProp = 1:n;
end
```

In the preceding code, the first assignment to `prop` and `varProp` is scalar, and their second assignment is to an array with the same base type. The size of `prop` has an upper bound of 98, making it an upper-bounded, variable-size property.

If `n` is unknown at compile time, `obj.varProp` is an unbounded variable-size property. If it is known, it is an upper-bounded, variable-size class property.

- If the class property is initialized with a variable-size array, the property is variable-size.

```
classdef varSizeProp2
    properties
        prop
    end
```

```

methods
function obj = varSizeProp2(inVar)
    % Assign incoming value to local variable
    locVar = inVar;

    % Declare the local variable to be a variable-sized column
    % vector with no size limit
    coder.varsize('locVar',[inf 1],[1 0]);

    % Assign value
    obj.prop = locVar;
end
end
end

```

In the preceding code, `inVar` is passed to the class constructor and stored in `locVar`. `locVar` is modified to be variable-size by `coder.varsize` and assigned to the class property `obj.prop`, which makes the property variable-size.

- If the input to the function call `varSizeProp2` is variable-size, `coder.varsize` is not required.

```

function z = constructCall(n)
    z = varSizeProp2(1:n);
end

```

- If the value of `n` is unknown at compile-time and has no specified bounds, `z.prop` is an unbounded variable-size class property.
- If the value of `n` is unknown at compile-time and has specified bounds, `z.prop` is an upper-bounded variable-size class property.
- If a property is constant and its value is an object, you cannot change the value of a property of that object. For example, suppose that:
 - `obj` is an object of `myClass1`.
 - `myClass1` has a constant property `p1` that is an object of `myClass2`.
 - `myClass2` has a property `p2`.

Code generation does not support the following code:

```
obj.p1.p2 = 1;
```

Inheritance from Built-In MATLAB Classes Not Supported

You cannot generate code for classes that inherit from built-in MATLAB classes. For example, you cannot generate code for the following class:

```
classdef myclass < double
```

An exception to this rule is the MATLAB enumeration class. You can generate code for enumeration classes that inherit from built-in MATLAB classes. See “Code Generation for Enumerations” (MATLAB Coder).

See Also

`coder.target`

Related Examples

- “Generate Standalone C/C++ Code That Detects and Reports Run-Time Errors” (MATLAB Coder)
- “Classes That Support Code Generation”

Classes That Support Code Generation

You can generate code for MATLAB value and handle classes and user-defined System objects. Your class can have multiple methods and properties and can inherit from multiple classes.

| To generate code for: | Example: |
|--|--|
| Value classes | "Generate Code for MATLAB Value Classes" on page 15-8 |
| Handle classes including user-defined System objects | "Generate Code for MATLAB Handle Classes and System Objects" on page 15-12 |

For more information, see:

- "Role of Classes in MATLAB"
- "MATLAB Classes Definition for Code Generation" on page 15-2

Generate Code for MATLAB Value Classes

This example shows how to generate code for a MATLAB value class and then view the generated code in the code generation report.

- 1 In a writable folder, create a MATLAB value class, Shape. Save the code as Shape.m.

```

classdef Shape
% SHAPE Create a shape at coordinates
% centerX and centerY
    properties
        centerX;
        centerY;
    end
    properties (Dependent = true)
        area;
    end
    methods
        function out = get.area(obj)
            out = obj.getarea();
        end
        function obj = Shape(centerX,centerY)
            obj.centerX = centerX;
            obj.centerY = centerY;
        end
    end
end
methods(Abstract = true)
    getarea(obj);
end
methods(Static)
    function d = distanceBetweenShapes(shape1,shape2)
        xDist = abs(shape1.centerX - shape2.centerX);
        yDist = abs(shape1.centerY - shape2.centerY);
        d = sqrt(xDist^2 + yDist^2);
    end
end
end
end

```

- 2 In the same folder, create a class, Square, that is a subclass of Shape. Save the code as Square.m.

```

classdef Square < Shape
% Create a Square at coordinates center X and center Y
% with sides of length of side
    properties
        side;
    end
    methods
        function obj = Square(side,centerX,centerY)
            obj@Shape(centerX,centerY);
            obj.side = side;
        end
        function Area = getarea(obj)
            Area = obj.side^2;
        end
    end
end
end
end

```

- 3 In the same folder, create a class, Rhombus, that is a subclass of Shape. Save the code as Rhombus.m.

```

classdef Rhombus < Shape
    properties
        diag1;
        diag2;
    end
    methods
        function obj = Rhombus(diag1,diag2,centerX,centerY)
            obj@Shape(centerX,centerY);
            obj.diag1 = diag1;
            obj.diag2 = diag2;
        end
        function Area = getarea(obj)
            Area = 0.5*obj.diag1*obj.diag2;
        end
    end
end
end

```

- 4 Write a function that uses this class.

```

function [TotalArea, Distance] = use_shape
    %#codegen
    s = Square(2,1,2);
    r = Rhombus(3,4,7,10);
    TotalArea = s.area + r.area;
    Distance = Shape.distanceBetweenShapes(s,r);

```

- 5 Generate a static library for use_shape and generate a code generation report.

```

codegen -config:lib -report use_shape

```

codegen generates a C static library with the default name, use_shape, and supporting files in the default folder, codegen/lib/use_shape.

- 6 Click the **View report** link.

- 7 To see the Rhombus class definition, on the **MATLAB Source** pane, under Rhombus.m, click Rhombus. The Rhombus class constructor is highlighted.

- 8 Click the **Variables** tab. You see that the variable obj is an object of the Rhombus class. To see its properties, expand obj.

The screenshot shows the MATLAB Coder Report Viewer interface. The 'MATLAB Source' pane on the left displays a 'Call Tree' view. The main pane shows the source code for the 'Rhombus' class constructor, which inherits from 'Shape'. The code includes properties for 'diag1' and 'diag2', and methods for 'obj' and 'Area'. The 'Generated Code' pane on the left shows a list of source files. The bottom pane displays a table with columns for Name, Type, Size, and Class.

| Name | Type | Size | Class |
|---------|--------|-------|---------|
| obj | Output | 1 × 1 | Rhombus |
| centerX | | 1 × 1 | double |
| centerY | | 1 × 1 | double |
| diag1 | | 1 × 1 | double |
| diag2 | | 1 × 1 | double |
| centerX | Input | 1 × 1 | double |
| centerY | Input | 1 × 1 | double |
| diag1 | Input | 1 × 1 | double |
| diag2 | Input | 1 × 1 | double |

- 9 In the **MATLAB Source** pane, click **Call Tree**.

The **Call Tree** view shows that `use_shape` calls the `Rhombus` constructor and that the `Rhombus` constructor calls the `Shape` constructor.

The image shows a close-up of the 'MATLAB SOURCE' pane with the 'Call Tree' view selected. The tree shows the following structure:

- fx use_shape
 - fx Square/Square
 - fx Rhombus/Rhombus
 - fx Shape/Shape
 - fx Shape/get.area
 - fx Shape/get.area
 - fx Shape/distanceBetweenShapes

- 10 In the code pane, in the `Rhombus` class constructor, move your pointer to this line:

```
obj@Shape(centerX,centerY)
```

The `Rhombus` class constructor calls the `Shape` method of the base `Shape` class. To view the `Shape` class definition, in `obj@Shape`, double-click `Shape`.


```
Shape.m
1 classdef Shape
2 % SHAPE Create a shape at coordinates
3 % centerX and centerY
4     properties
5         centerX;
6         centerY;
7     end
8     properties (Dependent = true)
9         area;
10    end
11    methods
12        function out = get.area(obj)
13            out = obj.getarea();
14        end
15        function obj = Shape(centerX,centerY)
16            obj.centerX = centerX;
17            obj.centerY = centerY;
18        end
19    end
20    methods(Abstract = true)
21        getarea(obj);
22    end
23    methods(Static)
24        function d = distanceBetweenShapes(shape1,shape2)
25            xDist = abs(shape1.centerX - shape2.centerX);
26            yDist = abs(shape1.centerY - shape2.centerY);
27            d = sqrt(xDist^2 + yDist^2);
28        end
29    end
30 end
31
32
```

Generate Code for MATLAB Handle Classes and System Objects

This example shows how to generate code for a user-defined System object and then view the generated code in the code generation report.

- 1 In a writable folder, create a System object, `AddOne`, which subclasses from `matlab.System`. Save the code as `AddOne.m`.

```
classdef AddOne < matlab.System
% ADDONE Compute an output value that increments the input by one

    methods (Access=protected)
        % stepImpl method is called by the step method
        function y = stepImpl(~,x)
            y = x+1;
        end
    end
end
```

- 2 Write a function that uses this System object.

```
function y = testAddOne(x)
%#codegen
    p = AddOne();
    y = p.step(x);
end
```

- 3 Generate a MEX function for this code.

```
codegen -report testAddOne -args {0}
```

The `-report` option instructs `codegen` to generate a code generation report, even if no errors or warnings occur. The `-args` option specifies that the `testAddOne` function takes one scalar double input.

- 4 Click the **View report** link.
- 5 In the **MATLAB Source** pane, click `testAddOne`. To see information about the variables in `testAddOne`, click the **Variables** tab.

The screenshot shows the MATLAB IDE interface. The top toolbar includes navigation and editing tools. The left pane shows the MATLAB SOURCE tree with 'testAddOne.m' selected. The main editor displays the code for 'testAddOne.m':

```

1 function y = test
2 %#codegen
3 p = AddOne();
4 y = p.step(x);
5 end

```

A tooltip for variable 'p' is visible, showing its properties: Size: 1 x 1, Class: AddOne. Below the code editor, a SUMMARY table is displayed:

| Name | Type | Size | Class |
|------|--------|-------|--------|
| y | Output | 1 x 1 | double |
| x | Input | 1 x 1 | double |
| p | Local | 1 x 1 | AddOne |

6 To view the class definition for addOne, in the **MATLAB Source** pane, click AddOne.

The screenshot shows the MATLAB IDE interface with the 'AddOne.m' file selected in the MATLAB SOURCE pane. The main editor displays the class definition for 'AddOne.m':

```

1 classdef AddOne < matlab.System
2 % ADDONE Compute an output value that increments the input by one
3
4 methods (Access=protected)
5 % stepImpl method is called by the step method
6 function y = stepImpl(~,x)
7     y = x+1;
8 end
9 end
10 end

```

See Also

More About

- “Code Generation for Handle Class Destructors” on page 15-15

Code Generation for Handle Class Destructors

You can generate code for MATLAB code that uses `delete` methods (destructors) for handle classes. To perform clean-up operations, such as closing a previously opened file before an object is destroyed, use a `delete` method. The generated code calls the `delete` method at the end of an object's lifetime, even if execution is interrupted by a run-time error. When System objects are destroyed, `delete` calls the `release` method, which in turn calls the user-defined `releaseImpl`. For more information on when to define a `delete` method in a MATLAB code, see "Handle Class Destructor".

Guidelines and Restrictions

When you write the MATLAB code, adhere to these guidelines and restrictions:

- Code generation does not support recursive calls of the `delete` method. Do not create an object of a certain class inside the `delete` method for the same class. This usage might cause a recursive call of `delete` and result in an error message.
- The generated code always calls the `delete` method, when an object goes out of scope. Code generation does not support explicit calls of the `delete` method.
- Initialize all properties of `MyClass` that the `delete` method of `MyClass` uses either in the constructor or as the default property value. If `delete` tries to access a property that has not been initialized in one of these two ways, the code generator produces an error message.
- Suppose a property `prop1` of `MyClass1` is itself an object (an instance of another class `MyClass2`). Initialize all properties of `MyClass2` that the `delete` method of `MyClass1` uses. Perform this initialization either in the constructor of `MyClass2` or as the default property value. If `delete` tries to access a property of `MyClass2` that has not been initialized in one of these two ways, the code generator produces an error message. For example, define the two classes `MyClass1` and `MyClass2`:

```
classdef MyClass1 < handle
    properties
        prop1
    end
    methods
        function h = MyClass1(index)
            h.prop1 = index;
        end
        function delete(h)
            fprintf('h.prop1.prop2 is: %1.0f\n',h.prop1.prop2);
        end
    end
end

classdef MyClass2 < handle
    properties
        prop2
    end
end
```

Suppose you try to generate code for this function:

```
function MyFunction
obj2 = MyClass2;
```

```
obj1 = MyClass1(obj2); % Assign obj1.prop1 to the input (obj2)
end
```

The code generator produces an error message because you have not initialized the property `obj2.prop2` that the `delete` method displays.

Behavioral Differences of Objects in Generated Code and in MATLAB

The behavior of objects in the generated code can be different from their behavior in MATLAB in these situations:

- The order of destruction of several independent objects might be different in MATLAB than in the generated code.
- The lifetime of objects in the generated code can be different from their lifetime in MATLAB. MATLAB calls the `delete` method when an object can no longer be reached from any live variable. The generated code calls the `delete` method when an object goes out of scope. In some situations, this difference causes `delete` to be called later on in the generated code than in MATLAB. For example, define the class:

```
classdef MyClass < handle
    methods
        function delete(h)
            global g
            % Destructor displays current value of global variable g
            fprintf('The global variable is: %1.0f\n',g);
        end
    end
end
```

Run the function:

```
function MyFunction
    global g
    g = 1;
    obj = MyClass;
    obj = MyClass;
    % MATLAB destroys the first object here
    g = 2;
    % MATLAB destroys the second object here
    % Generated code destroys both objects here
end
```

The first object can no longer be reached from any live variable after the second instance of `obj = MyClass` in `MyFunction`. MATLAB calls the `delete` method for the first object after the second instance of `obj = MyClass` in `MyFunction` and for the second object at the end of the function. The output is:

```
The global variable is: 1
The global variable is: 2
```

In the generated code, both `delete` method calls happen at the end of the function when the two objects go out of scope. Running `MyFunction_mex` results in a different output:

```
The global variable is: 2
The global variable is: 2
```

- In MATLAB, persistent objects are automatically destroyed when they cannot be reached from any live variable. In the generated code, you have to call the `terminate` function explicitly to destroy the persistent objects.
- The generated code does not destroy partially constructed objects. If a handle object is not fully constructed at run time, the generated code produces an error message but does not call the `delete` method for that object. For a System object, if there is a run-time error in `setupImpl`, the generated code does not call `releaseImpl` for that object.

MATLAB does call the `delete` method to destroy a partially constructed object.

See Also

More About

- “Generate Code for MATLAB Handle Classes and System Objects” on page 15-12
- “System Objects in MATLAB Code Generation” on page 15-23

Class Does Not Have Property

If a MATLAB class has a method, `mymethod`, that returns a handle class with a property, `myprop`, you cannot generate code for the following type of assignment:

```
obj.mymethod().myprop=...
```

For example, consider the following classes:

```
classdef MyClass < handle
    properties
        myprop
    end
    methods
        function this = MyClass
            this.myprop = MyClass2;
        end
        function y = mymethod(this)
            y = this.myprop;
        end
    end
end

classdef MyClass2 < handle
    properties
        aa
    end
end
```

You cannot generate code for function `foo`.

```
function foo

h = MyClass;

h.mymethod().aa = 12;
```

In this function, `h.mymethod()` returns a handle object of type `MyClass2`. In MATLAB, the assignment `h.mymethod().aa = 12;` changes the property of that object. Code generation does not support this assignment.

Solution

Rewrite the code to return the object and then assign a value to a property of the object.

```
function foo

h = MyClass;

b=h.mymethod();
b.aa=12;
```


Handle Object Limitations for Code Generation

The code generator statically determines the lifetime of a handle object. When you use handle objects, this static analysis has certain restrictions.

With static analysis the generated code can reuse memory rather than rely on a dynamic memory management scheme, such as reference counting or garbage collection. The code generator can avoid dynamic memory allocation and run-time automatic memory management. These generated code characteristics are important for some safety-critical and real-time applications.

For limitations, see:

- “A Variable Outside a Loop Cannot Refer to a Handle Object Allocated Inside a Loop” on page 15-19
- “A Handle Object That a Persistent Variable Refers To Must Be a Singleton Object” on page 15-20

The code generator analyzes whether all variables are defined prior to use. Undefined variables or data types cause an error during code generation. In certain circumstances, the code generator cannot determine if references to handle objects are defined. See “References to Handle Objects Can Appear Undefined” on page 15-21.

A Variable Outside a Loop Cannot Refer to a Handle Object Allocated Inside a Loop

Consider the handle class `mycls` and the function `usehandle1`.

```
classdef mycls < handle
    properties
        prop
    end

    methods
        function obj = mycls(x)
            obj.prop = x;
        end
    end
end

function y = usehandle1
    p = mycls(0); % Instance of mycls with prop value 10 created

    for i = 1:10
        p = mycls(i); % Handle object allocated inside loop
    end

    y = p.prop; % Handle object referenced outside loop
end
```

If you try to generate code for the `usehandle1` function, the code generator produces an error. The error occurs because:

- A handle object is allocated inside the `for` loop. The variable `p.prop` refers to this handle object.

- Outside the loop, the variable `x` refers to the property `prop` handle object.

A Handle Object That a Persistent Variable Refers To Must Be a Singleton Object

If a persistent variable refers to a handle object, the code generator allows only one instance of the object during the program's lifetime. The object must be a singleton object. To create a singleton handle object, enclose statements that create the object in the `if isempty()` guard for the persistent variable.

For example, consider the class `mycls` and the function `usehandle2`. The code generator reports an error for `usehandle2` because `p.prop` refers to the `mycls` object that the statement `inner = mycls` creates. This statement creates a `mycls` object for each invocation of `usehandle2`.

```
classdef mycls < handle
    properties
        prop
    end
end

function usehandle2(x)
    assert(isa(x, 'double'));
    persistent p;
    inner = mycls;
    inner.prop = x;
    if isempty(p)
        p = mycls;
        p.prop = inner;
    end
end
```

If you move the statements `inner = mycls` and `inner.prop = x` inside the `if isempty()` guard, code generation succeeds. The statement `inner = mycls` executes only once during the program's lifetime.

```
function usehandle2(x)
    assert(isa(x, 'double'));
    persistent p;
    if isempty(p)
        inner = mycls;
        inner.prop = x;
        p = mycls;
        p.prop = inner;
    end
end
```

Consider the function `usehandle3`. The code generator reports an error for `usehandle3` because the persistent variable `p` refers to the `mycls` object that the statement `myobj = mycls` creates. This statement creates a `mycls` object for each invocation of `usehandle3`.

```
function usehandle3(x)
    assert(isa(x, 'double'));
    myobj = mycls;
    myobj.prop = x;
    doinit(myobj);
    disp(myobj.prop);
    function doinit(obj)
        persistent p;
    end
end
```

```

if isempty(p)
    p = obj;
end

```

If you make `myobj` persistent and enclose the statement `myobj = mycls` inside an `if isempty()` guard, code generation succeeds. The statement `myobj = mycls` executes only once during the program's lifetime.

```

function usehandle3(x)
assert(isa(x, 'double'));
persistent myobj;
if isempty(myobj)
    myobj = mycls;
end

```

```

doinit(myobj);

```

```

function doinit(obj)
persistent p;
if isempty(p)
    p = obj;
end

```

References to Handle Objects Can Appear Undefined

Consider the function `refHandle` that copies a handle object property to another object. The function uses a simple handle class and value class. In MATLAB, the function runs without error.

```

function [out1, out2, out3] = refHandle()
    x = myHandleClass;
    y = x;
    v = myValueClass();
    v.prop = x;
    x.prop = 42;
    out1 = x.prop;
    out2 = y.prop;
    out3 = v.prop.prop;
end

```

```

classdef myHandleClass < handle
    properties
        prop
    end
end

```

```

classdef myValueClass
    properties
        prop
    end
end

```

During code generation, an error occurs:

```
Property 'v.prop.prop' is undefined on some execution paths.
```

Three variables reference the same memory location: `x`, `y`, and `v.prop`. The code generator determines that `x.prop` and `y.prop` share the same value. The code generator cannot determine

that the handle object property `v.prop.prop` shares its definition with `x.prop` and `y.prop`. To avoid the error, define `v.prop.prop` directly.

System Objects in MATLAB Code Generation

In this section...

“Usage Rules and Limitations for System Objects for Generating Code” on page 15-23

“System Objects in codegen” on page 15-25

“System Objects in the MATLAB Function Block” on page 15-25

“System Objects in the MATLAB System Block” on page 15-25

“System Objects and MATLAB Compiler Software” on page 15-25

You can generate C/C++ code in MATLAB from your system that contains System objects by using MATLAB Coder. You can generate efficient and compact code for deployment in desktop and embedded systems and accelerate fixed-point algorithms.

Usage Rules and Limitations for System Objects for Generating Code

The following usage rules and limitations apply to using System objects in code generated from MATLAB.

Object Construction and Initialization

- If objects are stored in persistent variables, initialize System objects once by embedding the object handles in an `if` statement with a call to `isempty()`.
- Set arguments to System object constructors as compile-time constants.
- Initialize all System objects properties that `releaseImpl` uses before the end of `setupImpl`.
- You cannot initialize System objects properties with other MATLAB class objects as default values in code generation. You must initialize these properties in the constructor.

Inputs and Outputs

- System objects accept a maximum of 1024 inputs. A maximum of eight dimensions per input is supported.
- The data type of the inputs should not change.
- The complexity of the inputs should not change.
- If you want the size of inputs to change, verify that support for variable-size is enabled. Code generation support for variable-size data also requires that variable-size support is enabled. By default in MATLAB, support for variable-size data is enabled.
- System objects predefined in the software do not support variable-size if their data exceeds the `DynamicMemoryAllocationThreshold` value.
- Do not set System objects to become outputs from the MATLAB Function block.
- Do not use the Save and Restore Simulation Operating Point option for any System object in a MATLAB Function block.
- Do not pass a System object as an example input argument to a function being compiled with `codegen`.
- Do not pass a System object to functions declared as extrinsic (functions called in interpreted mode) using the `coder.extrinsic` function. System objects returned from extrinsic functions and scope System objects that automatically become extrinsic can be used as inputs to another extrinsic function. But, these functions do not generate code.

Properties

- In MATLAB System blocks, you cannot use variable-size for discrete state properties of System objects. Private properties can be variable-size.
- Objects cannot be used as default values for properties.
- You can only assign values to nontunable properties once, including the assignment in the constructor.
- Nontunable property values must be constant.
- For fixed-point inputs, if a tunable property has dependent data type properties, you can set tunable properties only at construction time or after the object is locked.
- For `getNumInputsImpl` and `getNumOutputsImpl` methods, if you set the return argument from an object property, that object property must have the `Nontunable` attribute.

Global Variables

- Global variables are allowed in a System object, unless you are using that System object in Simulink via the MATLAB System block. See “Generate Code for Global Data” (MATLAB Coder).

Methods

- Code generation support is available only for these System object methods:
 - `get`
 - `getNumInputs`
 - `getNumOutputs`
 - `isDone` (for sources only)
 - `isLocked`
 - `release`
 - `reset`
 - `set` (for tunable properties)
 - `step`
- For System objects that you define, code generation support is available only for these methods:
 - `getDiscreteStateImpl`
 - `getNumInputsImpl`
 - `getNumOutputsImpl`
 - `infoImpl`
 - `isDoneImpl`
 - `isInputDirectFeedthroughImpl`
 - `outputImpl`
 - `processTunedPropertiesImpl`
 - `releaseImpl` — Code is not generated automatically for this method. To release an object, you must explicitly call the `release` method in your code.
 - `resetImpl`
 - `setupImpl`

- `stepImpl`
- `updateImpl`
- `validateInputsImpl`
- `validatePropertiesImpl`

System Objects in codegen

You can include System objects in MATLAB code in the same way you include any other elements. You can then compile a MEX file from your MATLAB code by using the `codegen` command, which is available if you have a MATLAB Coder license. This compilation process, which involves a number of optimizations, is useful for accelerating simulations. See “Get Started with MATLAB Coder” (MATLAB Coder) and “MATLAB Classes” (MATLAB Coder) for more information.

Note Most, but not all, System objects support code generation. Refer to the particular object’s reference page for information.

System Objects in the MATLAB Function Block

Using the MATLAB Function block, you can include any System object and any MATLAB language function in a Simulink model. This model can then generate embeddable code. System objects provide higher-level algorithms for code generation than do most associated blocks. For more information, see “Implement MATLAB Functions in Simulink with MATLAB Function Blocks”.

System Objects in the MATLAB System Block

Using the MATLAB System block, you can include in a Simulink model individual System objects that you create with a class definition file. The model can then generate embeddable code. For more information, see “MATLAB System Block”.

System Objects and MATLAB Compiler Software

MATLAB Compiler™ software supports System objects for use inside MATLAB functions. The compiler product does not support System objects for use in MATLAB scripts.

See Also

More About

- “Generate Code That Uses Row-Major Array Layout” (MATLAB Coder)

Specify Objects as Inputs

When you accelerate code by using `fiaccel`, to specify the type of an input that is a value class object, you can provide an example object with the `-args` option.

- 1 Define the value class. For example, define a class `myRectangle`.

```
classdef myRectangle
    properties
        length;
        width;
    end
    methods
        function obj = myRectangle(l,w)
            if nargin > 0
                obj.length = l;
                obj.width = w;
            end
        end
        function area = calcarea(obj)
            area = obj.length * obj.width;
        end
    end
end
```

- 2 Define a function that takes an object of the value class as an input. For example:

```
function z = getarea(r)
    %#codegen
    z = calcarea(r);
end
```

- 3 Define an object of the class.

```
rect_obj = myRectangle(fi(4),fi(5))

rect_obj =





    myRectangle with properties:

        length: [1x1 embedded.fi]
        width: [1x1 embedded.fi]
```

- 4 Pass the example object to `fiaccel` by using the `-args` option.

```
fiaccel getarea -args {rect_obj} -report
```

In the report, you see that `r` has the same properties, `length` and `width`, as the example object `rect_obj`.

| SUMMARY | | ALL MESSAGES (0) | | CODE INSIGHTS (0) | | VARIABLES | |
|---------|--------|------------------|-------------|-------------------|--|-----------|---|
| Name | Type | Size | Class | | | | |
| z | Output | 1 × 1 | embedded.fi | | | |  |
| r | Input | 1 × 1 | myRectangle | | | |  |
| length | | 1 × 1 | embedded.fi | | | |  |
| width | | 1 × 1 | embedded.fi | | | |  |

Instead of providing an example object, you can create a type for an object of the value class and provide the type with the `-args` option.

- 1 Define an object of the class:

```
rect_obj = myRectangle(fi(4), fi(5))
```

```
rect_obj =
```

```
  myRectangle with properties:
```

```
    length: [1×1 embedded.fi]
    width: [1×1 embedded.fi]
```

- 2 To create a type for an object of `myRectangle` that has the same property types as `rect_obj`, use `coder.typeof`. `coder.typeof` creates a `coder.ClassType` object that defines a type for a class.

```
t = coder.typeof(rect_obj)
```

```
t =
```

```
coder.ClassType
  1×1 myRectangle
    length: 1×1 embedded.fi
           DataTypeMode: Fixed-point: binary point scaling
           Signedness: Signed
           WordLength: 16
           FractionLength: 12

    width : 1×1 embedded.fi
           DataTypeMode: Fixed-point: binary point scaling
           Signedness: Signed
           WordLength: 16
           FractionLength: 12
```

- 3 Pass the type to `fiaccel` by using the `-args` option.

```
fiaccel getarea -args {t} -report
```

After you create the type, you can change the types of the properties.

```
t.Properties.length = coder.typeof(fi(0,1,32,29))
t.Properties.width = coder.typeof(fi(0,1,32,29))
```

You can also add or delete properties. For example, to add a property `newprop`:

```
t.Properties.newprop = coder.typeof(int16(1))
```

Consistency Between `coder.ClassType` Object and Class Definition File

When you accelerate code, the properties of the `coder.ClassType` object that you pass to `fiaccel` must be consistent with the properties in the class definition file. If the class definition file has properties that your code does not use, the `coder.ClassType` object does not have to include those properties. `fiaccel` removes properties that you do not use.

Limitations for Using Objects as Entry-Point Function Inputs

Entry-point function inputs that are objects have these limitations:

- An object that is an entry-point function input must be an object of a value class. Objects of handle classes cannot be entry-point function inputs. Therefore, a value class that contains a handle class cannot be an entry-point function input.
- An object cannot be a global variable.
- If an object has duplicate property names, you cannot use it with `coder.Constant`. Duplicate property names occur in an object of a subclass in these situations:
 - The subclass has a property with the same name as a property of the superclass.
 - The subclass derives from multiple superclasses that use the same name for a property.

For information about when MATLAB allows duplicate property names, see “Subclassing Multiple Classes”.

See Also

`coder.typeof`

More About

- “Specify Properties of Entry-Point Function Inputs” on page 31-2
- “MATLAB Classes Definition for Code Generation” on page 15-2

Work Around Language Limitation: Code Generation Does Not Support Object Arrays

Issue

In certain situations, your MATLAB algorithm uses an array of objects that are instances of the same class. But code generation does not support object arrays. When attempting to generate code for such MATLAB code, you get this or a similar error message:

Code generation does not support object arrays.

Possible Solutions

Use Cell Array of Objects

Code generation supports cell arrays of objects. In your MATLAB code, represent the collection of objects by using a cell array instead of an array.

For example, suppose that your MATLAB algorithm uses the class `Square`:

```
classdef Square
    properties(Access = private)
        side
    end

    methods(Access = public)
        function obj = Square(side)
            obj.side = side;
        end

        function area = calculateArea(obj)
            area = obj.side^2;
        end
    end
end
```

The function `addAreas` constructs and uses a 1-by-3 array of `Square` objects:

```
function y = addAreas(n)
    obj = Square(0);
    collection = [obj obj obj]; % collection is an array

    for i = 1:numel(collection)
        collection(i) = Square(n + i);
    end

    y = 0;
    for i = 1:numel(collection)
        y = y + collection(i).calculateArea;
    end
end
```

Attempt to generate a MEX function for `addAreas`. Code generation fails because the local variable `collection` is an object array.

```
codegen addAreas -args 0 -report
```

??? Code generation does not support object arrays.

```
Error in ==> addAreas Line: 3 Column: 14  
Code generation failed: View Error Report
```

Redefine `collection` to be a cell array instead. Modify the code to use cell array indexing to index into `collection`. Name the modified function `addAreas_new`.

```
function y = addAreas_new(n)  
obj = Square(0);  
collection = {obj obj obj}; % collection is a cell array  
  
for i = 1:numel(collection)  
    collection{i} = Square(n + i);  
end  
  
y = 0;  
for i = 1:numel(collection)  
    y = y + collection{i}.calculateArea;  
end  
end
```

Attempt to generate a MEX function for `addAreas_new`. Code generation succeeds and produces `addAreas_new_mex`.

```
codegen addAreas_new -args 0 -report
```

Code generation successful: View report

Verify that `addAreas_new` and `addAreas_new_mex` have the same runtime behavior.

```
disp([addAreas_new(0) addAreas_new_mex(0)])
```

```
14    14
```

For Assignment with Nonscalar Indexing, Use Curly Braces and deal

Suppose that your original MATLAB code performs assignment to the array of objects by using nonscalar indexing. For example, you might add this line after the first `for` loop in the `addAreas` function:

```
collection(1:2) = [Square(10) Square(20)];
```

In the modified function `addAreas_new`, index into the corresponding cell array by using curly braces `{}` and perform assignment by using the `deal` function. Replace the above line by:

```
[collection{1:2}] = deal(Square(10),Square(20));
```

See Also

More About

- “MATLAB Classes Definition for Code Generation” (MATLAB Coder)
- “What Is a Cell Array?”

Defining Data for Code Generation

- “Data Definition Considerations for Code Generation” on page 16-2
- “Code Generation for Complex Data” on page 16-8
- “Encoding of Characters in Code Generation” on page 16-12
- “Array Size Restrictions for Code Generation” on page 16-13
- “Code Generation for Constants in Structures and Arrays” on page 16-14
- “Code Generation for Strings” on page 16-16
- “Define String Scalar Inputs” on page 16-17
- “Code Generation for Sparse Matrices” on page 16-19
- “Specify Array Layout in Functions and Classes” on page 16-21
- “Code Design for Row-Major Array Layout” on page 16-25
- “Generate Code for Growing Arrays and Cell Arrays with end + 1 Indexing” on page 16-27

Data Definition Considerations for Code Generation

To generate efficient standalone code, you must define the following types and classes of data differently from when you run your code in MATLAB.

| Data | Type Considerations | More Information |
|------------------------|---|---|
| Arrays | Maximum number of elements is restricted. | "Array Size Restrictions for Code Generation" on page 16-13 |
| Numeric types | Assign numeric type variables a value before using them in operations or returning them as outputs. | "Best Practices for Defining Variables for C/C++ Code Generation" (MATLAB Coder) |
| Complex numbers | <ul style="list-style-type: none"> Set complexity of variables at the time of assignment and before first use. Expressions containing a complex number or variable evaluate to a complex result, even if the imaginary part of the result is zero. | "Code Generation for Complex Data" on page 16-8 |
| Characters and strings | <ul style="list-style-type: none"> Characters are restricted to 8 bits of precision. For code generation, string scalars do not support global variables, indexing with curly braces, missing values, or size changes by using the function <code>coder. varsizes</code>. | <ul style="list-style-type: none"> "Encoding of Characters in Code Generation" on page 16-12 "Code Generation for Strings" (MATLAB Coder) |
| Variable-Size data | After initial fixed-size assignment to a variable, attempts to grow the variable might cause a compilation error. | <ul style="list-style-type: none"> "Code Generation for Variable-Size Arrays" (MATLAB Coder) "Define Variable-Size Data for Code Generation" (MATLAB Coder) |
| Structures | <ul style="list-style-type: none"> Assign fields to structures in the same order on each control path. Assign corresponding fields in the structure array elements with same size, type, and complexity. | <ul style="list-style-type: none"> "Define Scalar Structures for Code Generation" (MATLAB Coder) "Define Arrays of Structures for Code Generation" (MATLAB Coder) |

| Data | Type Considerations | More Information |
|--------------------|--|---|
| Cell arrays | <ul style="list-style-type: none"> • Assign all cell array elements before passing the cell array to a function or returning it from a function. • Variable-size cell array elements must all have the same size, type, and complexity. | <ul style="list-style-type: none"> • “Code Generation for Cell Arrays” (MATLAB Coder) • “Cell Array Limitations for Code Generation” (MATLAB Coder) |
| Tables | <ul style="list-style-type: none"> • You must specify variable names by using the 'VariableNames' name-value argument when creating tables from input arrays. • Limited data type support when you preallocate a table by using the table function and the 'Size' name-value argument. • Table indices that specify variables must be compile time constant. • You cannot change the size of a table by assignments. • You cannot change the VariableNames, RowNames, DimensionNames, or UserData properties of a table after you create it. <p>Limitations that apply to classes also apply to tables.</p> | <ul style="list-style-type: none"> • “Code Generation for Tables” (MATLAB Coder) • “Table Limitations for Code Generation” (MATLAB Coder) |
| Categorical arrays | <p>Categorical arrays do not support these inputs and operations:</p> <ul style="list-style-type: none"> • Arrays of MATLAB objects • Sparse matrices • Duplicate category names • Growth by assignment • Adding a category • Deleting an element <p>Limitations that apply to classes also apply to categorical arrays.</p> | <ul style="list-style-type: none"> • “Code Generation for Categorical Arrays” (MATLAB Coder) • “Categorical Array Limitations for Code Generation” (MATLAB Coder) |

| Data | Type Considerations | More Information |
|-----------------|--|---|
| Datetime arrays | <p><code>datetime</code> arrays do not support these inputs and operations:</p> <ul style="list-style-type: none"> • Text inputs • The 'Format' name-value argument • The 'TimeZone' name-value argument and the 'TimeZone' property • Setting time component properties • Growth by assignment • Deleting an element <p>Limitations that apply to classes also apply to <code>datetime</code> arrays.</p> | <ul style="list-style-type: none"> • “Code Generation for Datetime Arrays” (MATLAB Coder) • “Datetime Array Limitations for Code Generation” (MATLAB Coder) |
| Duration arrays | <p>Duration arrays do not support these inputs and operations:</p> <ul style="list-style-type: none"> • Text inputs • Growth by assignment • Deleting an element • Converting duration values to text by using <code>char</code>, <code>cellstr</code>, or <code>string</code> functions <p>Limitations that apply to classes also apply to duration arrays.</p> | <ul style="list-style-type: none"> • “Code Generation for Duration Arrays” (MATLAB Coder) • “Duration Array Limitations for Code Generation” (MATLAB Coder) |

| Data | Type Considerations | More Information |
|-----------------|--|---|
| Timetables | <ul style="list-style-type: none"> • You must specify variable names by using the 'VariableNames' name-value argument when creating timetables from input arrays. • Limited data type support when you preallocate a table by using the timetable function and the 'Size' name-value argument. • Timetable indices that specify variables must be compile time constant. • You cannot change the size of a timetable by assignments. • You cannot change the VariableNames, DimensionNames, or UserData properties of a timetable after you create it. • If you create a regular timetable, and you attempt to set irregular row times, then an error is produced. • If you create an irregular timetable, then it remains irregular even if you set its sample rate or time step. <p>Limitations that apply to classes also apply to timetables.</p> | <ul style="list-style-type: none"> • “Code Generation for Timetables” (MATLAB Coder) • “Timetable Limitations for Code Generation” (MATLAB Coder) |
| Enumerated data | Supports integer-based enumerated types only. | “Enumerations” |

| Data | Type Considerations | More Information |
|------------------|--|--|
| MATLAB Classes | <ul style="list-style-type: none"> • Before generating code, it is a best practice to test class property validation by running a MEX function over the full range of input values. • If a property does not have an explicit initial value, the code generator assumes that it is undefined at the beginning of the constructor. The code generator does not assign an empty matrix as the default. • The <code>coder. varsize</code> function is not supported for class properties. • If the initial value of a property is an object, then the property must be constant. To make a property constant, declare the Constant attribute in the property block. | <ul style="list-style-type: none"> • “Generate C++ Classes for MATLAB Classes” (MATLAB Coder) • “MATLAB Classes Definition for Code Generation” (MATLAB Coder) |
| Function handles | <ul style="list-style-type: none"> • Assigning different function handles to the same variable can cause a compile-time error. • You cannot pass function handles to or from entry-point functions or extrinsic functions. • You cannot view function handles from the MATLAB Function Block debugger. | “Function Handles” |

| Data | Type Considerations | More Information |
|----------------------|---|--|
| Deep learning arrays | <p><code>dlarrays</code> do not support these inputs and operations:</p> <ul style="list-style-type: none"> • The data format argument must be a compile-time constant • Define <code>dlarray</code> variables inside the entry-point function. • The input to a <code>dlarray</code> must be fixed-size. • Code generation does not support creating a <code>dlarray</code> type object by using the <code>coder.typeof</code> function with upper bound size and variable dimensions specified. | <ul style="list-style-type: none"> • “Code Generation for <code>dlarray</code>” (MATLAB Coder) • “<code>dlarray</code> Limitations for Code Generation” (MATLAB Coder) |

The information in the preceding table is not an exhaustive list of considerations for each data type. See the topics in the More Information column.

See Also

Related Examples

- “Best Practices for Defining Variables for C/C++ Code Generation” (MATLAB Coder)
- “Reuse the Same Variable with Different Properties” (MATLAB Coder)
- “Eliminate Redundant Copies of Variables in Generated Code” (MATLAB Coder)

Code Generation for Complex Data

In this section...

“Restrictions When Defining Complex Variables” on page 16-8

“Code Generation for Complex Data with Zero-Valued Imaginary Parts” on page 16-8

“Results of Expressions That Have Complex Operands” on page 16-11

“Results of Complex Multiplication with Nonfinite Values” on page 16-11

Restrictions When Defining Complex Variables

For code generation, you must set the complexity of variables at the time of assignment. Assign a complex constant to the variable or use the `complex` function. For example:

```
x = 5 + 6i; % x is a complex number by assignment.
y = complex(5,6); % y is the complex number 5 + 6i.
```

After assignment, you cannot change the complexity of a variable. Code generation for the following function fails because $x(k) = 3 + 4i$ changes the complexity of x .

```
function x = test1( )
x = zeros(3,3); % x is real
for k = 1:numel(x)
    x(k) = 3 + 4i;
end
end
```

To resolve this issue, assign a complex constant to x .

```
function x = test1( )
x = zeros(3,3)+ 0i; %x is complex
for k = 1:numel(x)
    x(k) = 3 + 4i;
end
end
```

Code Generation for Complex Data with Zero-Valued Imaginary Parts

For code generation, complex data that has all zero-valued imaginary parts remains complex. This data does not become real. This behavior has the following implications:

- In some cases, results from functions that sort complex data by absolute value can differ from the MATLAB results. See “Functions That Sort Complex Values by Absolute Value” on page 16-8.
- For functions that require that complex inputs are sorted by absolute value, complex inputs with zero-valued imaginary parts must be sorted by absolute value. These functions include `ismember`, `union`, `intersect`, `setdiff`, and `setxor`.

Functions That Sort Complex Values by Absolute Value

Functions that sort complex values by absolute value include `sort`, `issorted`, `sortrows`, `median`, `min`, and `max`. These functions sort complex numbers by absolute value even when the imaginary parts are zero. In general, sorting the absolute values produces a different result than sorting the real parts. Therefore, when inputs to these functions are complex with zero-valued imaginary parts in

generated code, but real in MATLAB, the generated code can produce different results than MATLAB. In the following examples, the input to `sort` is real in MATLAB, but complex with zero-valued imaginary parts in the generated code:

- **You Pass Real Inputs to a Function Generated for Complex Inputs**

- 1 Write this function:

```
function myout = mysort(A)
myout = sort(A);
end
```

- 2 Call `mysort` in MATLAB.

```
A = -2:2;
mysort(A)

ans =
```

```
    -2    -1     0     1     2
```

- 3 Generate a MEX function for complex inputs.

```
A = -2:2;
codegen mysort -args {complex(A)} -report
```

- 4 Call the MEX Function with real inputs.

```
mysort_mex(A)

ans =
```

```
     0     1    -1     2    -2
```

You generated the MEX function for complex inputs, therefore, it treats the real inputs as complex numbers with zero-valued imaginary parts. It sorts the numbers by the absolute values of the complex numbers. Because the imaginary parts are zero, the MEX function returns the results to the MATLAB workspace as real numbers. See “Inputs and Outputs for MEX Functions Generated for Complex Arguments” on page 16-10.

- **Input to `sort` Is Output from a Function That Returns Complex in Generated Code**

- 1 Write this function:

```
function y = myfun(A)
x = eig(A);
y = sort(x, 'descend');
```

The output from `eig` is the input to `sort`. In generated code, `eig` returns a complex result. Therefore, in the generated code, `x` is complex.

- 2 Call `myfun` in MATLAB.

```
A = [2 3 5;0 5 5;6 7 4];
myfun(A)
```

```
ans =
```

```
    12.5777
     2.0000
    -3.5777
```

The result of `eig` is real. Therefore, the inputs to `sort` are real.

- 3 Generate a MEX function for complex inputs.

```
codegen myfun -args {complex(A)}
```

- 4 Call the MEX function.

```
myfun_mex(A)
```

```
ans =
```

```
12.5777
-3.5777
 2.0000
```

In the MEX function, `eig` returns a complex result. Therefore, the inputs to `sort` are complex. The MEX function sorts the inputs in descending order of the absolute values.

Inputs and Outputs for MEX Functions Generated for Complex Arguments

For MEX functions created by the `codegen` command, the `fiaccl` command, or the MATLAB Coder app:

- Suppose that you generate the MEX function for complex inputs. If you call the MEX function with real inputs, the MEX function transforms the real inputs to complex values with zero-valued imaginary parts.
- If the MEX function returns complex values that have all zero-valued imaginary parts, the MEX function returns the values to the MATLAB workspace as real values. For example, consider this function:

```
function y = foo()
    y = 1 + 0i; % y is complex with imaginary part equal to zero
end
```

If you generate a MEX function for `foo` and view the code generation report, you see that `y` is complex.

```
codegen foo -report
```

| Name | Type | Size | Class |
|------|--------|-------|----------------|
| y | Output | 1 × 1 | complex double |

If you run the MEX function, you see that in the MATLAB workspace, the result of `foo_mex` is the real value 1.

```
z = foo_mex
```

```
ans =
```

```
1
```

Results of Expressions That Have Complex Operands

In general, expressions that contain one or more complex operands produce a complex result in generated code, even if the value of the result is zero. Consider the following line of code:

```
z = x + y;
```

Suppose that at run time, x has the value $2 + 3i$ and y has the value $2 - 3i$. In MATLAB, this code produces the real result $z = 4$. During code generation, the types for x and y are known, but their values are not known. Because either or both operands in this expression are complex, z is defined as a complex variable requiring storage for a real and an imaginary part. z equals the complex result $4 + 0i$ in generated code, not 4 , as in MATLAB code.

Exceptions to this behavior are:

- When the imaginary parts of complex results are zero, MEX functions return the results to the MATLAB workspace as real values. See “Inputs and Outputs for MEX Functions Generated for Complex Arguments” on page 16-10.
- When the imaginary part of the argument is zero, complex arguments to extrinsic functions are real.

```
function y = foo()
    coder.extrinsic('sqrt')
    x = 1 + 0i; % x is complex
    y = sqrt(x); % x is real, y is real
end
```

- Functions that take complex arguments but produce real results return real values.

```
y = real(x); % y is the real part of the complex number x.
y = imag(x); % y is the real-valued imaginary part of x.
y = isreal(x); % y is false (0) for a complex number x.
```

- Functions that take real arguments but produce complex results return complex values.

```
z = complex(x,y); % z is a complex number for a real x and y.
```

Results of Complex Multiplication with Nonfinite Values

When an operand of a complex multiplication contains a nonfinite value, the generated code might produce a different result than the result that MATLAB produces. The difference is due to the way that code generation defines complex multiplication. For code generation:

- Multiplication of a complex value by a complex value $(a + bi)(c + di)$ is defined as $(ac - bd) + (ad + bc)i$. The complete calculation is performed, even when a real or an imaginary part is zero.
- Multiplication of a real value by a complex value $c(a + bi)$ is defined as $ca + cbi$.

Encoding of Characters in Code Generation

MATLAB represents characters in 16-bit Unicode. The code generator represents characters in an 8-bit codeset that the locale setting determines. Differences in character encoding between MATLAB and code generation have these consequences:

- Code generation of characters with numeric values greater than 255 produces an error.
- For some characters in the range 128–255, it might not be possible to represent the character in the codeset of the locale setting or to convert the character to an equivalent 16-bit Unicode character. Passing characters in this range between MATLAB and generated code can result in errors or different answers.
- For code generation, some toolbox functions accept only 7-bit ASCII characters.
- Casting a character that is not in the 7-bit ASCII codeset to a numeric type, such as double, can produce a different result in the generated code than in MATLAB. As a best practice, for code generation, avoid performing arithmetic with characters.

See Also

More About

- “Locale Setting Concepts for Internationalization”
- “Differences Between Generated Code and MATLAB Code” on page 19-6

Array Size Restrictions for Code Generation

For code generation, the maximum number of elements of an array is constrained by the code generator and the target hardware.

For fixed-size arrays and variable-size arrays that use static memory allocation, the maximum number of elements is the smaller of:

- `intmax('int32')`.
- The largest integer that fits in the C `int` data type on the target hardware.

For variable-size arrays that use dynamic memory allocation, the maximum number of elements is the smaller of:

- `intmax('int32')`.
- The largest power of 2 that fits in the C `int` data type on the target hardware.

These restrictions apply even on a 64-bit platform.

For a fixed-size array, if the number of elements exceeds the maximum, the code generator reports an error at compile time.

For a variable-size array, if the number of elements exceeds the maximum at run time, the error checking behavior depends on the code generation target:

- While running the code generated by using the `codegen` command, the `fiaccl` command, or the MATLAB Coder app, if run-time error checks are enabled, the generated code reports an error. By default, run-time error checks are enabled for MEX code and disabled for standalone C/C++ code.
- During simulation of a MATLAB Function block, the software reports an error. Generated standalone code for MATLAB Function blocks cannot report array size violations.

See Also

More About

- “Control Run-Time Checks” on page 12-48
- “Potential Differences Reporting” on page 19-23

Code Generation for Constants in Structures and Arrays

The code generator does not recognize constant structure fields or array elements in the following cases:

Fields or elements are assigned inside control constructs

In the following code, the code generator recognizes that the structure fields `s.a` and `s.b` are constants.

```
function y = mystruct()
s.a = 3;
s.b = 5;
y = zeros(s.a,s.b);
```

If any structure field is assigned inside a control construct, the code generator does not recognize the constant fields. This limitation also applies to arrays with constant elements. Consider the following code:

```
function y = mystruct(x)
s.a = 3;
if x > 1
    s.b = 4;
else
    s.b = 5;
end
y = zeros(s.a,s.b);
```

The code generator does not recognize that `s.a` and `s.b` are constant. If variable-sizing is enabled, `y` is treated as a variable-size array. If variable-sizing is disabled, the code generator reports an error.

Constants are assigned to array elements using non-scalar indexing

In the following code, the code generator recognizes that `a(1)` is constant.

```
function y = myarray()
a = zeros(1,3);
a(1) = 20;
y = coder.const(a(1));
```

In the following code, because `a(1)` is assigned using non-scalar indexing, the code generator does not recognize that `a(1)` is constant.

```
function y = myarray()
a = zeros(1,3);
a(1:2) = 20;
y = coder.const(a(1));
```

A function returns a structure or array that has constant and nonconstant elements

For an output structure that has both constant and nonconstant fields, the code generator does not recognize the constant fields. This limitation also applies to arrays that have constant and nonconstant elements. Consider the following code:

```
function y = mystruct_out(x)
s = create_structure(x);
y = coder.const(s.a);
```

```
function s = create_structure(x)
s.a = 10;
s.b = x;
```

Because `create_structure` returns a structure `s` that has one constant field and one nonconstant field, the code generator does not recognize that `s.a` is constant. The `coder.const` call fails because `s.a` is not constant.

Code Generation for Strings

Code generation supports 1-by-1 MATLAB string arrays. Code generation does not support string arrays that have more than one element.

A 1-by-1 string array, called a string scalar, contains one piece of text, represented as a 1-by-n character vector. An example of a string scalar is "Hello, world". For more information about strings, see "Text in String and Character Arrays".

Limitations

For string scalars, code generation does not support:

- Global variables
- Indexing with curly braces {}
- Missing values
- Defining input types programmatically by using preconditioning with `assert` statements, when generating code by using the `codegen` command, the `fiaccl` command, or the MATLAB Coder app
- Their use with `coder.varsize`, when generating code by using the `codegen` command, the `fiaccl` command, or the MATLAB Coder app
- Their use as Simulink signals, parameters, or data store memory

For code generation, limitations that apply to classes apply to strings. See "MATLAB Classes Definition for Code Generation" on page 15-2.

Differences Between Generated Code and MATLAB Code

- Converting a string that contains multiple unary operators to `double` can produce different results between MATLAB and the generated code. Consider this function:

```
function out = foo(op)
out = double(op + 1);
end
```

For an input value "--", the function converts the string "--1" to `double`. In MATLAB, the answer is NaN. In the generated code, the answer is 1.

- Double conversion for a string with misplaced commas (commas that are not used as thousands separators) can produce different results from MATLAB.

See Also

More About

- "Define String Scalar Inputs" on page 16-17

Define String Scalar Inputs

You can define string scalar inputs at the command line. Programmatic specification of string scalar input types by using preconditioning (`assert` statements) is not supported.

Define String Scalar Types at the Command Line

To define string scalar inputs at the command line, use one of these procedures:

- “Provide an Example String Scalar Input” on page 16-17
- “Provide a String Scalar Type” on page 16-17
- “Provide a Constant String Scalar Input” on page 16-17
- “Provide a Variable-Size String Scalar Input” on page 16-17

Alternatively, if you have a test file that calls your entry-point function with example inputs, you can determine the input types by using `coder.getArgTypes`.

Provide an Example String Scalar Input

To provide an example string scalar to `fiaccl`, use the `-args` option:

```
fiaccl myFunction -args {"Hello, world"}
```

Provide a String Scalar Type

To provide a type for a string scalar to `fiaccl`:

- 1 Define a string scalar. For example:


```
s = "mystring";
```
- 2 Create a type from `s`.


```
t = coder.typeof(s);
```
- 3 Pass the type to `fiaccl` by using the `-args` option.

```
fiaccl myFunction -args {t}
```

Provide a Constant String Scalar Input

To specify that a string scalar input is constant, use `coder.Constant` with the `-args` option:

```
fiaccl myFunction -args {coder.Constant("Hello, world")}
```

Provide a Variable-Size String Scalar Input

To specify that a string scalar input has a variable-size:

- 1 Define a string scalar. For example:


```
s = "mystring";
```
- 2 Create a type from `s`.


```
t = coder.typeof(s);
```

- 3 Assign the `StringLength` property of the type the upper bound of the string length and set `VariableStringLength` to `true`. For example, specify that type `t` is variable-size with an upper bound of 10.

```
t.StringLength = 10;  
t.VariableStringLength = true;
```

To specify that `t` is variable-size with no upper bound:

```
t.StringLength = Inf;
```

This automatically sets the `VariableStringLength` property to `true`.

- 4 Pass the type to `fiaccel` by using the `-args` option.

```
fiaccel myFunction -args {t}
```

See Also

`coder.Constant` | `coder.typeof`

More About

- “Code Generation for Strings” on page 16-16
- “Specify Properties of Entry-Point Function Inputs” on page 31-2

Code Generation for Sparse Matrices

Sparse matrices provide efficient storage in memory for arrays with many zero elements. Sparse matrices can provide improved performance and reduced memory usage for generated code. Computation time on sparse matrices scales only with the number of operations on nonzero elements.

Functions for creating and manipulating sparse matrices are listed in “Sparse Matrices”. To check if a function is supported for code generation, see the function reference page. Code generation does not support sparse matrix inputs created by using `sparse` for all functions.

Input Definition

You can use `coder.typeof` to initialize a sparse matrix input to your function. For sparse matrices, the code generator does not track upper bounds for variable-size dimensions. All variable-size dimensions are treated as unbounded.

You cannot define sparse input types programmatically by using `assert` statements.

Code Generation Guidelines

Initialize matrices by using sparse constructors to maximize your code efficiency. For example, to construct a 3-by-3 identity matrix, use `speye(3,3)` rather than `sparse(eye(3,3))`.

Indexed assignment into sparse matrices incurs an overhead compared to indexed assignment into full matrices. For example:

```
S = speye(10);
S(7,7) = 42;
```

As in MATLAB, sparse matrices are stored in compressed sparse column format. When you insert a new nonzero element into a sparse matrix, all subsequent nonzero elements must be shifted downward, column by column. These extra manipulations can slow performance. See “Accessing Sparse Matrices”.

Code Generation Limitations

To generate code that uses sparse matrices, dynamic memory allocation must be enabled. To store the changing number of nonzero elements, and their values, sparse matrices use variable-size arrays in the generated code. To change dynamic memory allocation settings, see “Control Memory Allocation for Variable-Size Arrays” on page 29-4. Because sparse matrices use variable-size arrays for dynamic memory allocation, limitations on “Variable-Size Data” also apply to sparse matrices.

You cannot assign sparse data to data that is not sparse. The generated code uses distinct data type representations for sparse and full matrices. To convert to and from sparse data, use the explicit `sparse` and `full` conversion functions.

You cannot define a sparse matrix with competing size specifications. The code generator fixes the size of the sparse matrix when it produces the corresponding data type definition in C/C++. As an example, the function `foo` causes an error in code generation:

```
function y = foo(n)
  %#codegen
```

```
if n > 0
    y = sparse(3,2);
else
    y = sparse(4,3);
end
```

Logical indexing into sparse matrices is not supported for code generation. For example, this syntax causes an error:

```
S = magic(3);
S(S > 7) = 42;
```

For sparse matrices, you cannot delete array elements by assigning empty arrays:

```
S(:,2) = [];
```

See Also

[sparse](#) | [full](#) | [coder.typeof](#) | [magic](#) | [speye](#)

More About

- “Sparse Matrices”
- “Use Dynamically Allocated C++ Arrays in Generated Function Interfaces” (MATLAB Coder)

Specify Array Layout in Functions and Classes

You can specialize individual MATLAB functions for row-major layout or column-major layout by inserting `coder.rowMajor` or `coder.columnMajor` calls into the function body. Using these function specializations, you can combine row-major data and column-major data in your generated code. You can also specialize classes for one specific array layout. Function and class specializations allow you to:

- Incrementally modify your code for row-major layout or column-major layout.
- Define array layout boundaries for applications that require different layouts in different components.
- Structure the inheritance of array layout between many different functions and classes.

For MATLAB Coder entry-point (top-level) functions, all inputs and outputs must use the same array layout. In the generated C/C++ code, the entry-point function interface accepts and returns data with the same array layout as the function array layout specification.

Note By default, code generation uses column-major array layout.

Specify Array Layout in a Function

For an example of a specialized function, consider `addMatrixRM`:

```
function [S] = addMatrixRM(A,B)
%#codegen
S = zeros(size(A));
coder.rowMajor; % specify row-major code
for row = 1:size(A,1)
    for col = 1:size(A,2)
        S(row,col) = A(row,col) + B(row,col);
    end
end
```

For MATLAB Coder, you can generate code for `addMatrixRM` by using the `codegen` command.

```
codegen addMatrixRM -args {ones(20,10),ones(20,10)} -config:lib -launchreport
```

Because of the `coder.rowMajor` call, the code generator produces code that uses data stored in row-major layout.

Other functions called from a row-major function or column-major function inherit the same array layout. If a called function has its own distinct `coder.rowMajor` or `coder.columnMajor` call, the local call takes precedence.

You can mix column-major and row-major functions in the same code. The code generator inserts transpose or conversion operations when passing data between row-major and column-major functions. These conversion operations ensure that array elements are stored as required by functions with different array layout specifications. For example, the inputs to a column-major function, called from a row-major function, are converted to column-major layout before being passed to the column-major function.

Query Array Layout of a Function

To query the array layout of a function at compile time, use `coder.isRowMajor` or `coder.isColumnMajor`. This query can be useful for specializing your generated code when it involves row-major and column-major functions. For example, consider this function:

```
function [S] = addMatrixRouted(A,B)
if coder.isRowMajor
    %execute this code if row-major
    S = addMatrixRM(A,B);
elseif coder.isColumnMajor
    %execute this code if column-major
    S = addMatrix_OptimizedForColumnMajor(A,B);
end
```

This function behaves differently depending on whether it is row-major or column-major. When `addMatrixRouted` is row-major, it calls the `addMatrixRM` function, which has efficient memory access for row-major data. When the function is column-major, it calls a version of the `addMatrixRM` function optimized for column-major data.

For example, consider this function definition. The algorithm iterates through the columns in the outer loop and the rows in the inner loop, in contrast to the `addMatrixRM` function.

```
function [S] = addMatrix_OptimizedForColumnMajor(A,B)
%#codegen
S = zeros(size(A));
for col = 1:size(A,2)
    for row = 1:size(A,1)
        S(row,col) = A(row,col) + B(row,col);
    end
end
```

Code generation for this function yields:

```
...
/* column-major layout */
for (col = 0; col < 10; col++) {
    for (row = 0; row < 20; row++) {
        S[row + 20 * col] = A[row + 20 * col] + B[row + 20 * col];
    }
}
...
```

The generated code has a stride length of only one element. Due to the specializing queries, the generated code for `addMatrixRouted` provides efficient memory access for either choice of array layout.

Specify Array Layout in a Class

You can specify array layout for a class so that object property variables are stored with a specific array layout. To specify the array layout, place a `coder.rowMajor` or `coder.columnMajor` call in the class constructor. If you assign an object with a specified array layout to the property of another object, the array layout of the assigned object takes precedence.

Consider the row-major class `rowMats` as an example. This class contains matrix properties and a method that consists of an element-wise addition algorithm. The algorithm in the method performs

more efficiently for data stored in row-major layout. By specifying `coder.rowMajor` in the class constructor, the generated code uses row-major layout for the property data.

```
classdef rowMats
    properties (Access = public)
        A;
        B;
        C;
    end
    methods
        function obj = rowMats(A,B)
            coder.rowMajor;
            if nargin == 0
                obj.A = 0;
                obj.B = 0;
                obj.C = 0;
            else
                obj.A = A;
                obj.B = B;
                obj.C = zeros(size(A));
            end
        end
        function obj = add(obj)
            for row = 1:size(obj.A,1)
                for col = 1:size(obj.A,2)
                    obj.C(row,col) = obj.A(row,col) + obj.B(row,col);
                end
            end
        end
    end
end
```

Use the class in a simple function `doMath`. The inputs and outputs of the entry-point function must all use the same array layout.

```
function [out] = doMath(in1,in2)
    %#codegen
    out = zeros(size(in1));
    myMats = rowMats(in1,in2);
    myMats = myMats.add;
    out = myMats.C;
end
```

For MATLAB Coder, you can generate code by entering:

```
A = rand(20,10);
B = rand(20,10);
cfg = coder.config('lib');
codegen -config cfg doMath -args {A,B} -launchreport
```

With default settings, the code generator assumes that the entry-point function inputs and outputs use column-major layout, because you do not specify row-major layout for the function `doMath`. Therefore, before calling the class constructor, the generated code converts `in1` and `in2` to row-major layout. Similarly, it converts the `doMath` function output back to column-major layout.

When designing a class for a specific array layout, consider:

- If you do not specify the array layout in a class constructor, objects inherit their array layout from the function that calls the class constructor, or from code generation configuration settings.
- You cannot specify the array layout in a nonstatic method by using `coder.rowMajor` or `coder.columnMajor`. Methods use the same array layout as the receiving object. Methods do not inherit the array layout of the function that calls them. For static methods, which are used similarly to ordinary functions, you can specify the array layout in the method.
- If you specify the array layout of a superclass, the subclass inherits this array layout specification. You cannot specify conflicting array layouts between superclasses and subclasses.

Code Design for Row-Major Array Layout

Outside of code generation, MATLAB uses column-major layout by default. Array layout specifications do not affect self-contained MATLAB code. To test the efficiency of your generated code or your MATLAB Function block, create separate versions with row-major layout and column-major layout. Then, compare their performance.

You can design your MATLAB code to avoid potential inefficiencies related to array layout. Inefficiencies can be caused by:

- Conversions between row-major layout and column-major layout.
- One-dimensional or linear indexing of row-major data.
- Reshaping or rearrangement of row-major data.

Array layout conversions are necessary when you mix row-major and column-major specifications in the same code or model, or when you use linear indexing on data that is stored in row-major. When you simulate a model or generate code for a model that uses column-major, and that contains a MATLAB Function block that uses row-major, then the software converts input data to row-major and output data back to column-major as needed, and vice versa.

Inefficiencies can be caused by functions or algorithms that are less optimized for a given choice of array layout. If a function or algorithm is more efficient for a different layout, you can enforce that layout by embedding it in another function with a `coder.rowMajor` or `coder.columnMajor` call.

Linear Indexing Uses Column-Major Array Layout

The code generator follows MATLAB column-major semantics for linear indexing. For more information on linear indexing in MATLAB, see “Array Indexing”.

To use linear indexing on row-major data, the code generator must first recalculate the data representation in column-major layout. This additional processing can slow performance. To improve code efficiency, avoid using linear indexing on row-major data, or use column-major layout for code that uses linear indexing.

For example, consider the function `sumShiftedProducts`, which accepts a matrix as an input and outputs a scalar value. The function uses linear indexing on the input matrix to sum up the product of each matrix element with an adjacent element. The output value of this operation depends on the order in which the input elements are stored.

```
function mySum = sumShiftedProducts(A)
%#codegen
mySum = 0;
% create linear vector of A elements
B = A(:);
% multiply B by B with elements shifted by one, and take sum
mySum = sum( B.*circshift(B,1) );
end
```

For MATLAB Coder, to generate code that uses row-major layout, enter:

```
codegen -config:mex sumShiftedProducts -args {ones(2,3)} -launchreport -rowmajor
```

For an example input, consider the matrix:

```
D = reshape(1:6,3,2)'
```

which yields:

```
D =  
    1    2    3  
    4    5    6
```

If you pass this matrix as input to the generated code, the elements of A are stored in the order:

```
    1    2    3    4    5    6
```

In contrast, because the vector B is obtained by linear indexing, it is stored in the order:

```
    1    4    2    5    3    6
```

The code generator must insert a reshaping operation to rearrange the data from row-major layout for A to column-major layout for B. This additional operation reduces the efficiency of the function for row-major layout. The inefficiency increases with the size of the array. Because linear indexing always uses column-major layout, the generated code for `sumShiftedProducts` produces the same output result whether generated with row-major layout or column-major layout.

In general, functions that compute indices or subscripts also use linear indexing, and produce results corresponding to data stored in column-major layout. These functions include:

- `ind2sub`
- `sub2ind`
- `colon`

Generate Code for Growing Arrays and Cell Arrays with end + 1 Indexing

Code generation supports growing either an array or a cell array in your MATLAB code by using end + 1 indexing. To use this functionality, make sure that the code generation configuration property EnableVariableSizing or the corresponding setting **Enable variable-sizing** in the MATLAB Coder app is enabled.

Grow Array with (end + 1) Indexing

To grow an array X, you can assign a value to X(end + 1). If you make this assignment in your MATLAB code, the code generator treats the dimension you grow as variable-size.

For example, you can generate code for this code snippet:

```
...
a = [1 2 3 4 5 6];
a(end + 1) = 7;

b = [1 2];
for i = 3:10
    b(end + 1) = i;
end
...
```

When you use (end + 1) to grow an array, follow these restrictions:

- Use only (end + 1). Do not use (end + 2), (end + 3), and so on.
- Use (end + 1) with vectors only. For example, the following code is not allowed because X is a matrix, not a vector.

```
...
X = [1 2; 3 4];
X(end + 1) = 5;
...
```

- You can grow empty arrays of size 1x0 by using (end + 1). Growing arrays of size 0x1 is not supported. Growing an array of size 0x0 is supported only if you create that array by using [].

Growing Variable-Size Column Array That is Initialized as Scalar at Run Time

In MATLAB execution, if you grow a scalar array by using (end+1) indexing, the array grows along the second dimension and produces a row vector. For example, define the function grow:

```
function z = grow(n, m)
n(end+1) = m;
z = n;
end
```

Call grow with example inputs:

```
grow(2,3)
```

```
ans =
     2     3
```

By contrast, in code generation, suppose that:

- You specify the array to be of variable-size column type (for example, `:Inf x 1`) at compile time, *and*
- Initialize this array as a scalar at run time.

In such situations, the generated code attempts to grow the scalar along the first dimension and therefore, produces a run-time error. For example, generate MEX code for `grow`. Specify the input `n` to be a `:Inf x 1` double array. Specify the input `m` to be a double scalar.

```
codegen grow -args {coder.typeof(0, [Inf 1], [1 0]), 1}
```

Code generation successful.

Run the generated MEX with the same inputs as before.

```
grow_mex(2,3)
```

```
Attempted to grow a scalar along the first dimension using end+1 indexing. This
behavior differs from MATLAB execution which grows a scalar along the second
dimension.
```

```
Error in grow (line 2)
n(end+1) = m;
```

How to Avoid This Error

To avoid this error and grow the array along the column dimension in both generated code and MATLAB execution, rewrite your MATLAB code in *either* of these ways:

- Grow the array using the concatenation operator instead of using `(end+1)`. For example, rewrite the `grow` function as:

```
function z = growCat(n, m)
n = [n;m];
z = n;
end
```

- In your function, create a temporary variable by transposing the variable that you want to grow. Then, grow this temporary variable by using `(end+1)` indexing. Finally, take a second transpose of this temporary variable. For example, rewrite the `grow` function as:

```
function z = growTransposed(n, m)
temp = n';
temp(end+1) = m;
z = temp';
end
```

Grow Cell Array with `{end + 1}` Indexing

To grow a cell array `X`, you can use `X{end + 1}`. For example:

```
...
X = {1 2};
```



```
X{end + 1} = 'a';
...
```

When you use {end + 1} to grow a cell array, follow these restrictions:

- Use only {end + 1}. Do not use {end + 2}, {end + 3}, and so on.
- Use {end + 1} with vectors only. For example, the following code is not allowed because X is a matrix, not a vector:

```
...
X = {1 2; 3 4};
X{end + 1} = 5;
```

- Use {end + 1} only with a variable. In the following code, {end + 1} does not cause {1 2 3} to grow. In this case, the code generator treats {end + 1} as an out-of-bounds index into X{2}.

```
...
X = {'a' { 1 2 3 }};
X{2}{end + 1} = 4;
...
```

- When {end + 1} grows a cell array in a loop, the cell array must be variable-size. Therefore, the cell array must be homogeneous on page 30-2.

This code is allowed because X is homogeneous.

```
...
X = {1 2};
for i=1:n
    X{end + 1} = 3;
end
...
```

This code is not allowed because X is heterogeneous.

```
...
X = {1 'a' 2 'b'};
for i=1:n
    X{end + 1} = 3;
end
...
```

- For a coding pattern that causes a difference in behavior between generated code and MATLAB, see “Growing Variable-Size Column Cell Array That is Initialized as Scalar at Run Time” on page 19-17.

Defining Functions for Code Generation

- “Code Generation for Variable Length Argument Lists” on page 17-2
- “Generate Code for arguments Block That Validates Input and Output Arguments” on page 17-3
- “Specify Number of Entry-Point Function Input or Output Arguments to Generate” on page 17-7
- “Code Generation for Anonymous Functions” on page 17-9
- “Code Generation for Nested Functions” on page 17-10

Code Generation for Variable Length Argument Lists

When you use `varargin` and `varargout` for code generation, there are these restrictions:

- If you use `varargin` to define an argument to an entry-point function, the code generator produces the function with a fixed number of arguments. This fixed number of arguments is based on the number of arguments that you specify when you generate code.
- You cannot write to `varargin`. If you want to write to input arguments, copy the values into a local variable.
- To index into `varargin` and `varargout`, use curly braces `{}`, not parentheses `()`.
- The code generator must be able to determine the value of the index into `varargin` or `varargout`.

See Also

More About

- “Nonconstant Index into `varargin` or `varargout` in a for-Loop” on page 49-43
- “Specify Number of Entry-Point Function Input or Output Arguments to Generate” on page 17-7

Generate Code for arguments Block That Validates Input and Output Arguments

You can generate code for `arguments` blocks that perform input and output argument validation in your MATLAB function. Using argument validation, you can constrain the class, size, and other aspects of function input and output values without writing code in the body of the function to perform these tests. See “Function Argument Validation”.

```
function myFunction(inputArg)
    arguments
        inputArg (dim1,dim2,...) ClassName {fcn1,fcn2,...} = defaultValue
    end
    % Function code
end
```

The diagram illustrates the structure of an `arguments` block. The input argument `inputArg` is followed by three validation specifications: `(dim1,dim2,...)` (labeled **Size**), `ClassName` (labeled **Class**), and `{fcn1,fcn2,...}` (labeled **Functions**). These specifications are followed by an equals sign and a default value `= defaultValue`.

Supported Features

Code generation supports most features of `arguments` blocks, including size and class validation, validation functions, and default values. Code generation also supports the `namedargs2cell` function.

Code generation does not support these features of `arguments` blocks:

- Size validation, class validation, and validation functions for repeating arguments
- Multiple repeating input arguments
- Name-value input arguments at entry-point functions
- Name-value input arguments from class properties using the `structName.?ClassName` syntax
- Size validation for `table`, `timetable`, or `dlarray` objects.

Names Must Be Compile-Time Constants

Suppose that `foo` is a function that uses name-value argument validation. When you call `foo` from another function `bar`, the code generator must be able to determine the names that you provide to `foo` at compile time.

For example, code generation succeeds for the entry-point function `myNamedArg_valid`. This function contains two calls to the function `local`. For both these calls, the argument name `'x'` is known during code generation.

```
function [out1,out2] = myNamedArg_valid(in1,in2)
    out1 = local(x=in1);
    out2 = local('x',in2);
end
```

```
function out = local(args)
arguments
    args.x
end

out = args.x;
end

codegen myNamedArg_valid -args {0,0}
```

Code generation successful.

By contrast, code generation fails for the entry-point function `myNamedArg_invalid` because the argument name for the function `local` is supplied at run time.

```
function out = myNamedArg_invalid(value, inputName)
out = local(inputName, value);
end

function out = local(args)
arguments
    args.x
end

out = args.x;
end

codegen myNamedArg_invalid -args {0,coder.typeof('a')}
```

Error calling 'myNamedArg_invalid/local'. This call-site passes more inputs to this function than caused by: This argument is not constant, and therefore does not match against a name-value argument 'myNamedArg_invalid/local' during code generation. Code generation might fail or produce results if a name passed at a call site is not known during code generation.

Error in ==> myNamedArg_invalid Line: 2 Column: 17
Code generation failed: View Error Report

In certain situations, the code generator assigns the name that you passed to an optional positional or repeating input argument. In such situations, code generation succeeds with a warning and the generated code might produce results that are different from MATLAB execution. See “Passing Input Argument Name at Run Time” (MATLAB Coder).

Using the Structure That Holds Name-Value Arguments

Suppose that your MATLAB function for which you intend to generate code uses a structure named `NameValueArgs` to define two name-value arguments, `Name1` and `Name2`:

```
function result = myFunction(NameValueArgs)
arguments
    NameValueArgs.Name1
    NameValueArgs.Name2
end
...
```

In the body of your function, directly use the structure fields `NameValueArgs.Name1` and `NameValueArgs.Name2` to read or write data.

Do not use the whole structure variable `NameValueArgs` itself (without the dot syntax), except in these situations:

- Use the `isfield` function to check if the caller has supplied a value for a certain name-value argument. For example, to provide a default value for `NameValueArgs.Name2` outside of the arguments block, you can use this code snippet:

```
if ~isfield(NameValueArgs, 'Name2')
    NameValueArgs.Name2 = defaultValue;
end
```

- Use the `namedargs2cell` function to forward the name-value arguments to another function. For example:

```
argsCell = namedargs2cell(NameValueArgs);
foo(argsCell{:});
```

Any use of the whole structure variable `NameValueArgs` (including the above two special cases) is not supported inside loops, anonymous functions, or nested functions.

Differences Between Generated Code and MATLAB Code

Certain unusual code patterns might cause the code generated for argument validation to behave differently from MATLAB. To learn about some of these differences, see these links:

- “Empty Repeating Input Argument” (MATLAB Coder)
- “Passing Input Argument Name at Run Time” (MATLAB Coder)
- “Output Argument Validation of Conditionally-Assigned Outputs” (MATLAB Coder)

Input Type Specification and arguments blocks

Using function argument validation (`arguments` blocks) to specify input types of entry-point functions is not supported. Even if your entry-point function contains `arguments` blocks that validate the input arguments, you must specify the properties of these input arguments by using one of the three approaches listed in “Methods for Defining Properties of Primary Inputs” (MATLAB Coder).

Default Values for Entry-Point Function Inputs in Generated Code

The `arguments` block allows you to specify default values for one or more positional input arguments. Specifying a default value in the argument declaration makes a positional argument optional because MATLAB can use the default value when you do not pass a value in the function call. When you generate code by using the `codegen` command or accelerate fixed-point code by using the `fiaccel` command, you can choose to not specify the properties of one or more optional positional arguments that have constant default values. In such situations, the default values of these optional arguments are hard-coded in the generated code and these arguments do not appear in the generated code interface. For examples, see the following table.

| MATLAB Code | Generated Code |
|--|---|
| <pre>function out = useDefaults_1(a,b,c) arguments a (1,1) double = 3 b (1,1) double = 5 c (1,1) double = 7 end out = a + b + c; end</pre> | <pre>codegen command: codegen -config:lib -c useDefaults_1 -args {} -report Generated code: double useDefaults_1(void) { return 15.0; }</pre> |
| <pre>function out = useDefaults_2(a,b,c) arguments a (1,1) double b (1,1) double = 5 c (1,1) double = 7 end out = a + b + c; end</pre> | <pre>codegen command: codegen -config:lib -c useDefaults_2 -args 0 -report Generated code: double useDefaults_2(double a) { return (a + 5.0) + 7.0; }</pre> |
| | <pre>codegen command: codegen -config:lib -c useDefaults_2 -args {0,0} -rep Generated code: double useDefaults_2(double a, double b) { return (a + b) + 7.0; }</pre> |

See Also

arguments

Related Examples

- “Function Argument Validation”
- “Methods for Defining Properties of Primary Inputs” on page 31-4

Specify Number of Entry-Point Function Input or Output Arguments to Generate

You can control the number of input or output arguments in a generated entry-point function. From one MATLAB function, you can generate entry-point functions that have different signatures.

Control Number of Input Arguments

If your entry-point function uses `varargin`, specify the properties for the arguments that you want in the generated function.

Consider this function:

```
function [x, y] = myops(varargin)
%#codegen
if (nargin > 1)
    x = varargin{1} + varargin{2};
    y = varargin{1} * varargin{2};
else
    x = varargin{1};
    y = -varargin{1};
end
```

To generate a function that takes only one argument, provide one argument with `-args`.

```
fiaccl myops -args {fi(3, 1, 16, 13)} -report
```

You can also control the number of input arguments when the MATLAB function does not use `varargin`.

Consider this function:

```
function [x, y] = myops(a,b)
%#codegen
if (nargin > 1)
    x = a + b;
    y = a * b;
else
    x = a;
    y = -a;
end
```

To generate a function that takes only one argument, provide one argument with `-args`.

```
fiaccl myops -args {fi(3, 1, 16, 13)} -report
```

Control the Number of Output Arguments

When you use `fiaccl`, you can specify the number of output arguments by using the `-nargout` option.

Consider this function:

```
function [x, y] = myops(a,b)
%#codegen
```

```
x = a + b;
y = a * b;
end
```

Generate a function that has one output argument.

```
fiaccel myops -args {fi(3,1,16,13) fi(3,1,16,13)} -nargout 1 -report
```

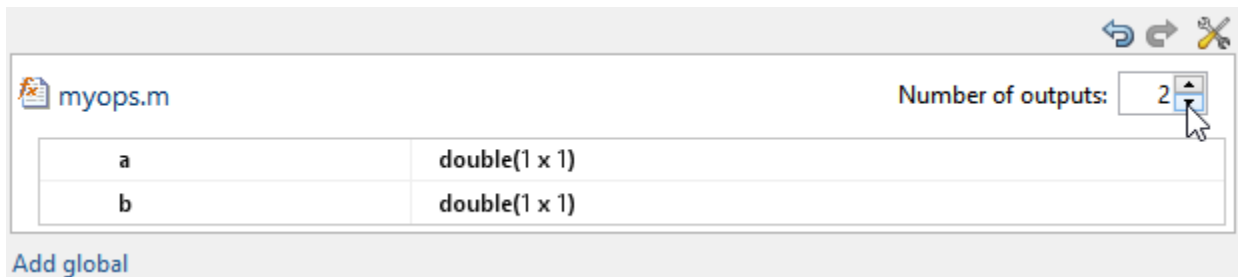
You can also use `-nargout` to specify the number of output arguments for an entry-point function that uses `varargout`.

Rewrite `myops` to use `varargout`.

```
function varargout = myops(a,b)
%#codegen
varargout{1} = a + b;
varargout{2} = a * b;
end
```

Generate code for one output argument.

```
fiaccel myops -args {fi(3,1,16,13) fi(3,1,16,13)} -nargout 1 -report
```



See Also

More About

- “Code Generation for Variable Length Argument Lists” on page 17-2
- “Specify Properties of Entry-Point Function Inputs” on page 31-2

Code Generation for Anonymous Functions

You can use anonymous functions in MATLAB code intended for code generation. For example, you can generate code for the following MATLAB code that defines an anonymous function that finds the square of a number.

```
sqr = @(x) x.^2;  
a = sqr(5);
```

Anonymous functions are useful for creating a function handle to pass to a MATLAB function that evaluates an expression over a range of values. For example, this MATLAB code uses an anonymous function to create the input to the `fzero` function:

```
b = 2;  
c = 3.5;  
x = fzero(@(x) x^3 + b*x + c,0);
```

Anonymous Function Limitations for Code Generation

Anonymous functions have the code generation limitations of value classes and cell arrays.

See Also

More About

- “MATLAB Classes Definition for Code Generation” on page 15-2
- “Cell Array Limitations for Code Generation” on page 30-7
- “Parameterizing Functions”

Code Generation for Nested Functions

You can generate code for MATLAB functions that contain nested functions. For example, you can generate code for the function `parent_fun`, which contains the nested function `child_fun`.

```
function parent_fun
x = 5;
child_fun

    function child_fun
        x = x + 1;
    end

end
```

Nested Function Limitations for Code Generation

When you generate code for nested functions, you must adhere to the code generation restrictions for value classes, cell arrays, and handle classes. You must also adhere to these restrictions:

- If the parent function declares a persistent variable, it must assign the persistent variable before it calls a nested function that uses the persistent variable.
- A nested recursive function cannot refer to a variable that the parent function uses.
- If a nested function refers to a structure variable, you must define the structure by using `struct`.
- If a nested function uses a variable defined by the parent function, you cannot use `coder.varsize` with the variable in either the parent or the nested function.

See Also

More About

- “MATLAB Classes Definition for Code Generation” on page 15-2
- “Handle Object Limitations for Code Generation” on page 15-19
- “Cell Array Limitations for Code Generation” on page 30-7
- “Code Generation for Recursive Functions” on page 14-12

Defining MATLAB Variables for C/C++ Code Generation

- “Variables Definition for Code Generation” on page 18-2
- “Best Practices for Defining Variables for C/C++ Code Generation” on page 18-3
- “Eliminate Redundant Copies of Variables in Generated Code” on page 18-7
- “Reassignment of Variable Properties” on page 18-9
- “Reuse the Same Variable with Different Properties” on page 18-10
- “Supported Variable Types” on page 18-13
- “Edit and Represent Coder Type Objects and Properties” on page 18-14

Variables Definition for Code Generation

In MATLAB, variables can change their properties dynamically at run time so you can use the same variable to hold a value of any class, size, or complexity. For example, the following code works in MATLAB:

```
function x = foo(c) %#codegen
if(c>0)
    x = 0;
else
    x = [1 2 3];
end
disp(x);
end
```

However, statically-typed languages like C must be able to determine variable properties at compile time. Therefore, for C/C++ code generation, you must explicitly define the class, size, and complexity of variables in MATLAB source code before using them. For example, rewrite the above source code with a definition for `x`:

```
function x = foo(c) %#codegen
x = zeros(1,3);
if(c>0)
    x = 0;
else
    x = [1 2 3];
end
disp(x);
end
```

See Also

Related Examples

- “Best Practices for Defining Variables for C/C++ Code Generation” on page 18-3

Best Practices for Defining Variables for C/C++ Code Generation

In this section...

“Explicitly Define Variables Before Using Them” on page 18-3
 “Use Caution When Reassigning Variable Properties” on page 18-4
 “Define Variable Numeric Data Types” on page 18-5
 “Define Matrices Before Assigning Indexed Variables” on page 18-5
 “Index Arrays by Using Constant Value Vectors” on page 18-5

MATLAB code used for code generation must adhere to additional restrictions. When you define variables, follow these best practices to optimize variable use and avoid generating errors in your code.

Explicitly Define Variables Before Using Them

For C/C++ code generation, you must explicitly define variable values and properties before using them in operations or returning them as outputs. Doing this prevents errors that occur when the variable is not defined.

Note When you define variables, they are local by default and do not persist between function calls. To make variables persistent, use the `persistent` function.

Initializing a new variable sometimes results in redundant copies in the generated code. For more information, see “Eliminate Redundant Copies of Variables in Generated Code” on page 18-7.

Define Variables on All Execution Paths

You must define a variable on all execution paths, such as execution paths dictated by `if` statements. Consider this MATLAB code that defines a variable before using it as an input to a function:

```

...
if c <= 0
    x = 11;
end
% Later in your code ...
if c > 0
    % Use x in the function, foo
    foo(x);
end
...
  
```

The code assigns `x` to a value only if `c <= 0` and uses `x` only when `c > 0`. Depending on the value for `c`, this code can work in MATLAB without errors. However, the MATLAB code generates a compilation error when you generate C/C++ code from it because the code generator detects that `x` is undefined on the execution path when `c > 0`.

To make this code suitable for code generation, define `x` before using it:

```

x = 0;
...
  
```

```
if c <= 0
    x = 11;
end
% Later in your code ...
if c > 0
% Use x in the function, foo
    foo(x);
end
...
```

Define Fields in a Structure

You must also define each structure field for all execution paths. Consider this MATLAB code:

```
...
if c > 0
    s.a = 11;
    disp(s);
else
    s.a = 12;
    s.b = 12;
end
% Use s in the function, foo
foo(s);
...
```

The first part of the `if` statement uses only the field `a`, and the `else` statement uses fields `a` and `b`. This code works in MATLAB, but generates a compilation error during C/C++ code generation. To prevent this error, do not add fields to a structure after you use the structure. For more information, see “Structure Definition for Code Generation” on page 25-2.

To make this code suitable for C/C++ code generation, define the fields of `s` before using them.

```
...
% Define fields in structure s
s = struct('a',0, 'b', 0);
if c > 0
    s.a = 11;
    disp(s);
else
    s.a = 12;
    s.b = 12;
end
% Use s in the function, foo
foo(s);
...
```

Use Caution When Reassigning Variable Properties

You can reassign certain variables after the initial assignment with a value of different class, size, or complexity, as described in “Reassignment of Variable Properties” on page 18-9. However, if you reassign the variable type after the initial assignment, the code often returns a compilation error during code generation. In general, assign each variable a specific class, size, type, and complexity.

Define Variable Numeric Data Types

`double` is the default numeric data type in MATLAB. To define variables of other data types, you must explicitly define the data type in the definition with the correct prefix or operator. Be mindful of the data types you use, because using variables assigned to different data types in your code can cause type mismatch errors.

For example, this code defines the variable `x` as a double and `y` as an 8-bit integer:

```
x = 15;
y = uint8(x);
```

For more information on supported types in MATLAB, see “Numeric Types”.

Define Matrices Before Assigning Indexed Variables

Growing a variable by writing an element beyond its current size results in a compile-time or run-time error. You must define the array before assigning values to its elements.

The only situation in which code generation supports growing arrays via indexing is when you use `end + 1`. See “Generate Code for Growing Arrays and Cell Arrays with `end + 1` Indexing” (MATLAB Coder).

For example, this assignment results in an error:

```
g(3,2) = 14.6;
```

Correct this code by defining the matrix with a specific size first.

```
g = zeros(5,5);
g(3,2) = 14.6;
```

For more information about indexing matrices, see “Incompatibility with MATLAB in Matrix Indexing Operations for Code Generation” on page 29-22.

You can grow the size of variables if you define the variable as variable size first. See “Define Variable-Size Data Explicitly by Using `coder.varsizes`” on page 29-12.

Index Arrays by Using Constant Value Vectors

Although code generation can support variable size arrays, variable size arrays requires additional memory, which can slow performance. When possible, use constant value vectors as you index arrays.

If you need to index through an array, be careful when using the colon operator. In some cases the code generator does not correctly determine whether an array indexed with a the colon operator is fixed size or variable size. As a result, you may define an array that does not change size, and should therefore be fixed-sized, but the code generator specifies the array as variable-sized.

For example, this code creates the array `out` by using the variable `i` indexed through the random row vector `A`:

```
...
% extract elements -1+2*i through 5+2*i for processing
```

```
A = rand(1,10);  
out = A(-1+2*i:5+2*i);
```

In this example, *A* is always the same size. If *i* is a compile-time constant value, the code generator produces a fixed size object for *out*. If *i* is unknown at compile time, the code generator produces a variable size array for *out*.

By contrast, this example produces a fixed-sized array if *i* is known or unknown at compile time:

```
...  
% extract elements i through 2+i for processing  
A = rand(1,10);  
out = A(i:2+i);
```

Occasionally, you can rewrite code to generate fixed size arrays. In this example, the first array is variable-sized, and second is fixed-sized, despite both producing the same array:

```
...  
width = 25;  
A = A(j-width:j+width); % A is variable-size, if j is unknown at compile time  
fsA = A(j+(-width:width)); % This makes A fixed-size, even if j is unknown at compile time  
...
```

See Also

`coder.nullcopy` | `persistent`

Related Examples

- “Eliminate Redundant Copies of Variables in Generated Code” on page 18-7
- “Structure Definition for Code Generation” on page 25-2
- “Incompatibility with MATLAB in Matrix Indexing Operations for Code Generation” on page 29-22
- “Avoid Data Copies of Function Inputs in Generated Code” (MATLAB Coder)

Eliminate Redundant Copies of Variables in Generated Code

In this section...

“When Redundant Copies Occur” on page 18-7

“How to Eliminate Redundant Copies by Defining Uninitialized Variables” on page 18-7

“Defining Uninitialized Variables” on page 18-7

When Redundant Copies Occur

During C/C++ code generation, the code generator checks for statements that attempt to access uninitialized memory. If it detects execution paths where a variable is used but is potentially not defined, it generates a compile-time error. To prevent these errors, define variables by assignment before using them in operations or returning them as function outputs.

Note, however, that variable assignments not only copy the properties of the assigned data to the new variable, but also initialize the new variable to the assigned value. This forced initialization sometimes results in redundant copies in C/C++ code. To eliminate redundant copies, define uninitialized variables by using the `coder.nullcopy` function, as described in “How to Eliminate Redundant Copies by Defining Uninitialized Variables” on page 18-7.

How to Eliminate Redundant Copies by Defining Uninitialized Variables

- 1 Define the variable with `coder.nullcopy`.
- 2 Initialize the variable before reading it.

When the uninitialized variable is an array, you must initialize all of its elements before passing the array as an input to a function or operator — even if the function or operator does not read from the uninitialized portion of the array.

What happens if you access uninitialized data?

Uninitialized memory contains arbitrary values. Therefore, accessing uninitialized data may lead to segmentation violations or nondeterministic program behavior (different runs of the same program may yield inconsistent results).

Defining Uninitialized Variables

In the following code, the assignment statement `X = zeros(1,N)` not only defines `X` to be a 1-by-5 vector of real doubles, but also initializes each element of `X` to zero.

```
function X = withoutNullcopy %#codegen

N = 5;
X = zeros(1,N);
for i = 1:N
    if mod(i,2) == 0
        X(i) = i;
    elseif mod(i,2) == 1
        X(i) = 0;
```

```
end  
end
```

This forced initialization creates an extra copy in the generated code. To eliminate this overhead, use `coder.nullcopy` in the definition of `X`:

```
function X = withNullcopy %#codegen  
  
N = 5;  
X = coder.nullcopy(zeros(1,N));  
for i = 1:N  
    if mod(i,2) == 0  
        X(i) = i;  
    else  
        X(i) = 0;  
    end  
end
```

See Also

`coder.nullcopy`

More About

- “Avoid Data Copies of Function Inputs in Generated Code” (MATLAB Coder)

Reassignment of Variable Properties

For C/C++ code generation, there are certain variables that you can reassign after the initial assignment with a value of different class, size, or complexity:

Dynamically sized variables

A variable can hold values that have the same class and complexity but different sizes. If the size of the initial assignment is not constant, the variable is dynamically sized in generated code. For more information, see “Variable-Size Data”.

Variables reused in the code for different purposes

You can reassign the type (class, size, and complexity) of a variable after the initial assignment if each occurrence of the variable can have only one type. In this case, the variable is renamed in the generated code to create multiple independent variables. For more information, see “Reuse the Same Variable with Different Properties” on page 18-10.

Reuse the Same Variable with Different Properties

In this section...

“When You Can Reuse the Same Variable with Different Properties” on page 18-10

“When You Cannot Reuse Variables” on page 18-10

“Limitations of Variable Reuse” on page 18-11

For C/C++ code generation, there are certain variables that you can reassign after the initial assignment with a value of different class, size, or complexity. A variable can hold values that have the same class and complexity but different sizes. If the size of the initial assignment is not constant, the variable is dynamically sized in generated code. For more information, see “Variable-Size Data”.

You can reassign the type (class, size, and complexity) of a variable after the initial assignment if each occurrence of the variable can have only one type. In this case, the variable is renamed in the generated code to create multiple independent variables.

When You Can Reuse the Same Variable with Different Properties

You can reuse (reassign) an input, output, or local variable with different class, size, or complexity if the code generator can unambiguously determine the properties of each occurrence of this variable during C/C++ code generation. If so, MATLAB creates separate uniquely named local variables in the generated code. You can view these renamed variables in the code generation report.

A common example of variable reuse is in `if-elseif-else` or `switch-case` statements. For example, the following function `example1` first uses the variable `t` in an `if` statement, where it holds a scalar double, then reuses `t` outside the `if` statement to hold a vector of doubles.

```
function y = example1(u) %#codegen
if all(all(u>0))
    % First, t is used to hold a scalar double value
    t = mean(mean(u)) / numel(u);
    u = u - t;
end
% t is reused to hold a vector of doubles
t = find(u > 0);
y = sum(u(t(2:end-1)));
```

When You Cannot Reuse Variables

You cannot reuse (reassign) variables if it is not possible to determine the class, size, and complexity of an occurrence of a variable unambiguously during code generation. In this case, variables cannot be renamed and a compilation error occurs.

For example, the following `example2` function assigns a fixed-point value to `x` in the `if` statement and reuses `x` to store a matrix of doubles in the `else` clause. It then uses `x` after the `if-else` statement. This function generates a compilation error because after the `if-else` statement, variable `x` can have different properties depending on which `if-else` clause executes.

```
function y = example2(use_fixpoint, data) %#codegen
if use_fixpoint
    % x is fixed-point
    x = fi(data, 1, 12, 3);
```

```

else
    % x is a matrix of doubles
    x = data;
end
% When x is reused here, it is not possible to determine its
% class, size, and complexity
t = sum(sum(x));
y = t > 0;
end

```

Example 18.1. Variable Reuse in an if Statement

To see how MATLAB renames a reused variable `t`:

- 1 Create a MATLAB file `example1.m` containing the following code.

```

function y = example1(u) %#codegen
if all(all(u>0))
    % First, t is used to hold a scalar double value
    t = mean(mean(u)) / numel(u);
    u = u - t;
end
% t is reused to hold a vector of doubles
t = find(u > 0);
y = sum(u(t(2:end-1)));
end

```

- 2 Generate a MEX function for `example1` and produce a code generation report.

```
codegen -o example1x -report example1.m -args {ones(5,5)}
```

- 3 Open the code generation report.

On the **Variables** tab, you see two uniquely named local variables `t>1` and `t>2`.

| SUMMARY | ALL MESSAGES (0) | BUILD LOGS | | COD |
|---------|------------------|------------|--------|-----|
| Name | Type | Size | Class | |
| y | Output | 1 × 1 | double | |
| u | Input | 5 × 5 | double | |
| t > 1 | Local | 1 × 1 | double | |
| t > 2 | Local | :25 × 1 | double | |

- 4 In the list of variables, click `t>1`. The report highlights the instances of the variable `t` that are inside of the `if` statement. These instances of `t` are scalar double.
- 5 Click `t>2`. The code generation report highlights the instances of `t` that are outside of the `if` statement. These instances of `t` are variable-size column vectors with an upper bound of 25.

Limitations of Variable Reuse

The following variables cannot be renamed in generated code:

- Persistent variables.
- Global variables.

- Variables passed to C code using `coder.ref`, `coder.rref`, `coder.wref`.
- Variables whose size is set using `coder.varsize`.
- Variables whose names are controlled using `coder.cstructname`, when generating code by using MATLAB Coder.
- The index variable of a `for`-loop when it is used inside the loop body.
- The block outputs of a MATLAB Function block in a Simulink model.
- Chart-owned variables of a MATLAB function in a Stateflow chart.

Supported Variable Types

You can use the following data types for C/C++ code generation from MATLAB:

| Type | Description |
|-------------------------------|---|
| char | Character array |
| complex | Complex data. Cast function takes real and imaginary components |
| double | Double-precision floating point |
| int8, int16, int32, int64 | Signed integer |
| logical | Boolean true or false |
| single | Single-precision floating point |
| struct | Structure |
| uint8, uint16, uint32, uint64 | Unsigned integer |
| Fixed-point | Fixed-point data types |

Edit and Represent Coder Type Objects and Properties

Passing an object to `coder.typeof` or passing a class name as a string scalar to `coder.newtype` creates an object that represents the type of object for code generation.

The coder type object displays a succinct description of the object properties while excluding internal state values. Nonconstant properties display their type and size, while constant properties display only their values.

To create a coder type object, pass a compatible object to `coder.typeof`. For example:

```
t = categorical({'r','g','b'});
tType = coder.typeof(t)
```

The representation of variable `t` is stored in coder type object `tType`.

```
tType =
    matlab.coder.type.CategoricalType
    1x3 categorical
    Categories : 3x1 homogeneous cell
    Ordinal : 1x1 logical
    Protected : 1x1 logical
```

Object Properties

You can edit the properties of coder type objects. You can assign scalar values to the object properties. Values are implicitly converted to the corresponding coder type values when they are assigned to coder type object properties. The code generator implicitly converts constants assigned to coder type object properties to `coder.Constant` values. You can resize objects themselves

Resize Objects by Using `coder.resize`

You can resize most objects by using `coder.resize`. You can resize objects, its properties and create arrays within the properties.

For example, for a `timetable` coder object, you can resize the object:

```
t = timetable((1:5)',(11:15)', 'SampleRate',1);
tType = coder.typeof(t);
tType = coder.resize(tType, [10 2],[1 0])
```

This code resizes the `timetable` to a `:10x2` object.

```
tType =
    matlab.coder.type.RegularTimetableType
    :10x2 timetable
           Data : 1x2 homogeneous cell
           Description : 1x0 char
           UserData : 0x0 double
           DimensionNames : {'Time'}    {'Variables'}
           VariableNames : {'Var1'}    {'Var2'}
           VariableDescriptions : 1x2 homogeneous cell
           VariableUnits : 1x2 homogeneous cell
           VariableContinuity : 1x2 matlab.internal.coder.tabular.Continuity
```

```

StartTime : 1x1 matlab.coder.type.DurationType
SampleRate : 1x1 double
TimeStep : 1x1 matlab.coder.type.DurationType

```

The constant properties of `tType` display their values. The nonconstant properties display only their type and size.

Note Not all types representing MATLAB classes are compatible with `coder.resize`.

Resize Objects by Editing Object Properties

You can resize the objects by editing the properties themselves. For a `duration` coder type object `x`, edit the `Size` property to change the size as needed.

```

x = coder.typeof(duration((1:3),0,0));
x.Size = [10 10]

```

This code changes the size of the coder type object.

```

x =
    matlab.coder.type.DurationType
    10x10 duration
    Format : 1x8 char

```

You can also make the coder type object variable-size by setting the `VarDims` flag:

```

x.VarDims(2) = true

```

The second dimension of the coder type object is upper-bounded at 10.

```

x =
    matlab.coder.type.DurationType
    10x:10 duration
    Format : 1x8 char

```

Legacy Representation of Coder Type Objects

In R2021a, calling `coder.typeof` no longer returns a `coder.ClassType` object. If your workflow requires the legacy representation of coder type objects, use the `getCoderType` function on the variable that has the new representation of your class or object. For example, to get the legacy representation of a `datetime` variable, use the variable that has the new representation `tt` to call the `getCoderType` function:

```

t = datetime;
tt = coder.typeof(t);
ttLegacy = tt.getCoderType()

```

In the Coder Type Editor, the code generator includes the function `getCoderType` for coder type objects. Use this function to return the legacy representation of coder types. See, “Create and Edit Input Types by Using the Coder Type Editor” (MATLAB Coder)

Certain MATLAB data types provide customized type representations for MATLAB code generation. In other cases, the type is represented using `coder.ClassType`.

See Also

`coder.resize` | `coder.newtype` | `coder.typeof` | “Code Generation for Variable-Size Arrays”
(MATLAB Coder)

Design Considerations for C/C++ Code Generation

- “When to Generate Code from MATLAB Algorithms” on page 19-2
- “Which Code Generation Feature to Use” on page 19-3
- “Prerequisites for C/C++ Code Generation from MATLAB” on page 19-4
- “MATLAB Code Design Considerations for Code Generation” on page 19-5
- “Differences Between Generated Code and MATLAB Code” on page 19-6
- “Potential Differences Reporting” on page 19-23
- “Potential Differences Messages” on page 19-25
- “MATLAB Language Features Supported for C/C++ Code Generation” on page 19-29

When to Generate Code from MATLAB Algorithms

Generating code from MATLAB algorithms for desktop and embedded systems allows you to perform your software design, implementation, and testing completely within the MATLAB workspace. You can:

- Verify that your algorithms are suitable for code generation
- Generate efficient, readable, and compact C/C++ code automatically, which eliminates the need to manually translate your MATLAB algorithms and minimizes the risk of introducing errors in the code.
- Modify your design in MATLAB code to take into account the specific requirements of desktop and embedded applications, such as data type management, memory use, and speed.
- Test the generated code and easily verify that your modified algorithms are functionally equivalent to your original MATLAB algorithms.
- Generate MEX functions to:
 - Accelerate MATLAB algorithms in certain applications.
 - Speed up fixed-point MATLAB code.
- Generate hardware description language (HDL) from MATLAB code.

When Not to Generate Code from MATLAB Algorithms

Do not generate code from MATLAB algorithms for the following applications. Use the recommended MathWorks product instead.

| To: | Use: |
|---|--|
| Deploy an application that uses handle graphics | MATLAB Compiler |
| Use Java | MATLAB Compiler SDK™ |
| Use toolbox functions that do not support code generation | Toolbox functions that you rewrite for desktop and embedded applications |
| Deploy MATLAB based GUI applications on a supported MATLAB host | MATLAB Compiler |
| Deploy web-based or Windows applications | MATLAB Compiler SDK |
| Interface C code with MATLAB | MATLAB mex function |

Which Code Generation Feature to Use

| To... | Use... | Required Product | To Explore Further... |
|--|-----------------------|------------------------|--|
| Generate MEX functions for verifying generated code | codegen function | MATLAB Coder | Try this in “Accelerate MATLAB Algorithm by Generating MEX Function” (MATLAB Coder). |
| Produce readable, efficient, and compact code from MATLAB algorithms for deployment to desktop and embedded systems. | MATLAB Coder app | MATLAB Coder | Try this in “Generate C Code by Using the MATLAB Coder App” (MATLAB Coder). |
| | codegen function | MATLAB Coder | Try this in “Generate C Code at the Command Line” (MATLAB Coder). |
| Generate MEX functions to accelerate MATLAB algorithms | MATLAB Coder app | MATLAB Coder | See “Accelerate MATLAB Algorithms” (MATLAB Coder). |
| | codegen function | MATLAB Coder | |
| Integrate MATLAB code into Simulink | MATLAB Function block | Simulink | Try this in “Call MATLAB Function Files in MATLAB Function Blocks”. |
| Speed up fixed point MATLAB code | fiaccel function | Fixed-Point Designer | Learn more in “Code Acceleration and Code Generation from MATLAB” on page 12-2. |
| Integrate custom C code into MATLAB and generate efficient, readable code | codegen function | MATLAB Coder | Learn more in “Call Custom C/C++ Code from the Generated Code” (MATLAB Coder). |
| Integrate custom C code into code generated from MATLAB | coder.ceval function | MATLAB Coder | Learn more in coder.ceval. |
| Generate HDL from MATLAB code | MATLAB Function block | Simulink and HDL Coder | Learn more at www.mathworks.com/products/slhdlcoder . |

Prerequisites for C/C++ Code Generation from MATLAB

To generate C/C++ or MEX code from MATLAB algorithms, you must install the following software:

- MATLAB Coder product
- C/C++ compiler

MATLAB Code Design Considerations for Code Generation

When writing MATLAB code that you want to convert into efficient, standalone C/C++ code, you must consider the following:

- Data types

C and C++ use static typing. To determine the types of your variables before use, MATLAB Coder requires a complete assignment to each variable.

- Array sizing

Variable-size arrays and matrices are supported for code generation. You can define inputs, outputs, and local variables in MATLAB functions to represent data that varies in size at run time.

- Memory

You can choose whether the generated code uses static or dynamic memory allocation.

With dynamic memory allocation, you potentially use less memory at the expense of time to manage the memory. With static memory, you get better speed, but with higher memory usage. Most MATLAB code takes advantage of the dynamic sizing features in MATLAB, therefore dynamic memory allocation typically enables you to generate code from existing MATLAB code without modifying it much. Dynamic memory allocation also allows some programs to compile even when upper bounds cannot be found.

Static allocation reduces the memory footprint of the generated code, and therefore is suitable for applications where there is a limited amount of available memory, such as embedded applications.

- Speed

Because embedded applications must run in real time, the code must be fast enough to meet the required clock rate.

To improve the speed of the generated code:

- Choose a suitable C/C++ compiler. Do not use the default compiler that MathWorks supplies with MATLAB for Windows 64-bit platforms.
- Consider disabling run-time checks.

By default, for safety, the code generated for your MATLAB code contains memory integrity checks and responsiveness checks. Generally, these checks result in more generated code and slower simulation. Disabling run-time checks usually results in streamlined generated code and faster simulation. Disable these checks only if you have verified that array bounds and dimension checking is unnecessary.

See Also

- “Data Definition” (MATLAB Coder) “Numeric Types”
- “Code Generation for Variable-Size Arrays” on page 29-2
- “Control Run-Time Checks” on page 12-48

Differences Between Generated Code and MATLAB Code

To convert MATLAB code to efficient C/C++ code, the code generator introduces optimizations that intentionally cause the generated code to behave differently, and sometimes produce different results, than the original source code.

Here are some of the differences:

- “Functions that have Multiple Possible Outputs” on page 19-7
- “Passing Input Argument Name at Run Time” (MATLAB Coder)
- “Empty Repeating Input Argument” (MATLAB Coder)
- “Output Argument Validation of Conditionally-Assigned Outputs” (MATLAB Coder)
- “Writing to ans Variable” on page 19-10
- “Logical Short-Circuiting” on page 19-11
- “Loop Index Overflow” on page 19-11
- “Indexing for Loops by Using Single Precision Operands” (MATLAB Coder)
- “Index of an Unentered for Loop” (MATLAB Coder)
- “Character Size” on page 19-14
- “Order of Evaluation in Expressions” on page 19-14
- “Name Resolution While Constructing Function Handles” on page 19-15
- “Termination Behavior” on page 19-16
- “Size of Variable-Size N-D Arrays” on page 19-16
- “Size of Empty Arrays” on page 19-17
- “Size of Empty Array That Results from Deleting Elements of an Array” on page 19-17
- “Growing Variable-Size Column Cell Array That is Initialized as Scalar at Run Time” on page 19-17
- “Binary Element-Wise Operations with Single and Double Operands” on page 19-18
- “Floating-Point Numerical Results” on page 19-19
- “NaN and Infinity” on page 19-19
- “Negative Zero” on page 19-19
- “Code Generation Target” on page 19-20
- “MATLAB Class Property Initialization” on page 19-20
- “MATLAB Classes in Nested Property Assignments That Have Set Methods” on page 19-20
- “MATLAB Handle Class Destructors” on page 19-20
- “Variable-Size Data” on page 19-21
- “Complex Numbers” on page 19-21
- “Converting Strings with Consecutive Unary Operators to double” on page 19-21
- “Display Function” (MATLAB Coder)

These differences are applicable for:

- MEX and standalone C/C++ code generation by using the `codegen` command or the MATLAB Coder app.

- Fixed-point code acceleration by generating MEX using the `fiaccl` command.
- MATLAB Function block simulation using Simulink.

When you run your generated `fiaccl` MEX, C/C++ MEX or standalone C/C++ code, run-time error checks can detect some of these differences. By default, run-time error checks are enabled for MEX code and disabled for standalone C/C++ code. To help you identify and address differences before you deploy code, the code generator reports a subset of the differences as potential differences on page 19-23.

Functions that have Multiple Possible Outputs

Certain mathematical operations, such as singular value decomposition and eigenvalue decomposition of a matrix, can have multiple answers. Two different algorithms implementing such an operation can return different outputs for identical input values. Two different implementations of the same algorithm can also exhibit the same behavior.

For such mathematical operations, the corresponding functions in the generated code and MATLAB might return different outputs for identical input values. To see if a function has this behavior, in the corresponding function reference page, see the **C/C++ Code Generation** section under **Extended Capabilities**. Examples of such functions include `svd` and `eig`.

Passing Input Argument Name at Run Time

Suppose that `foo` is a function that uses name-value argument validation. When you call `foo` from another function `bar`, the code generator must be able to determine the names that you provide to `foo` at compile time.

If the argument names are passed at run time, code generation fails in most situations. See “Names Must Be Compile-Time Constants” (MATLAB Coder).

In certain situations, the code generator assigns the name that you passed to an optional positional or repeating input argument. In such situations, code generation succeeds with a warning and the generated code might produce results that are different from MATLAB execution. For example, consider this function:

```
function out = myNamedArg_warns(a,b)
out = local(a,b);
end

function out = local(varargin,args)
arguments (Repeating)
    varargin
end

arguments
    args.x
    args.y
end

if isfield(args,'x') && isfield(args,'y')
    out = args.x / args.y;
elseif isfield(args,'x')
    out = args.x;
```

```
else
    out = varargin{1};
end
end
```

Behavior of MATLAB Execution

If you call `myNamedArg_warns` with `'x'` as the first input argument, MATLAB matches it against the first name-value argument of the function `local`.

```
myNamedArg_warns('x',5)
```

```
ans =
     5
```

By contrast, if you call `myNamedArg_warns` with `'z'` as the first input argument (that does not match with either name-value argument of `local`), MATLAB assigns the inputs into elements of `varargin`.

```
myNamedArg_warns('z',5)
```

```
ans =
    'z'
```

Behavior of Generated Code

Attempt to generate a MEX by running the `codegen` command. Specify the type of the first argument to be a character scalar and the second argument to be a double scalar. Code generation succeeds with a warning.

```
codegen myNamedArg_warns -args {'x',2}
```

```
Warning: This argument is not constant, and therefore does not match against a name-value argument during code generation. Code generation might fail or produce results that do not agree with MATLAB. This warning is not known during code generation.
```

```
Warning in ==> myNamedArg_warns Line: 2 Column: 13
Code generation successful (with warnings): View report
```

Irrespective of whether you pass `'x'` or `'z'` as the first input argument, the generated MEX assigns it to the first cell of `varargin`.

```
myNamedArg_warns_mex('x',5)
```

```
ans =
    'x'
```

```
myNamedArg_warns_mex('z',5)
```

```
ans =
    'z'
```

Workaround

To enable the code generator to match the first input against the name-value arguments of the function `local`, declare the first input to be a compile-time constant with value `'x'`. You can do this by using the `coder.Constant` function with the `-args` option of the `codegen` command.

```
codegen myNamedArg_warns -args {coder.Constant('x'),2}
```

```
Code generation successful.
```

Now, the behavior of the generated MEX agrees with MATLAB, although the MEX is unable to accept any value other than `'x'` for the first input.

```
myNamedArg_warns_mex('x',5)
```

```
ans =
```

```
    5
```

```
myNamedArg_warns_mex('z',5)
```

```
Constant function parameter 'a' has a different run-time value than the compile-time value.
```

```
Error in myNamedArg_warns_mex
```

Empty Repeating Input Argument

In code generation, if a repeating input argument (that is declared in an `arguments` block) is empty at run time, the size of that argument is 0×0 . By contrast, in MATLAB execution, the size of an empty repeating input argument is 1×0 .

For example, consider this function:

```
function out = testVararginSize
out = local;
end

function out = local(varargin)
arguments (Repeating)
    varargin
end
out = size(varargin);
end
```

Running `testVararginSize` in MATLAB returns `[1 0]`. If you generate a MEX for `testVararginSize` and run the generated MEX, you get `[0 0]`. However, iterating over elements of `varargin` by using `length(varargin)` or `numel(varargin)` produces the same behavior across MATLAB and code generation.

Output Argument Validation of Conditionally-Assigned Outputs

The code generator validates an output argument if the argument is assigned a type during code generation. By contrast, MATLAB execution validates an output argument if the argument is assigned a value when the MATLAB function returns.

In most situations, this underlying behavioral difference does not cause your generated code to behave differently than MATLAB. Here is an example function for which you do see this difference:

```
function outerFunc(in)
innerFunc(in);
end

function out = innerFunc(inputVal)
arguments (Output)
    out {mustBePositive}
end
if inputVal
    out = inputVal;
end
end
```

In MATLAB, the execution of `func` succeeds for all double inputs. If the input is positive, `out` is assigned this positive value and the validator `mustBePositive` runs without assertion. If the input is negative or zero, `out` is not assigned and is not validated.

Attempt to generate code for `func`. Specify the input type to be a double scalar.

```
codegen outerFunc -args 0
```

```
Variable 'out' is not fully defined on some execution paths.
```

```
Error in ==> outerFunc Line: 7 Column: 5
Code generation failed: View Error Report
```

Because the variable `out` is assigned a double scalar value on one execution path, code generation assigns a double scalar type to `out` at compile time. The code generator then attempts to perform validation on `out` and discovers that it is not fully defined if the `if` condition fails.

Writing to ans Variable

When you run MATLAB code that returns an output without specifying an output argument, MATLAB implicitly writes the output to the `ans` variable. If the variable `ans` already exists in the workspace, MATLAB updates its value to the output returned.

The code generated from such MATLAB code does not implicitly write the output to an `ans` variable.

For example, define the MATLAB function `foo` that explicitly creates an `ans` variable in the first line. The function then implicitly updates the value of `ans` when the second line executes.

```
function foo %#codegen
ans = 1;
2;
disp(ans);
end
```

Run `foo` at the command line. The final value of `ans`, which is 2, is displayed at the command line.

```
foo
```

```
2
```

Generate a MEX function from `foo`.

```
codegen foo
```

Run the generated MEX function `foo_mex`. This function explicitly creates the `ans` variable and assigns the value 1 to it. But `foo_mex` does not implicitly update the value of `ans` to 2.

```
foo_mex
```

```
1
```

Logical Short-Circuiting

Suppose that your MATLAB code has the logical operators `&` and `|` placed inside square brackets (`[]` and `]`). For such code patterns, the generated code does not employ short-circuiting behavior for these logical operators, but some MATLAB execution employs short-circuiting behavior. See “Tips” and “Tips”.

For example, define the MATLAB function `foo` that uses the `&` operator inside square brackets in the conditional expression of an `if . . . end` block.

```
function foo
if [returnsFalse() & hasSideEffects()]
end
end
```

```
function out = returnsFalse
out = false;
end
```

```
function out = hasSideEffects
out = true;
disp('This is my string');
end
```

The first argument of the `&` operator is always `false` and determines the value of the conditional expression. So, in MATLAB execution, short-circuiting is employed and the second argument is not evaluated. So, `foo` does not call the `hasSideEffects` function during execution and does not display anything at the command line.

Generate a MEX function for `foo`. Call the generated MEX function `foo_mex`.

```
foo_mex
```

```
This is my string
```

In the generated code, short-circuiting is not employed. So, the `hasSideEffects` function is called and the string is displayed at the command line.

Loop Index Overflow

Suppose that a `for`-loop end value is equal to or close to the maximum or minimum value for the loop index data type. In the generated code, the last increment or decrement of the loop index might cause the index variable to overflow. The index overflow might result in an infinite loop.

When memory integrity checks are enabled, if the code generator detects that the loop index might overflow, it reports an error. The software error checking is conservative. It might incorrectly report a loop index overflow. By default, memory-integrity checks are enabled for MEX code and disabled for

standalone C/C++ code. See “Why Test MEX Functions in MATLAB?” (MATLAB Coder) and “Generate Standalone C/C++ Code That Detects and Reports Run-Time Errors” (MATLAB Coder).

To avoid a loop index overflow, use the workarounds in this table.

| Loop Conditions Causing the Potential Overflow | Workaround |
|---|--|
| <ul style="list-style-type: none"> The loop index increments by 1. The end value equals the maximum value of the integer type. | <p>If the loop does not have to cover the full range of the integer type, rewrite the loop so that the end value is not equal to the maximum value of the integer type. For example, replace:</p> <pre>N=intmax('int16') for k=N-10:N</pre> <p>with:</p> <pre>for k=1:10</pre> |
| <ul style="list-style-type: none"> The loop index decrements by 1. The end value equals the minimum value of the integer type. | <p>If the loop does not have to cover the full range of the integer type, rewrite the loop so that the end value is not equal to the minimum value of the integer type. For example, replace:</p> <pre>N=intmin('int32') for k=N+10:-1:N</pre> <p>with:</p> <pre>for k=10:-1:1</pre> |
| <ul style="list-style-type: none"> The loop index increments or decrements by 1. The start value equals the minimum or maximum value of the integer type. The end value equals the maximum or minimum value of the integer type. | <p>If the loop must cover the full range of the integer type, cast the type of the loop start, step, and end values to a bigger integer or to double. For example, rewrite:</p> <pre>M= intmin('int16'); N= intmax('int16'); for k=M:N % Loop body end</pre> <p>as:</p> <pre>M= intmin('int16'); N= intmax('int16'); for k=int32(M):int32(N) % Loop body end</pre> |
| <ul style="list-style-type: none"> The loop index increments or decrements by a value not equal to 1. On the last loop iteration, the loop index is not equal to the end value. | <p>Rewrite the loop so that the loop index in the last loop iteration is equal to the end value.</p> |

Indexing for Loops by Using Single Precision Operands

Suppose in your MATLAB code, you are indexing a `for` loop that has a colon operator, where at least one of the colon operands is a single type operand and the number of iterations is greater than `flintmax('single') = 16777216`. When all these conditions are true, code generation might generate run-time or compile-time errors because the generated code calculates different values for the loop index variable than the values that MATLAB calculates.

For example, consider this MATLAB code:

```
function j = singlePIndex
n = flintmax('single') + 2;
j = single(0);
for i = single(1):single(n)
    j = i;
end
end
```

This code snippet executes in MATLAB, but it causes a compile-time or run-time error because the value of the loop index variable, `i`, is calculated differently in the generated code. The code generator displays a compile-time or run-time error and stops code generation or execution to prevent this discrepancy.

To avoid this discrepancy, replace the single type operands with double type or integer type operands.

For more information on run-time errors, see “Generate Standalone C/C++ Code That Detects and Reports Run-Time Errors” (MATLAB Coder).

Index of an Unentered for Loop

In your MATLAB code and generated code, after a `for` loop execution is complete, the value of the index variable is equal to its value during the final iteration of the `for` loop.

In MATLAB, if the loop does not execute, the value of the index variable is stored as `[]` (empty matrix). In generated code, if the loop does not execute, the value of the index variable is different than the MATLAB index variable.

- If you provide the `for` loop start and end variables at run time, the value of the index variable is equal to the start of the range. For example, consider this MATLAB code:

```
function out = indexTest(a,b)
for i = a:b
end
out = i;
end
```

Suppose that `a` and `b` are passed as `1` and `-1`. The `for` loop does not execute. In MATLAB, `out` is assigned `[]`. In the generated code, `out` is assigned the value of `a`, which is `1`.

- If you provide the `for` loop start and end values before compile time, the value of the index variable is assigned `[]` in both MATLAB and the generated code. Consider this MATLAB code:

```
function out = indexTest
for i = 1:-1
end
out = i;
end
```

In both MATLAB and the generated code, `out` is assigned `[]`.

Character Size

MATLAB supports 16-bit characters, but the generated code represents characters in 8 bits, the standard size for most embedded languages like C. See “Encoding of Characters in Code Generation” on page 16-12.

Order of Evaluation in Expressions

Generated code does not enforce the order of evaluation in expressions. For most expressions, the order of evaluation is not significant. For expressions that have side effects, the generated code might produce the side effects in a different order from the original MATLAB code. Expressions that produce side effects include those that:

- Modify persistent or global variables
- Display data to the screen
- Write data to files
- Modify the properties of handle class objects

In addition, the generated code does not enforce order of evaluation of logical operators that do not short circuit.

For more predictable results, it is good coding practice to split expressions that depend on the order of evaluation into multiple statements.

- Rewrite

```
A = f1() + f2();
```

as

```
A = f1();  
A = A + f2();
```

so that the generated code calls `f1` before `f2`.

- Assign the outputs of a multi-output function call to variables that do not depend on one another. For example, rewrite

```
[y, y.f, y.g] = foo;
```

as

```
[y, a, b] = foo;  
y.f = a;  
y.g = b;
```

- When you access the contents of multiple cells of a cell array, assign the results to variables that do not depend on one another. For example, rewrite

```
[y, y.f, y.g] = z{:};
```

as

```
[y, a, b] = z{:};
y.f = a;
y.g = b;
```

Name Resolution While Constructing Function Handles

MATLAB and code generation follow different precedence rules for resolving names that follow the symbol @. These rules do not apply to anonymous functions. The precedence rules are summarized in this table.

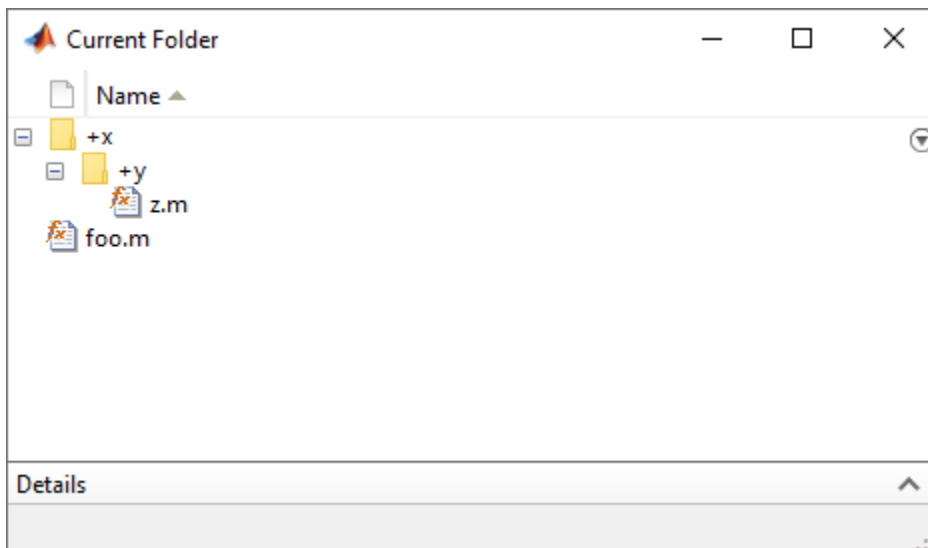
| Expression | Precedence Order in MATLAB | Precedence Order in Code Generation |
|--|--|--|
| An expression that does not contain periods, for example @x | Nested function, local function, private function, path function | Local variable, nested function, local function, private function, path function |
| An expression that contains exactly one period, for example @x.y | Local variable, path function | Local variable, path function (Same as MATLAB) |
| An expression that contains more than one period, for example @x.y.z | Path function | Local variable, path function |

If `x` is a local variable that is itself a function handle, generated code and MATLAB interpret the expression `@x` differently:

- MATLAB produces an error.
- Generated code interprets `@x` as the function handle of `x` itself.

Here is an example that shows this difference in behavior for an expression that contains two periods.

Suppose that your current working folder contains a package `x`, which contains another package `y`, which contains the function `z`. The current working folder also contains the entry-point function `foo` for which you want to generate code.



This is the definition for the file `foo`:

```
function out = foo
    x.y.z = @() 'x.y.z is an anonymous function';
    out = g(x);
end

function out = g(x)
    f = @x.y.z;
    out = f();
end
```

This is the definition for function `z`:

```
function out = z
    out = 'x.y.z is a package function';
end
```

Generate a MEX function for `foo`. Separately call both the generated MEX function `foo_mex` and the MATLAB function `foo`.

```
codegen foo
foo_mex
foo

ans =

    'x.y.z is an anonymous function'

ans =

    'x.y.z is a package function'
```

The generated code produces the first output. MATLAB produces the second output. Code generation resolves `@x.y.z` to the local variable `x` that is defined in `foo`. MATLAB resolves `@x.y.z` to `z`, which is within the package `x.y`.

Termination Behavior

Generated code does not match the termination behavior of MATLAB source code. For example, if infinite loops do not have side effects, optimizations remove them from generated code. As a result, the generated code can possibly terminate even though the corresponding MATLAB code does not.

Size of Variable-Size N-D Arrays

For variable-size N-D arrays, the `size` function might return a different result in generated code than in MATLAB source code. The `size` function sometimes returns trailing ones (singleton dimensions) in generated code, but always drops trailing ones in MATLAB. For example, for an N-D array `X` with dimensions `[4 2 1 1]`, `size(X)` might return `[4 2 1 1]` in generated code, but always returns `[4 2]` in MATLAB. See “Incompatibility with MATLAB in Determining Size of Variable-Size N-D Arrays” on page 29-19.

Size of Empty Arrays

The size of an empty array in generated code might be different from its size in MATLAB source code. See “Incompatibility with MATLAB in Determining Size of Empty Arrays” on page 29-19.

Size of Empty Array That Results from Deleting Elements of an Array

Deleting all elements of an array results in an empty array. The size of this empty array in generated code might differ from its size in MATLAB source code.

| Case | Example Code | Size of Empty Array in MATLAB | Size of Empty Array in Generated Code |
|---|---|-------------------------------|---------------------------------------|
| Delete all elements of an m-by-n array by using the colon operator (:). | <code>coder.ysize('X',[4,4],[1,1]); X = zeros(2); X(:) = [];</code> | 0-by-0 | 1-by-0 |
| Delete all elements of a row vector by using the colon operator (:). | <code>coder.ysize('X',[1,4],[0,1]); X = zeros(1,4); X(:) = [];</code> | 0-by-0 | 1-by-0 |
| Delete all elements of a column vector by using the colon operator (:). | <code>coder.ysize('X',[4,1],[1,0]); X = zeros(4,1); X(:) = [];</code> | 0-by-0 | 0-by-1 |
| Delete all elements of a column vector by deleting one element at a time. | <code>coder.ysize('X',[4,1],[1,0]); X = zeros(4,1); for i = 1:4 X(i) = []; end</code> | 1-by-0 | 0-by-1 |

Growing Variable-Size Column Cell Array That is Initialized as Scalar at Run Time

In MATLAB execution, if you grow a scalar cell array by using `{end+1}` indexing, the cell array grows along the second dimension and produces a row cell array. For example, define the function `growCell`:

```
function z = growCell(n, m)
for i = 1:m
    n{end+1} = m;
end
z = n;
end
```

Call `growCell` with example inputs:

```
growCell({2}, 3)
```

```
ans =
```

```
1×4 cell array
```

```
{[2]} {[3]} {[3]} {[3]}
```

By contrast, in code generation, suppose that:

- You specify the cell array to be of variable-size column type (for example, `:Inf x 1`) at compile time, *and*
- Initialize this cell array as a scalar at run time.

In such situations, the generated code grows the scalar cell array along the first dimension and produces a column cell array. For example, generate MEX code for `growCell`. Specify the input `n` to be a `:Inf x 1` cell array with double as the underlying type. Specify the input `m` to be of double scalar type.

```
codegen growCell -args {coder.typeof({0}, [Inf 1], [1 0]), 1}
```

Code generation successful.

Run the generated MEX with the same inputs as before.

```
growCell_mex({2}, 3)
```

```
ans =
```

```
4x1 cell array
```

```
{[2]}
{[3]}
{[3]}
{[3]}
```

Binary Element-Wise Operations with Single and Double Operands

If your MATLAB code contains a binary element-wise operation that involves a single type operand and a double type operand, the generated code might not produce the same result as MATLAB.

For such an operation, MATLAB casts both operands to double type and performs the operation with the double types. MATLAB then casts the result to single type and returns it.

The generated code casts the double type operand to single type. It then performs the operation with the two single types and returns the result.

For example, define a MATLAB function `foo` that calls the binary element-wise operation `plus`.

```
function out = foo(a,b)
out = a + b;
end
```

Define a variable `s1` of single type and a variable `v1` of double type. Generate a MEX function for `foo` that accepts a single type input and a double type input.

```
s1 = single(1.4e32);
d1 = -5.305e+32;
codegen foo -args {s1, d1}
```

Call both `foo` and `foo_mex` with inputs `s1` and `d1`. Compare the two results.

```
m1 = foo(s1,d1);
m1c = foo_mex(s1,d1);
m1 == m1c
```

```
ans =  
  
    logical  
  
    0
```

The output of the comparison is a logical 0, which indicates that the generated code and MATLAB produces different results for these inputs.

Floating-Point Numerical Results

The generated code might not produce the same floating-point numerical results as MATLAB in these:

When computer hardware uses extended precision registers

Results vary depending on how the C/C++ compiler allocates extended precision floating-point registers. Computation results might not match MATLAB calculations because of different compiler optimization settings or different code surrounding the floating-point calculations.

For certain advanced library functions

The generated code might use different algorithms to implement certain advanced library functions, such as `fft`, `svd`, `eig`, `mldivide`, and `mrdivide`.

For example, the generated code uses a simpler algorithm to implement `svd` to accommodate a smaller footprint. Results might also vary according to matrix properties. For example, MATLAB might detect symmetric or Hermitian matrices at run time and switch to specialized algorithms that perform computations faster than implementations in the generated code.

For implementation of BLAS library functions

For implementations of BLAS library functions, generated C/C++ code uses reference implementations of BLAS functions. These reference implementations might produce different results from platform-specific BLAS implementations in MATLAB.

NaN and Infinity

The generated code might not produce exactly the same pattern of NaN and Inf values as MATLAB code when these values are mathematically meaningless. For example, if MATLAB output contains a NaN, output from the generated code should also contain a NaN, but not necessarily in the same place.

The bit pattern for NaN can differ between MATLAB code output and generated code output because the C99 language standard that is used to generate code does not specify a unique bit pattern for NaN across all implementations. Avoid comparing bit patterns across different implementations, for example, between MATLAB output and SIL or PIL output.

Negative Zero

In a floating-point type, the value 0 has either a positive sign or a negative sign. Arithmetically, 0 is equal to -0, but some operations are sensitive to the sign of a 0 input. Examples include `rdivide`, `atan2`, `atan2d`, and `angle`. Division by 0 produces Inf, but division by -0 produces -Inf. Similarly, `atan2d(0, -1)` produces 180, but `atan2d(-0, -1)` produces -180.

If the code generator detects that a floating-point variable takes only integer values of a suitable range, then the code generator can use an integer type for the variable in the generated code. If the code generator uses an integer type for the variable, then the variable stores -0 as $+0$ because an integer type does not store a sign for the value 0 . If the generated code casts the variable back to a floating-point type, the sign of 0 is positive. Division by 0 produces Inf , not $-\text{Inf}$. Similarly, `atan2d(0, -1)` produces 180 , not -180 .

There are other contexts in which the generated code might treat -0 differently than MATLAB. For example, suppose that your MATLAB code computes the minimum of two scalar doubles x and y by using `z = min(x, y)`. The corresponding line in the generated C code might be `z = fmin(x, y)`. The function `fmin` is defined in the runtime math library of the C compiler. Because the comparison operation `0.0 == -0.0` returns `true` in C/C++, the compiler's implementation of `fmin` might return either `0.0` or `-0.0` for `fmin(0.0, -0.0)`.

Code Generation Target

The `coder.target` function returns different values in MATLAB than in the generated code. The intent is to help you determine whether your function is executing in MATLAB or has been compiled for a simulation or code generation target. See `coder.target`.

MATLAB Class Property Initialization

Before code generation, at class loading time, MATLAB computes class default values. The code generator uses the values that MATLAB computes. It does not recompute default values. If the property definition uses a function call to compute the initial value, the code generator does not execute this function. If the function has side effects such as modifying a global variable or a persistent variable, then it is possible that the generated code might produce different results than MATLAB. For more information, see “Defining Class Properties for Code Generation” on page 15-3.

MATLAB Classes in Nested Property Assignments That Have Set Methods

When you assign a value to a handle object property, which is itself a property of another object, and so on, then the generated code can call set methods for handle classes that MATLAB does not call.

For example, suppose that you define a set of variables such that x is a handle object, pa is an object, pb is a handle object, and pc is a property of pb . Then you make a nested property assignment, such as:

```
x.pa.pb.pc = 0;
```

In this case, the generated code calls the set method for the object pb and the set method for x . MATLAB calls only the set method for pb .

MATLAB Handle Class Destructors

The behavior of handle class destructors in the generated code can be different from the behavior in MATLAB in these situations:

- The order of destruction of several independent objects might be different in MATLAB than in the generated code.

- The lifetime of objects in the generated code can be different from their lifetime in MATLAB.
- The generated code does not destroy partially constructed objects. If a handle object is not fully constructed at run time, the generated code produces an error message but does not call the `delete` method for that object. For a System object, if there is a run-time error in `setupImpl`, the generated code does not call `releaseImpl` for that object.

MATLAB does call the `delete` method to destroy a partially constructed object.

For more information, see “Code Generation for Handle Class Destructors” on page 15-15.

Variable-Size Data

See “Incompatibilities with MATLAB in Variable-Size Support for Code Generation” on page 29-18.

Complex Numbers

See “Code Generation for Complex Data” on page 16-8.

Converting Strings with Consecutive Unary Operators to double

Converting a string that contains multiple, consecutive unary operators to `double` can produce different results between MATLAB and the generated code. Consider this function:

```
function out = foo(op)
out = double(op + 1);
end
```

For an input value `--`, the function converts the string `--1` to `double`. In MATLAB, the answer is `NaN`. In the generated code, the answer is `1`.

Display Function

Statements and expressions in MATLAB code that omit the semicolon implicitly invoke the `display` function. You can also explicitly invoke `display` as shown here:

```
display(2+3);

5
```

The MEX code generated for MATLAB code that invokes the `display` function preserves calls to this function and shows the output. In standalone code generated for targets that do not have access to MATLAB Runtime, implicit and explicit calls to `display` are removed. This includes calls to overridden class methods of `display`.

To display text in code generated for other targets, override the `disp` function in your MATLAB classes. For example:

```
%MATLAB Class

classdef foo
    methods
        function obj = foo
    end
```

```
        function disp(self)
            disp("Overridden disp");
        end
    end
end
```

%Entry-point Function

```
function callDisp
a = foo;
disp(a);
end
```

The generated code for the entry-point function is shown here:

```
/* Include Files */
#include "callDisp.h"
#include <stdio.h>

/* Function Definitions */
/*
 * Arguments      : void
 * Return Type    : void
 */
void callDisp(void)
{
    printf("%s\n", "Overridden disp");
    fflush(stdout);
}
```

Function Handle Difference

Invoking `display` through a function handle in MATLAB prints the name of the variable as well. For example, running this function in MATLAB results in the following output:

```
function displayDiff
z = 10;
f = @display;
f(z)
end
```

```
z =
    10
```

However, the generated code for this snippet only outputs the value 10.

See Also

More About

- “Potential Differences Reporting” on page 19-23
- “Potential Differences Messages” on page 19-25

Potential Differences Reporting

Generation of efficient C/C++ code from MATLAB code sometimes results in behavior differences between the generated code and the MATLAB code on page 19-6. When you run your program, run-time error checks can detect some of these differences. By default, run-time error checks are:

- Enabled for MEX code generated by using `codegen`, `fiaccl`, or the MATLAB Coder app.
- Disabled for standalone C/C++ code generated by using `codegen` or the MATLAB Coder app.

To help you identify and address differences before you deploy code, the code generator reports a subset of the differences as potential differences. A potential difference is a difference that occurs at run time only under certain conditions.

Addressing Potential Differences Messages

If the code generator detects a potential difference, it displays a message for the difference on the **Potential Differences** tab of the report. If you use the MATLAB Coder app to generate code, you can view the message in the **Potential Differences** tab of the app itself. To highlight the MATLAB code that corresponds to the message, click the message.

The presence of a potential difference message does not necessarily mean that the difference will occur when you run the generated code. To determine whether the potential difference affects your application:

- Analyze the behavior of your MATLAB code for the range of data for your application.
- Test a MEX function generated from your MATLAB code. Use the range of data that your application uses. If the difference occurs, the MEX function reports an error.

If your analysis or testing confirms the reported difference, consider modifying your code. Some potential differences messages provide a workaround. For additional information about some of the potential differences messages, see “Potential Differences Messages” on page 19-25. Even if you modify your code to prevent a difference from occurring at run time, the code generator might still report the potential difference.

The set of potential differences that the code generator detects is a subset of the differences that MEX functions report as errors. It is a best practice to test a MEX function over the full range of application data.

Disabling and Enabling Potential Differences Reporting for MATLAB Coder

By default, potential differences reporting is enabled for:

- Code generation with the `codegen` command
- The **Check for Run-Time Issues** step in the MATLAB Coder app

To disable potential differences reporting:

- In a code configuration object, set `ReportPotentialDifferences` to `false`.
- In the MATLAB Coder app, in the **Debugging** settings, clear the **Report differences from MATLAB** check box.

By default, potential differences reporting is disabled for the **Generate code** step and the code generation report in the MATLAB Coder app. To enable potential differences reporting, in the **Debugging** settings, select the **Report differences from MATLAB** check box.

Disabling and Enabling Potential Differences Reporting for Fixed-Point Designer

By default, potential differences reporting is enabled for code acceleration with `fiaccel`. To disable it, in a code acceleration configuration object, set `ReportPotentialDifferences` to `false`.

See Also

More About

- “Potential Differences Messages” on page 19-25
- “Incompatibilities with MATLAB in Variable-Size Support for Code Generation” on page 29-18
- “Differences Between Generated Code and MATLAB Code” on page 19-6

Potential Differences Messages

When you enable potential differences on page 19-23 reporting, the code generator reports potential differences between the behavior of the generated code and the behavior of the MATLAB code. Reviewing and addressing potential differences before you generate standalone code helps you to avoid errors and incorrect answers in generated code.

Here are some of the potential differences messages:

- “Automatic Dimension Incompatibility” on page 19-25
- “mtimes No Dynamic Scalar Expansion” on page 19-25
- “Matrix-Matrix Indexing” on page 19-26
- “Vector-Vector Indexing” on page 19-26
- “Loop Index Overflow” (MATLAB Coder)

Automatic Dimension Incompatibility

In the generated code, the dimension to operate along is selected automatically, and might be different from MATLAB. Consider specifying the working dimension explicitly as a constant value.

This restriction applies to functions that take the working dimension (the dimension along which to operate) as input. In MATLAB and in code generation, if you do not supply the working dimension, the function selects it. In MATLAB, the function selects the first dimension whose size does not equal 1. For code generation, the function selects the first dimension that has a variable size or that has a fixed size that does not equal 1. If the working dimension has a variable size and it becomes 1 at run time, then the working dimension is different from the working dimension in MATLAB. Therefore, when run-time error checks are enabled, an error can occur.

For example, suppose that X is a variable-size matrix with dimensions $1 \times 3 \times 5$. In the generated code, `sum(X)` behaves like `sum(X,2)`. In MATLAB, `sum(X)` behaves like `sum(X,2)` unless `size(X,2)` is 1. In MATLAB, when `size(X,2)` is 1, `sum(X)` behaves like `sum(X,3)`.

To avoid this issue, specify the intended working dimension explicitly as a constant value. For example, `sum(X,2)`.

mtimes No Dynamic Scalar Expansion

The generated code performs a general matrix multiplication. If a variable-size matrix operand becomes a scalar at run time, dimensions must still agree. There will not be an automatic switch to scalar multiplication.

Consider the multiplication $A*B$. If the code generator is aware that A is scalar and B is a matrix, the code generator produces code for scalar-matrix multiplication. However, if the code generator is aware that A and B are variable-size matrices, it produces code for a general matrix multiplication. At run time, if A turns out to be scalar, the generated code does not change its behavior. Therefore, when run-time error checks are enabled, a size mismatch error can occur.

Matrix-Matrix Indexing

For indexing a matrix with a matrix, `matrix1(matrix2)`, the code generator assumed that the result would have the same size as `matrix2`. If `matrix1` and `matrix2` are vectors at run time, their orientations must match.

In matrix-matrix indexing, you use one matrix to index into another matrix. In MATLAB, the general rule for matrix-matrix indexing is that the size and orientation of the result match the size and orientation of the index matrix. For example, if `A` and `B` are matrices, `size(A(B))` equals `size(B)`. When `A` and `B` are vectors, MATLAB applies a special rule. The special vector-vector indexing rule is that the orientation of the result is the orientation of the data matrix. For example, if `A` is 1-by-5 and `B` is 3-by-1, then `A(B)` is 1-by-3.

The code generator applies the same matrix-matrix indexing rules as MATLAB. If `A` and `B` are variable-size matrices, to apply the matrix-matrix indexing rules, the code generator assumes that `size(A(B))` equals `size(B)`. If, at run time, `A` and `B` become vectors and have different orientations, then the assumption is incorrect. Therefore, when run-time error checks are enabled, an error can occur.

To avoid this issue, force your data to be a vector by using the colon operator for indexing. For example, suppose that your code intentionally toggles between vectors and regular matrices at run time. You can do an explicit check for vector-vector indexing.

```
...
if isvector(A) && isvector(B)
    C = A(:);
    D = C(B(:));
else
    D = A(B);
end
...
```

The indexing in the first branch specifies that `C` and `B(:)` are compile-time vectors. Therefore, the code generator applies the indexing rule for indexing one vector with another vector. The orientation of the result is the orientation of the data vector, `C`.

Vector-Vector Indexing

For indexing a vector with a vector, `vector1(vector2)`, the code generator assumed that the result would have the same orientation as `vector1`. If `vector1` is a scalar at run time, the orientation of `vector2` must match `vector1`.

In MATLAB, the special rule for vector-vector indexing is that the orientation of the result is the orientation of the data vector. For example, if `A` is 1-by-5 and `B` is 3-by-1, then `A(B)` is 1-by-3. If, however, the data vector `A` is a scalar, then the orientation of `A(B)` is the orientation of the index vector `B`.

The code generator applies the same vector-vector indexing rules as MATLAB. If `A` and `B` are variable-size vectors, to apply the indexing rules, the code generator assumes that the orientation of `B` matches the orientation of `A`. At run time, if `A` is scalar and the orientation of `A` and `B` do not match,

then the assumption is incorrect. Therefore, when run-time error checks are enabled, a run-time error can occur.

To avoid this issue, make the orientations of the vectors match. Alternatively, index single elements by specifying the row and column. For example, `A(row, column)`.

Loop Index Overflow

The generated code assumes the loop index does not overflow on the last iteration of the loop. If the loop index overflows, an infinite loop can occur.

Suppose that a `for`-loop end value is equal to or close to the maximum or minimum value for the loop index data type. In the generated code, the last increment or decrement of the loop index might cause the index variable to overflow. The index overflow might result in an infinite loop.

When memory integrity checks are enabled, if the code generator detects that the loop index might overflow, it reports an error. The software error checking is conservative. It might incorrectly report a loop index overflow. By default, memory-integrity checks are enabled for MEX code and disabled for standalone C/C++ code. See “Why Test MEX Functions in MATLAB?” (MATLAB Coder) and “Generate Standalone C/C++ Code That Detects and Reports Run-Time Errors” (MATLAB Coder).

To avoid a loop index overflow, use the workarounds in this table.

| Loop Conditions Causing the Potential Overflow | Workaround |
|--|--|
| <ul style="list-style-type: none"> The loop index increments by 1. The end value equals the maximum value of the integer type. | <p>If the loop does not have to cover the full range of the integer type, rewrite the loop so that the end value is not equal to the maximum value of the integer type. For example, replace:</p> <pre>N=intmax('int16') for k=N-10:N</pre> <p>with:</p> <pre>for k=1:10</pre> |
| <ul style="list-style-type: none"> The loop index decrements by 1. The end value equals the minimum value of the integer type. | <p>If the loop does not have to cover the full range of the integer type, rewrite the loop so that the end value is not equal to the minimum value of the integer type. For example, replace:</p> <pre>N=intmin('int32') for k=N+10:-1:N</pre> <p>with:</p> <pre>for k=10:-1:1</pre> |

| Loop Conditions Causing the Potential Overflow | Workaround |
|---|--|
| <ul style="list-style-type: none"> • The loop index increments or decrements by 1. • The start value equals the minimum or maximum value of the integer type. • The end value equals the maximum or minimum value of the integer type. | <p>If the loop must cover the full range of the integer type, cast the type of the loop start, step, and end values to a bigger integer or to double. For example, rewrite:</p> <pre>M= intmin('int16'); N= intmax('int16'); for k=M:N % Loop body end</pre> <p>as:</p> <pre>M= intmin('int16'); N= intmax('int16'); for k=int32(M):int32(N) % Loop body end</pre> |
| <ul style="list-style-type: none"> • The loop index increments or decrements by a value not equal to 1. • On the last loop iteration, the loop index is not equal to the end value. | <p>Rewrite the loop so that the loop index in the last loop iteration is equal to the end value.</p> |

See Also

More About

- “Potential Differences Reporting” on page 19-23
- “Differences Between Generated Code and MATLAB Code” on page 19-6
- “Incompatibilities with MATLAB in Variable-Size Support for Code Generation” on page 29-18

MATLAB Language Features Supported for C/C++ Code Generation

MATLAB Features That Code Generation Supports

Code generation from MATLAB code supports many major language features including:

- n-dimensional arrays (see “Array Size Restrictions for Code Generation” on page 16-13)
- matrix operations, including deletion of rows and columns
- variable-size data (see “Code Generation for Variable-Size Arrays” on page 29-2)
- subscripting (see “Incompatibility with MATLAB in Matrix Indexing Operations for Code Generation” on page 29-22)
- complex numbers (see “Code Generation for Complex Data” on page 16-8)
- numeric classes (see “Supported Variable Types” on page 18-13)
- double-precision, single-precision, and integer math
- enumerations (see “Code Generation for Enumerations” on page 20-2)
- fixed-point arithmetic (see “Code Acceleration and Code Generation from MATLAB” on page 12-2)
- program control statements `if`, `switch`, `for`, `while`, and `break`
- arithmetic, relational, and logical operators
- local functions
- persistent variables
- global variables
- structures (see “Structure Definition for Code Generation” on page 25-2)
- cell arrays (see “Cell Arrays”)
- tables (see “Code Generation for Tables” on page 27-2)
- timetables (see “Code Generation for Timetables” on page 28-2)
- characters (see “Encoding of Characters in Code Generation” on page 16-12)
- string scalars (see “Code Generation for Strings” on page 16-16)
- categorical arrays (see “Code Generation for Categorical Arrays” on page 21-2)
- `datetime` arrays (see “Code Generation for Datetime Arrays” on page 22-2)
- `duration` arrays (see “Code Generation for Duration Arrays” on page 23-2)
- sparse matrices (see “Code Generation for Sparse Matrices” on page 16-19)
- function handles (see “Function Handle Limitations for Code Generation” on page 24-2)
- anonymous functions (see “Code Generation for Anonymous Functions” on page 17-9)
- recursive functions (see “Code Generation for Recursive Functions” on page 14-12)
- nested functions (see “Code Generation for Nested Functions” on page 17-10)
- variable length input and output argument lists (see “Code Generation for Variable Length Argument Lists” on page 17-2)
- function argument validation (see “Generate Code for arguments Block That Validates Input and Output Arguments” on page 17-3)

- subset of MATLAB toolbox functions (see “Functions and Objects Supported for C/C++ Code Generation” on page 26-2)
- subset of functions and System objects in several toolboxes (see “Functions and Objects Supported for C/C++ Code Generation” on page 26-2)
- function calls (see “Resolution of Function Calls for Code Generation” on page 14-2)
- class aliasing
- MATLAB classes (see “MATLAB Classes Definition for Code Generation” on page 15-2)

MATLAB Language Features That Code Generation Does Not Support

Code generation from MATLAB does not support the following frequently used MATLAB features (this list is not exhaustive):

- scripts
- GPU arrays

MATLAB Coder does not support GPU arrays. However, if you have GPU Coder™, you can generate CUDA® MEX code that takes GPU array inputs.

- `calendarDuration` arrays
- Java
- Map containers
- time series objects
- tall arrays
- `try/catch` statements
- `import` statements
- pattern arrays

Code Generation for Enumerated Data

- “Code Generation for Enumerations” on page 20-2
- “Customize Enumerated Types in Generated Code” on page 20-7

Code Generation for Enumerations

Enumerations represent a fixed set of named values. Enumerations help make your MATLAB code more readable.

For code generation, when you use enumerations, adhere to these restrictions:

- Calls to methods of enumeration classes are not supported.
- Passing strings or character vectors to constructors of enumerations is not supported.
- The enumeration class must derive from one of these base types: `int8`, `uint8`, `int16`, `uint16`, `int32`, or `uint32`. See “Define Enumerations for Code Generation” on page 20-2.
- For members of `uint32` enumerations, code generation supports values that are less than or equal to `intmax("int32")`.
- You can use only a limited set of operations on enumerations. See “Allowed Operations on Enumerations” on page 20-4.
- Use enumerations with functions that support enumerated types for code generation. See “MATLAB Toolbox Functions That Support Enumerations” on page 20-5.

Define Enumerations for Code Generation

For code generation, the enumeration class must derive from one of these base types: `int8`, `uint8`, `int16`, `uint16`, `int32`, or `uint32`. For example:

```
classdef PrimaryColors < int32
    enumeration
        Red(1),
        Blue(2),
        Yellow(4)
    end
end
```

If you use MATLAB Coder to generate C/C++ code, you can use the base type to control the size of an enumerated type in the generated code. You can:

- Represent an enumerated type as a fixed-size integer that is portable to different targets.
- Reduce memory usage.
- Interface with legacy code.
- Match company standards.

Representation of Enumerated Type in Generated Code

The representation of the enumerated type in generated C/C++ code depends on the following:

- The base type of the MATLAB enumeration
- The target language (C or C++)
- If the target language is C++, the target language standard (C++03 or C++11)

Base Type is Native Integer Type

If the base type is the native integer type for the target platform (for example, `int32`), the code generator produces a C/C++ enumerated type. Consider this MATLAB enumerated type definition:

```

classdef LEDcolor < int32
    enumeration
        GREEN(1),
        RED(2)
    end
end

```

If you generate C code or C++03 code, the generated enumeration is:

```

enum LEDcolor
{
    GREEN = 1,
    RED
};

```

If you generate C++11 code, the generated code contains an enumeration class (by default) that explicitly defines the underlying type:

```

enum class LEDcolor : int
{
    GREEN = 1,
    RED
};

```

Base Type is Different from the Native Integer Type

Suppose that built-in integer base type for the enumeration is different from the native integer type for the target platform. For example, consider this MATLAB enumerated type definition:

```

classdef LEDcolor < int16
    enumeration
        GREEN(1),
        RED(2)
    end
end

```

- If you generate C code, the code generator produces a `typedef` statement for the enumerated type and `#define` statements for the enumerated values. For example, the enumerated type definition `LEDcolor` produces this C code:

```

typedef short LEDcolor;
#define GREEN ((LEDcolor)1)
#define RED ((LEDcolor)2)

```

- If you generate C++03 code, the enumeration members are converted to constants. These constants belong to the namespace that contains the enumeration type definition in the generated C++ code.

For example, suppose that you place the enumerated type definition `LEDcolor` inside the package `pkg`. The default behavior of the code generator is to convert MATLAB packages to C++ namespaces. The generated C++ code is placed inside the namespace `pkg`:

```

namespace pkg {
    typedef short LEDcolor;

    // enum pkg_LEDcolor
    const LEDcolor GREEN{1};
    const LEDcolor RED{2};
}

```

```
}

```

- C++11 allows you to specify the underlying type of an enumeration, just like MATLAB does. If you generate C++11 code, the MATLAB enumeration class is converted to a C++ enumeration class (by default) that explicitly defines the underlying type.

For example, suppose that you place the enumerated type definition `LEDcolor` inside the package `pkg`. The default behavior of the code generator is to convert MATLAB packages to C++ namespaces. The generated C++11 code is placed inside the namespace `pkg`:

```
namespace pkg {
enum class LEDcolor : short
{
    GREEN = 1, // Default value
    RED
};
}

```

The C/C++ type in the `typedef` statement or the underlying type of the C++11 enumeration depends on:

- The integer sizes defined for the production hardware in the hardware implementation object or the project settings. See `coder.HardwareImplementation`.
- The setting that determines the use of built-in C types or MathWorks typedefs in the generated code. See “Specify Data Types Used in Generated Code” (MATLAB Coder) and “Mapping MATLAB Types to Types in Generated Code” (MATLAB Coder).

Generate C++11 Code That Contains Ordinary C Enumerations

You can change the default behavior of the code generator to produce ordinary C enumerations in the generated C++11 code. Do one of the following:

- In the code generation configuration object, set the `CppGenerateEnumClass` property to `false`.
- In the MATLAB Coder app, in the **Generate** step, on the **Code Appearance** tab, clear the **Generate C++ enum class from MATLAB enumeration** check box.

To instruct the code generator to produce ordinary C enumeration for a particular MATLAB enumeration class in your code, include the static method `generateEnumClass` that returns `false` in the implementation of that MATLAB enumeration class. See “Customize Enumerated Types in Generated Code” on page 20-7.

Allowed Operations on Enumerations

For code generation, you are restricted to the operations on enumerations listed in this table.

| Operation | Example | Notes |
|------------------------|---------|-------|
| assignment operator: = | | — |

| Operation | Example | Notes |
|--|--|--|
| relational operators: < > <= >= == ~= | <code>xon == xoff</code> | Code generation does not support using == or ~= to test equality between an enumeration member and a string array, a character array, or a cell array of character arrays. |
| cast operation | <code>double(LEDcolor.RED)</code> | — |
| conversion to character array or string | <code>y = char(LEDcolor.RED); y1 = cast(LEDcolor.RED, 'char'); y2 = string(LEDcolor.RED);</code> | <ul style="list-style-type: none"> You can convert only compile-time scalar valued enumerations. For example, this code runs in MATLAB, but produces an error in code generation: <code>y2 = string(repmat(LEDcolor.RED,1,2));</code> The code generator preserves enumeration names when the conversion inputs are constants. For example, consider this enumerated type definition: <pre>classdef AnEnum < int32 enumeration zero(0), two(2), otherTwo(2) end end</pre>Generated code produces "two" for <code>y = string(AnEnum.two)</code> and "otherTwo" for <code>y = string(AnEnum.two)</code> |
| indexing operation | <code>m = [1 2] n = LEDcolor(m) p = n(LEDcolor.GREEN)</code> | — |
| control flow statements: if, switch, while | <code>if state == sysMode.ON led = LEDcolor.GREEN; else led = LEDcolor.RED; end</code> | — |

MATLAB Toolbox Functions That Support Enumerations

For code generation, you can use enumerations with these MATLAB toolbox functions:

- `cast`
- `cat`
- `char`
- `circshift`

- enumeration
- fliplr
- flipud
- histc
- intersect
- ipermute
- isequal
- isequaln
- isfinite
- isinf
- ismember
- isnan
- issorted
- length
- permute
- repmat
- reshape
- rot90
- setdiff
- setxor
- shiftdim
- sort
- sortrows
- squeeze
- string
- union
- unique

See Also

Related Examples

- “Customize Enumerated Types in Generated Code” on page 20-7

Customize Enumerated Types in Generated Code

For code generation, to customize an enumeration, in the static methods section of the class definition, include customized versions of the methods listed in this table.

| Method | Description | Default Value Returned or Specified | When to Use |
|--------------------------------------|--|--|---|
| <code>getDefaultValue</code> | Returns the default enumerated value. | First value in the enumeration class definition. | For a default value that is different than the first enumeration value, provide a <code>getDefaultValue</code> method that returns the default value that you want. See “Specify a Default Enumeration Value” on page 20-8. |
| <code>getHeaderFile</code> | Specifies the file that defines an externally defined enumerated type. | <code>''</code> | To use an externally defined enumerated type, provide a <code>getHeaderFile</code> method that returns the path to the header file that defines the type. In this case, the code generator does not produce the class definition. See “Specify a Header File” on page 20-8. |
| <code>addClassNameToEnumNames</code> | Specifies whether the class name becomes a prefix in the generated code. | <code>false</code> — prefix is not used. | <p>If you want the class name to become a prefix in the generated code, set the return value of the <code>addClassNameToEnumNames</code> method to <code>true</code>. See “Include Class Name Prefix in Generated Enumerated Type Value Names” on page 20-9.</p> <p>Note When generating C++11 enumeration classes, the code generator ignores this static method.</p> |

| Method | Description | Default Value Returned or Specified | When to Use |
|-------------------|---|--|--|
| generateEnumClass | Specifies whether to generate C++11 enumeration classes | true — enumeration classes are generated in C++11 code | When generating C++11 code, to instruct the code generator to produce ordinary C enumeration for a particular MATLAB enumeration, set the return value of generateEnumClass method to false. See “Generate C++11 Code Containing Ordinary C Enumeration” (MATLAB Coder). |

Specify a Default Enumeration Value

If the value of a variable that is cast to an enumerated type does not match one of the enumerated type values:

- Generated MEX reports an error.
- Generated C/C++ code replaces the value of the variable with the enumerated type default value.

Unless you specify otherwise, the default value for an enumerated type is the first value in the enumeration class definition. To specify a different default value, add your own `getDefaultValue` method to the methods section. In this example, the first enumeration member value is `LEDcolor.GREEN`, but the `getDefaultValue` method returns `LEDcolor.RED`:

```
classdef LEDcolor < int32
    enumeration
        GREEN(1),
        RED(2)
    end

    methods (Static)
        function y = getDefaultValue()
            y = LEDcolor.RED;
        end
    end
end
```

Specify a Header File

To specify that an enumerated type is defined in an external file, provide a customized `getHeaderFile` method. This example specifies that `LEDcolor` is defined in the external file `my_LEDcolor.h`.

```
classdef LEDcolor < int32
    enumeration
        GREEN(1),
        RED(2)
    end
end
```

```

end

methods(Static)
function y=getHeaderFile()
    y='my_LEDcolor.h';
end
end
end

```

You must provide `my_LEDcolor.h`. For example:

```

enum LEDcolor
{
    GREEN = 1,
    RED
};
typedef enum LEDcolor LEDcolor;

```

If you place the MATLAB enumeration `LEDcolor` inside the package `pkg` and generate C++ code, code generation preserves the name of this enumeration and places it inside the namespace `pkg` in the generated code. Therefore, in the header file that you provide, you must define this enumeration inside the namespace `pkg`.

Include Class Name Prefix in Generated Enumerated Type Value Names

By default, the generated enumerated type value name does not include the class name prefix. For example:

```

enum LEDcolor
{
    GREEN = 1,
    RED
};

typedef enum LEDcolor LEDcolor;

```

To include the class name prefix, provide an `addClassNameToEnumNames` method that returns `true`. For example:

```

classdef LEDcolor < int32
    enumeration
        GREEN(1),
        RED(2)
    end

    methods(Static)
        function y = addClassNameToEnumNames()
            y=true;
        end
    end
end

```

In the generated type definition, the enumerated value names include the class prefix `LEDcolor`.

```

enum LEDcolor
{

```

```
        LEDcolor_GREEN = 1,  
        LEDcolor_RED  
};  
  
typedef enum LEDcolor LEDcolor;
```

Generate C++11 Code Containing Ordinary C Enumeration

When you generate C++11 code, your MATLAB enumeration class is converted to a C++11 enumeration class. For example:

```
enum class MyEnumClass16 : short  
{  
    Orange = 0, // Default value  
    Yellow,  
    Pink  
};
```

To generate an ordinary C enumeration instead, provide a `generateEnumClass` method that returns `false`. For example:

```
classdef MyEnumClass16 < int16  
    enumeration  
        Orange(0),  
        Yellow(1),  
        Pink(2)  
    end  
  
    % particular enum opting out  
    methods(Static)  
        function y = generateEnumClass()  
            y = false;  
        end  
    end  
end
```

Now the generated C++11 code contains an ordinary C enumeration.

```
enum MyEnumClass16 : short  
{  
    Orange = 0, // Default value  
    Yellow,  
    Pink  
};
```

See Also

More About

- [Modifying Superclass Methods and Properties](#)
- [“Code Generation for Enumerations” on page 20-2](#)

Code Generation for Categorical Arrays

Code Generation for Categorical Arrays

In this section...

“Define Categorical Arrays for Code Generation” on page 21-2

“Allowed Operations on Categorical Arrays” on page 21-2

“MATLAB Toolbox Functions That Support Categorical Arrays” on page 21-3

Categorical arrays store data with values from a finite set of discrete categories. You can specify an order for the categories, but it is not required. A categorical array provides efficient storage and manipulation of nonnumeric data, while also maintaining meaningful names for the values.

When you use categorical arrays with code generation, adhere to these restrictions:

Define Categorical Arrays for Code Generation

For code generation, use the `categorical` function to create categorical arrays. For example, suppose the input argument to your MATLAB function is a numeric array of arbitrary size whose elements have values of either 1, 2, or 3. You can convert these values to the categories `small`, `medium`, and `large` and turn the input array into a categorical array, as shown in this code.

```
function c = foo(x) %#codegen
    c = categorical(x,1:3,{'small','medium','large'});
end
```

Allowed Operations on Categorical Arrays

For code generation, you are restricted to the operations on categorical arrays listed in this table.

| Operation | Example | Notes |
|---------------------------------------|---|---|
| assignment operator: = | <pre>c = categorical(1:3,1:3,{'small','medium','large'}); c(1) = 'large';</pre> | Code generation does not support using the assignment operator = to: <ul style="list-style-type: none"> Delete an element. Expand the size of a categorical array. Add a new category, even when the array is not protected. |
| relational operators: < > <= >= == ~= | <pre>c = categorical(1:3,'Ordinal',true); tf = c(1) < c(2);</pre> | Code generation supports all relational operators. |
| cast to numeric type | <pre>c = categorical(1:3); double(c(1));</pre> | Code generation supports casting categorical arrays to arrays of double- or single-precision floating-point numbers, or to integers. |

| Operation | Example | Notes |
|--------------------|---|---|
| conversion to text | <pre>c = categorical(1:3,1:3,{'small','medium','large'}); c1 = cellstr(c(1)); % One element c2 = cellstr(c); % Entire array</pre> | <p>Code generation does not support using the <code>char</code> or <code>string</code> functions to convert categorical values to text.</p> <p>To convert one or more elements of a categorical array to text, use the <code>cellstr</code> function.</p> |
| indexing operation | <pre>c = categorical(1:3,1:3,{'small','medium','large'}); idx = [1 2]; c(idx); idx = logical([1 1 0]); c(idx);</pre> | Code generation supports indexing by position, linear indexing, and logical indexing. |
| concatenation | <pre>c1 = categorical(1:3,1:3,{'small','medium','large'}); c2 = categorical(4:6,[2 1 4],{'medium','small','extra-large'}); c = [c1 c2];</pre> | Code generation supports concatenation of categorical arrays along any dimension. |

MATLAB Toolbox Functions That Support Categorical Arrays

For code generation, you can use categorical arrays with these MATLAB toolbox functions:

- `addcats`
- `cat`
- `categorical`
- `categories`
- `cellstr`
- `countcats`
- `ctranspose`
- `double`
- `eq`
- `ge`
- `gt`
- `histcounts`
- `horzcat`
- `int8`
- `int16`
- `int32`
- `int64`
- `intersect`
- `iscategory`
- `iscolumn`
- `isempty`
- `isequal`

- `isequaln`
- `ismatrix`
- `ismember`
- `isordinal`
- `isprotected`
- `isrow`
- `isscalar`
- `issorted`
- `issortedrows`
- `isundefined`
- `isvector`
- `le`
- `length`
- `lt`
- `max`
- `mergocats`
- `min`
- `ndims`
- `ne`
- `numel`
- `permute`
- `removecats`
- `renamecats`
- `reordercats`
- `reshape`
- `setcats`
- `setdiff`
- `setxor`
- `single`
- `size`
- `sort`
- `sortrows`
- `transpose`
- `uint8`
- `uint16`
- `uint32`
- `uint64`
- `union`
- `unique`

- `vertcat`

See Also

More About

- “Define Categorical Array Inputs” on page 21-6
- “Categorical Array Limitations for Code Generation” on page 21-8

Define Categorical Array Inputs

You can define categorical array inputs at the command line. Programmatic specification of categorical input types by using preconditioning (assert statements) is not supported.

Define Categorical Array Inputs at the Command Line

Use one of these procedures:

- “Provide an Example Categorical Array Input” on page 21-6
- “Provide a Categorical Array Type” on page 21-6
- “Provide a Constant Categorical Array Input” on page 21-6

Alternatively, if you have a test file that calls your entry-point function with example inputs, you can determine the input types by using `coder.getArgTypes`.

Provide an Example Categorical Array Input

Use the `-args` option:

```
C = categorical({'r','g','b'});  
fiaccl myFunction -args {C}
```

Provide a Categorical Array Type

To provide a type for a categorical array to `fiaccl`:

- 1 Define a categorical array. For example:

```
C = categorical({'r','g','b'});
```
- 2 Create a type from `C`.

```
t = coder.typeof(C);
```
- 3 Pass the type to `fiaccl` by using the `-args` option.

```
fiaccl myFunction -args {t}
```

Provide a Constant Categorical Array Input

To specify that a categorical array input is constant, use `coder.Constant` with the `-args` option:

```
C = categorical({'r','g','b'});  
fiaccl myFunction -args {coder.Constant(C)}
```

Representation of Categorical Arrays

A coder type object for a categorical array describes the object and its properties. Use `coder.typeof` or pass `categorical` as a string scalar to `coder.newtype`.

The coder type object displays a succinct description of the object properties while excluding internal state values. Nonconstant properties display their type and size, while constant properties display only their values. For example:

```
t = categorical({'r','g','b'});
tType = coder.typeof(t)
```

The representation of variable `t` is stored in coder type object `tType`.

```
tType =
    matlab.coder.type.CategoricalType
    1x3 categorical
    Categories : 3x1 homogeneous cell
    Ordinal : 1x1 logical
    Protected : 1x1 logical
```

If your workflow requires the legacy representation of coder type objects, use the `getCoderType` function on the variable that has the new representation of your class or object. See “Legacy Representation of Coder Type Objects” (MATLAB Coder).

Resize Object Properties by Using `coder.resize`

You can resize most objects by using `coder.resize`. You can resize objects, its properties and create arrays within the properties.

For a `categorical` coder object, you can resize the object properties:

```
t = categorical({'r','g','b'});
tType = coder.typeof(t);
tType.Categories = coder.resize(tType.Categories, [3 1],[1 0])
```

This code resizes the `Categories` property to be upper-bounded at 3 for the first dimension.

```
tType =
    matlab.coder.type.CategoricalType
    1x3 categorical
    Categories : :3x1 homogeneous cell
    Ordinal : 1x1 logical
    Protected : 1x1 logical
```

You can also resize the object by using `coder.resize`. See “Edit and Represent Coder Type Objects and Properties” (MATLAB Coder).

See Also

`categorical` | `coder.Constant` | `coder.typeof`

More About

- “Code Generation for Categorical Arrays” on page 21-2
- “Categorical Array Limitations for Code Generation” on page 21-8

Categorical Array Limitations for Code Generation

When you create categorical arrays in MATLAB code that you intend for code generation, you must specify the categories and elements of each categorical array by using the `categorical` function. See “Categorical Arrays”.

For categorical arrays, code generation does not support the following inputs and operations:

- Arrays of MATLAB objects.
- Sparse matrices.
- Duplicate category names when you specify them using the `categoryNames` input argument of the `categorical` function.
- Growth by assignment. For example, assigning a value beyond the end of an array produces an error.

```
function c = foo() %#codegen
    c = categorical(1:3,1:3,{'small','medium','large'});
    c(4) = 'medium';
end
```

- Adding a category. For example, specifying a new category by using the `=` operator produces an error, even when the categorical array is unprotected.

```
function c = foo() %#codegen
    c = categorical(1:3,1:3,{'small','medium','large'});
    c(1) = 'extra-large';
end
```

- Deleting an element. For example, assigning an empty array to an element produces an error.

```
function c = foo() %#codegen
    c = categorical(1:3,1:3,{'small','medium','large'});
    c(1) = [];
end
```

- Converting categorical values to text by using the `char` or `string` functions. To convert elements of a categorical array to text, use the `cellstr` function.

Limitations that apply to classes also apply to categorical arrays. For more information, see “MATLAB Classes Definition for Code Generation” (MATLAB Coder).

See Also

`categorical` | `cellstr`

More About

- “Code Generation for Categorical Arrays” on page 21-2
- “Define Categorical Array Inputs” on page 21-6

Code Generation for Datetime Arrays

- “Code Generation for Datetime Arrays” on page 22-2
- “Define Datetime Array Inputs” on page 22-5
- “Datetime Array Limitations for Code Generation” on page 22-7

Code Generation for Datetime Arrays

In this section...

“Define Datetime Arrays for Code Generation” on page 22-2

“Allowed Operations on Datetime Arrays” on page 22-2

“MATLAB Toolbox Functions That Support Datetime Arrays” on page 22-2

The values in a `datetime` array represent points in time using the proleptic ISO calendar.

When you use `datetime` arrays with code generation, adhere to these restrictions.

Define Datetime Arrays for Code Generation

For code generation, use the `datetime` function to create `datetime` arrays. For example, suppose the input arguments to your MATLAB function are numeric arrays whose values indicate the year, month, day, hour, minute, and second components for a point in time. You can create a `datetime` array from these input arrays.

```
function d = foo(y,mo,d,h,mi,s) %#codegen
    d = datetime(y,mo,d,h,mi,s);
end
```

Allowed Operations on Datetime Arrays

For code generation, you are restricted to the operations on `datetime` arrays listed in this table.

| Operation | Example | Notes |
|---------------------------------------|---|--|
| Assignment operator: = | <pre>d = datetime(2019,1:12,1,12,0,0); d(1) = datetime(2019,1,31);</pre> | Code generation does not support using the assignment operator = to: <ul style="list-style-type: none"> Delete an element. Expand the size of a <code>datetime</code> array. |
| Relational operators: < > <= >= == ~= | <pre>d = datetime(2019,1:12,1,12,0,0); tf = d(1) < d(2);</pre> | Code generation supports relational operators. |
| Indexing operation | <pre>d = datetime(2019,1:12,1,12,0,0); idx = [1 2]; d(idx); idx = logical([1 1 0]); d(idx);</pre> | Code generation supports indexing by position, linear indexing, and logical indexing. |
| Concatenation | <pre>d1 = datetime(2019,1:6,1,12,0,0); d2 = datetime(2019,7:12,1,12,0,0); d = [d1 d2];</pre> | Code generation supports concatenation of <code>datetime</code> arrays. |

MATLAB Toolbox Functions That Support Datetime Arrays

For code generation, you can use `datetime` arrays with these MATLAB toolbox functions:

- `cat`

- colon
- ctranspose
- datetime
- datevec
- diff
- eq
- ge
- gt
- hms
- horzcat
- hour
- interp1
- intersect
- iscolumn
- isempty
- isequal
- isequaln
- isfinite
- isinf
- ismatrix
- ismember
- isnat
- isreal
- isrow
- isscalar
- issorted
- issortedrows
- isvector
- le
- length
- linspace
- lt
- max
- mean
- min
- minus
- minute
- NaT
- ndims

- `ne`
- `numel`
- `permute`
- `plus`
- `posixtime`
- `repmat`
- `reshape`
- `setdiff`
- `setxor`
- `size`
- `sort`
- `sortrows`
- `topkrows`
- `transpose`
- `union`
- `unique`
- `vertcat`
- `ymd`

See Also

More About

- “Define Datetime Array Inputs” on page 22-5
- “Datetime Array Limitations for Code Generation” on page 22-7

Define Datetime Array Inputs

You can define `datetime` array inputs at the command line. Programmatic specification of `datetime` input types by using preconditioning (`assert` statements) is not supported.

Define Datetime Array Inputs at the Command Line

Use one of these procedures:

- “Provide an Example Datetime Array Input” on page 22-5
- “Provide a Datetime Array Type” on page 22-5
- “Provide a Constant Datetime Array Input” on page 22-5

Alternatively, if you have a test file that calls your entry-point function with example inputs, you can determine the input types by using `coder.getArgTypes`.

Provide an Example Datetime Array Input

Use the `-args` option:

```
D = datetime(2019,1:12,1,12,0,0);
fiaccl myFunction -args {D}
```

Provide a Datetime Array Type

To provide a type for a `datetime` array to `fiaccl`:

- 1 Define a `datetime` array. For example:


```
D = datetime(2019,1:12,1,12,0,0);
```
- 2 Create a type from `D`.


```
t = coder.typeof(D);
```
- 3 Pass the type to `fiaccl` by using the `-args` option.


```
fiaccl myFunction -args {t}
```

Provide a Constant Datetime Array Input

To specify that a `datetime` array input is constant, use `coder.Constant` with the `-args` option:

```
D = datetime(2019,1:12,1,12,0,0);
fiaccl myFunction -args {coder.Constant(C)}
```

Representation of Datetime Arrays

A `coder` type object for a `datetime` array describes the object and its properties. Use `coder.typeof` or pass `datetime` as a string scalar to `coder.newtype`.

The `coder` type object displays a succinct description of the object properties while excluding internal state values. Nonconstant properties display their type and size, while constant properties display only their values. For example:

```
t = datetime(2019,1:12,1,12,0,0);  
tType = coder.typeof(t)
```

The representation of variable `t` is stored in coder type object `tType`.

```
tType =  
  
    matlab.coder.type.DatetimeType  
    1x12 datetime  
    Format : 1x0 char  
    TimeZone : 1x0 char
```

If your workflow requires the legacy representation of coder type objects, use the `getCoderType` function on the variable that has the new representation of your class or object. See “Legacy Representation of Coder Type Objects” (MATLAB Coder).

Resize Object Properties by Using `coder.resize`

You can resize most objects by using `coder.resize`. You can resize objects, its properties and create arrays within the properties.

For a `datetime` coder object, you can resize the object properties:

```
t = datetime(2019,1:12,1,12,0,0);  
tType = coder.typeof(t)  
tType.Format = coder.resize(tType.Format, [1 12])
```

This code resizes the `Format` property to be a `1x12 char` property.

```
tType =  
  
    matlab.coder.type.DatetimeType  
    1x12 datetime  
    Format : 1x12 char  
    TimeZone : 1x0 char
```

You can also resize the object by using `coder.resize`. See “Edit and Represent Coder Type Objects and Properties” (MATLAB Coder).

See Also

`datetime` | `NaN` | `coder.Constant` | `coder.typeof`

More About

- “Code Generation for Datetime Arrays” on page 22-2
- “Datetime Array Limitations for Code Generation” on page 22-7

Datetime Array Limitations for Code Generation

When you create `datetime` arrays in MATLAB code that you intend for code generation, you must specify the values by using the `datetime` function. See “Dates and Time”.

For `datetime` arrays, code generation does not support the following inputs and operations:

- Text inputs. For example, specifying a character vector as the input argument produces an error.

```
function d = foo() %#codegen
    d = datetime('2019-12-01');
end
```

- The 'Format' name-value pair argument. You cannot specify the display format by using the `datetime` function, or by setting the `Format` property of a `datetime` array. To use a specific display format, create a `datetime` array in MATLAB, then pass it as an input argument to a function that is intended for code generation.
- The 'TimeZone' name-value pair argument and the `TimeZone` property. When you use `datetime` arrays in code that is intended for code generation, they must be unzoned.
- Setting time component properties. For example, setting the `Hour` property in the following code produces an error:

```
d = datetime;
d.Hour = 2;
```

- Growth by assignment. For example, assigning a value beyond the end of an array produces an error.

```
function d = foo() %#codegen
    d = datetime(2019,1:12,1,12,0,0);
    d(13) = datetime(2020,1,1,12,0,0);
end
```

- Deleting an element. For example, assigning an empty array to an element produces an error.

```
function d = foo() %#codegen
    d = datetime(2019,1:12,1,12,0,0);
    d(1) = [];
end
```

- Converting `datetime` values to text by using the `char`, `cellstr`, or `string` functions.

Limitations that apply to classes also apply to `datetime` arrays. For more information, see “MATLAB Classes Definition for Code Generation” (MATLAB Coder).

See Also

`datetime` | `NaT`

More About

- “Code Generation for Datetime Arrays” on page 22-2
- “Define Datetime Array Inputs” on page 22-5

Code Generation for Duration Arrays

- “Code Generation for Duration Arrays” on page 23-2
- “Define Duration Array Inputs” on page 23-6
- “Duration Array Limitations for Code Generation” on page 23-8

Code Generation for Duration Arrays

In this section...

“Define Duration Arrays for Code Generation” on page 23-2

“Allowed Operations on Duration Arrays” on page 23-2

“MATLAB Toolbox Functions That Support Duration Arrays” on page 23-3

The values in a duration array represent elapsed times in units of fixed length, such as hours, minutes, and seconds. You can create elapsed times in terms of fixed-length (24-hour) days and fixed-length (365.2425-day) years.

You can add, subtract, sort, compare, concatenate, and plot duration arrays.

When you use duration arrays with code generation, adhere to these restrictions.

Define Duration Arrays for Code Generation

For code generation, use the `duration` function to create duration arrays. For example, suppose the input arguments to your MATLAB function are three numeric arrays of arbitrary size whose elements specify lengths of time as hours, minutes, and seconds. You can create a duration array from these three input arrays.

```
function d = foo(h,m,s) %#codegen
    d = duration(h,m,s);
end
```

You can use the `years`, `days`, `hours`, `minutes`, `seconds`, and `milliseconds` functions to create duration arrays in units of years, days, hours, minutes, or seconds. For example, you can create an array of hours from an input numeric array.

```
function d = foo(h) %#codegen
    d = hours(h);
end
```

Allowed Operations on Duration Arrays

For code generation, you are restricted to the operations on duration arrays listed in this table.

| Operation | Example | Notes |
|---------------------------------------|---|--|
| assignment operator: = | <pre>d = duration(1:3,0,0); d(1) = hours(5); d = duration(1:3,0,0); d(1) = hours(5);</pre> | <p>Code generation does not support using the assignment operator = to:</p> <ul style="list-style-type: none"> Delete an element. Expand the size of a duration array. |
| relational operators: < > <= >= == ~= | <pre>d = duration(1:3,0,0); tf = d(1) < d(2); d = duration(1:3,0,0); tf = d(1) < d(2);</pre> | <p>Code generation supports relational operators.</p> |

| Operation | Example | Notes |
|--------------------|---|---|
| indexing operation | <pre>d = duration(1:3,0,0); idx = [1 2]; d(idx); idx = logical([1 1 0]); d(idx); d = duration(1:3,0,0); idx = [1 2]; d(idx); idx = logical([1 1 0]); d(idx);</pre> | Code generation supports indexing by position, linear indexing, and logical indexing. |
| concatenation | <pre>d1 = duration(1:3,0,0); d2 = duration(4,30,0); d = [d1 d2]; d1 = duration(1:3,0,0); d2 = duration(4,30,0); d = [d1 d2];</pre> | Code generation supports concatenation of duration arrays. |

MATLAB Toolbox Functions That Support Duration Arrays

For code generation, you can use duration arrays with these MATLAB toolbox functions:

- abs
- cat
- ceil
- colon
- cummax
- cummin
- cumsum
- ctranspose
- datevec
- days
- diff
- duration
- eps
- eq
- floor
- ge
- gt
- hms
- horzcat
- hours
- interp1
- intersect
- iscolumn

- isempty
- isequal
- isequaln
- isfinite
- isinf
- ismatrix
- ismember
- isnan
- isreal
- isrow
- isscalar
- issorted
- issortedrows
- isvector
- ldivide
- le
- length
- linspace
- lt
- max
- mean
- median
- milliseconds
- min
- minus
- minutes
- mldivide
- mode
- mrdivide
- mod
- mtimes
- ndims
- ne
- nnz
- numel
- permute
- plus
- repmat
- rdivide

- rem
- reshape
- seconds
- setdiff
- setxor
- sign
- size
- sort
- sortrows
- std
- sum
- times
- transpose
- uminus
- union
- unique
- uplus
- vertcat
- years

See Also

More About

- “Define Duration Array Inputs” on page 23-6
- “Duration Array Limitations for Code Generation” on page 23-8

Define Duration Array Inputs

You can define duration array inputs at the command line. Programmatic specification of duration input types by using preconditioning (`assert` statements) is not supported.

Define Duration Array Inputs at the Command Line

Use one of these procedures:

- “Provide an Example Duration Array Input” on page 23-6
- “Provide a Duration Array Type” on page 23-6
- “Provide a Constant Duration Array Input” on page 23-6

Alternatively, if you have a test file that calls your entry-point function with example inputs, you can determine the input types by using `coder.getArgTypes`.

Provide an Example Duration Array Input

Use the `-args` option:

```
D = duration(1:3,0,0);  
fiaccl myFunction -args {D}
```

Provide a Duration Array Type

To provide a type for a duration array to `fiaccl`:

- 1 Define a duration array. For example:

```
D = duration(1:3,0,0);
```
- 2 Create a type from `D`.

```
t = coder.typeof(D);
```
- 3 Pass the type to `fiaccl` by using the `-args` option.

```
fiaccl myFunction -args {t}
```

Provide a Constant Duration Array Input

To specify that a duration array input is constant, use `coder.Constant` with the `-args` option:

```
D = duration(1:3,0,0);  
fiaccl myFunction -args {coder.Constant(C)}
```

Representation of Duration Arrays

A `coder` type object for a duration array describes the object and its properties. Use `coder.typeof` or pass `duration` as a string scalar to `coder.newtype`.

The `coder` type object displays a succinct description of the object properties while excluding internal state values. Nonconstant properties display their type and size, while constant properties display only their values. For example:

```
tType = coder.newtype('duration')
```

A representation of an empty duration variable is stored in coder type object `tType`.

```
tType =  
  
    matlab.coder.type.DurationType  
        1x1 duration  
        Format : 1x8 char
```

If your workflow requires the legacy representation of coder type objects, use the `getCoderType` function on the variable that has the new representation of your class or object. See “Legacy Representation of Coder Type Objects” (MATLAB Coder).

Resize duration Properties by Editing Object Properties

You can resize most objects by editing the object properties. You can resize duration objects, its properties and create arrays within the properties.

For a duration coder object, you can resize the object properties:

```
t = duration((1:3),0,0);  
tType = coder.typeof(t)  
tType.Format = 'DD/MM/YYYY'
```

This code resizes the `Format` property to be a 1x10 char property.

```
tType =  
  
    matlab.coder.type.DurationType  
        1x3 duration  
        Format : 1x10 char
```

You can also resize the object by using `coder.resize`. See “Edit and Represent Coder Type Objects and Properties” (MATLAB Coder).

See Also

`duration` | `coder.Constant` | `coder.typeof`

More About

- “Code Generation for Duration Arrays” on page 23-2
- “Duration Array Limitations for Code Generation” on page 23-8

Duration Array Limitations for Code Generation

When you create duration arrays in MATLAB code that you intend for code generation, you must specify the durations by using the `duration`, `years`, `days`, `hours`, `minutes`, `seconds`, or `milliseconds` functions. See “Dates and Time”.

For duration arrays, code generation does not support the following inputs and operations:

- Text inputs. For example, specifying a character vector as the input argument produces an error.

```
function d = foo() %#codegen
    d = duration('01:30:00');
end
```

- Growth by assignment. For example, assigning a value beyond the end of an array produces an error.

```
function d = foo() %#codegen
    d = duration(1:3,0,0);
    d(4) = hours(4);
end
```

- Deleting an element. For example, assigning an empty array to an element produces an error.

```
function d = foo() %#codegen
    d = duration(1:3,0,0);
    d(1) = [];
end
```

- Converting duration values to text by using the `char`, `cellstr`, or `string` functions.

Limitations that apply to classes also apply to duration arrays. For more information, see “MATLAB Classes Definition for Code Generation” (MATLAB Coder).

See Also

`duration` | `years` | `days` | `hours` | `minutes` | `seconds` | `milliseconds`

More About

- “Code Generation for Duration Arrays” on page 23-2
- “Define Duration Array Inputs” on page 23-6

Code Generation for Function Handles

Function Handle Limitations for Code Generation

When you use function handles in MATLAB code intended for code generation, adhere to the following restrictions:

Do not use the same bound variable to reference different function handles

In some cases, using the same bound variable to reference different function handles causes a compile-time error. For example, this code does not compile:

```
function y = foo(p)
x = @plus;
if p
    x = @minus;
end
y = x(1, 2);
```

Do not pass function handles to or from `coder.ceval`

You cannot pass function handles as inputs to or outputs from `coder.ceval`. For example, suppose that `f` and `str.f` are function handles:

```
f = @sin;
str.x = pi;
str.f = f;
```

The following statements result in compilation errors:

```
coder.ceval('foo', @sin);
coder.ceval('foo', f);
coder.ceval('foo', str);
```

Do not associate a function handle with an extrinsic function

You cannot create a function handle that references an extrinsic MATLAB function.

Do not pass function handles to or from extrinsic functions

You cannot pass function handles to or from `feval` and other extrinsic MATLAB functions.

Do not pass function handles to or from entry-point functions

You cannot pass function handles as inputs to or outputs from entry-point functions. For example, consider this function:

```
function x = plotFcn(fhandle, data)

assert(isa(fhandle,'function_handle') && isa(data,'double'));

plot(data, fhandle(data));
x = fhandle(data);
```

In this example, the function `plotFcn` receives a function handle and its data as inputs. `plotFcn` attempts to call the function referenced by the `fhandle` with the input `data` and plot the results. However, this code generates a compilation error. The error indicates that the function `isa` does not recognize `'function_handle'` as a class name when called inside a MATLAB function to specify properties of inputs.

See Also

More About

- “Use the `coder.extrinsic` Construct” on page 14-7

Code Generation for MATLAB Structures

- “Structure Definition for Code Generation” on page 25-2
- “Structure Operations Allowed for Code Generation” on page 25-3
- “Define Scalar Structures for Code Generation” on page 25-4
- “Define Arrays of Structures for Code Generation” on page 25-6
- “Index Substructures and Fields” on page 25-8
- “Assign Values to Structures and Fields” on page 25-10
- “Pass Large Structures as Input Parameters” on page 25-11

Structure Definition for Code Generation

To generate efficient standalone code for structures, you must define and use structures differently than you normally would when running your code in the MATLAB environment:

| What's Different | More Information |
|---|---|
| Use a restricted set of operations. | "Structure Operations Allowed for Code Generation" on page 25-3 |
| Observe restrictions on properties and values of scalar structures. | "Define Scalar Structures for Code Generation" on page 25-4 |
| Make structures uniform in arrays. | "Define Arrays of Structures for Code Generation" on page 25-6 |
| Reference structure fields individually during indexing. | "Index Substructures and Fields" on page 25-8 |
| Avoid type mismatch when assigning values to structures and fields. | "Assign Values to Structures and Fields" on page 25-10 |

Structure Operations Allowed for Code Generation

To generate efficient standalone code for MATLAB structures, you are restricted to the following operations:

- Index structure fields using dot notation
- Define primary function inputs as structures
- Pass structures to local functions

Define Scalar Structures for Code Generation

In this section...

“Restrictions When Defining Scalar Structures by Assignment” on page 25-4
 “Adding Fields in Consistent Order on Each Control Flow Path” on page 25-4
 “Restriction on Adding New Fields After First Use” on page 25-4

Restrictions When Defining Scalar Structures by Assignment

When you define a scalar structure by assigning a variable to a preexisting structure, you do not need to define the variable before the assignment. However, if you already defined that variable, it must have the same class, size, and complexity as the structure you assign to it. In the following example, `p` is defined as a structure that has the same properties as the predefined structure `S`:

```
...
S = struct('a', 0, 'b', 1, 'c', 2);
p = S;
...
```

Adding Fields in Consistent Order on Each Control Flow Path

When you create a structure, you must add fields in the same order on each control flow path. For example, the following code generates a compiler error because it adds the fields of structure `x` in a different order in each `if` statement clause:

```
function y = fcn(u) %#codegen
if u > 0
    x.a = 10;
    x.b = 20;
else
    x.b = 30; % Generates an error (on variable x)
    x.a = 40;
end
y = x.a + x.b;
```

In this example, the assignment to `x.a` comes before `x.b` in the first `if` statement clause, but the assignments appear in reverse order in the `else` clause. Here is the corrected code:

```
function y = fcn(u) %#codegen
if u > 0
    x.a = 10;
    x.b = 20;
else
    x.a = 40;
    x.b = 30;
end
y = x.a + x.b;
```

Restriction on Adding New Fields After First Use

You cannot add fields to a structure after you perform the following operations on the structure:

- Reading from the structure
- Indexing into the structure array
- Passing the structure to a function

For example, consider this code:

```
...
x.c = 10; % Defines structure and creates field c
y = x; % Reads from structure
x.d = 20; % Generates an error
...
```

In this example, the attempt to add a new field `d` after reading from structure `x` generates an error.

This restriction extends across the structure hierarchy. For example, you cannot add a field to a structure after operating on one of its fields or nested structures, as in this example:

```
function y = fcn(u) %#codegen

x.c = 10;
y = x.c;
x.d = 20; % Generates an error
```

In this example, the attempt to add a new field `d` to structure `x` after reading from the structure's field `c` generates an error.

Define Arrays of Structures for Code Generation

In this section...

“Ensuring Consistency of Fields” on page 25-6

“Using repmat to Define an Array of Structures with Consistent Field Properties” on page 25-6

“Defining an Array of Structures by Using struct” on page 25-6

“Defining an Array of Structures Using Concatenation” on page 25-7

Ensuring Consistency of Fields

For code generation, when you create an array of MATLAB structures, corresponding fields in the array elements must have the same size, type, and complexity.

Once you have created the array of structures, you can make the structure fields variable-size by using `coder. varsize`. See “Declare Variable-Size Structure Fields” (MATLAB Coder).

Using repmat to Define an Array of Structures with Consistent Field Properties

You can create an array of structures from a scalar structure by using the MATLAB `repmat` function, which replicates and tiles an existing scalar structure:

- 1 Create a scalar structure, as described in “Define Scalar Structures for Code Generation” on page 25-4.
- 2 Call `repmat`, passing the scalar structure and the dimensions of the array.
- 3 Assign values to each structure using standard array indexing and structure dot notation.

For example, the following code creates `X`, a 1-by-3 array of scalar structures. Each element of the array is defined by the structure `s`, which has two fields, `a` and `b`:

```
...  
s.a = 0;  
s.b = 0;  
X = repmat(s,1,3);  
X(1).a = 1;  
X(2).a = 2;  
X(3).a = 3;  
X(1).b = 4;  
X(2).b = 5;  
X(3).b = 6;  
...
```

Defining an Array of Structures by Using struct

To create an array of structures using the `struct` function, specify the field value arguments as cell arrays. Each cell array element is the value of the field in the corresponding structure array element. For code generation, corresponding fields in the structures must have the same type. Therefore, the elements in a cell array of field values must have the same type.

For example, the following code creates a 1-by-3 structure array. For each structure in the array of structures, `a` has type `double` and `b` has type `char`.

```
s = struct('a', {1 2 3}, 'b', {'a' 'b' 'c'});
```

Defining an Array of Structures Using Concatenation

To create a small array of structures, you can use the concatenation operator, square brackets (`[]`), to join one or more structures into an array. See “Creating, Concatenating, and Expanding Matrices”. For code generation, the structures that you concatenate must have the same size, class, and complexity.

For example, the following code uses concatenation and a local function to create the elements of a 1-by-3 structure array:

```
...  
W = [ sab(1,2) sab(2,3) sab(4,5) ];  
  
function s = sab(a,b)  
    s.a = a;  
    s.b = b;  
...
```

See Also

MATLAB Function

Related Examples

- “Define Scalar Structures for Code Generation” on page 25-4
- “Define and Use Structure Parameters”
- “Create Structures in MATLAB Function Blocks”

Index Substructures and Fields

Use these guidelines when indexing substructures and fields for code generation:

Reference substructure field values individually using dot notation

For example, the following MATLAB code uses dot notation to index fields and substructures:

```

...
substruct1.a1 = 15.2;
substruct1.a2 = int8([1 2;3 4]);

mystruct = struct('ele1',20.5,'ele2',single(100),
                 'ele3',substruct1);

substruct2 = mystruct;
substruct2.ele3.a2 = 2*(substruct1.a2);
...

```

The generated code indexes elements of the structures in this example by resolving symbols as follows:

| Dot Notation | Symbol Resolution |
|-------------------------|--|
| substruct1.a1 | Field a1 of local structure substruct1 |
| substruct2.ele3.a1 | Value of field a1 of field ele3, a substructure of local structure substruct2 |
| substruct2.ele3.a2(1,1) | Value in row 1, column 1 of field a2 of field ele3, a substructure of local structure substruct2 |

Reference field values individually in structure arrays

To reference the value of a field in a structure array, you must index into the array to the structure of interest and then reference that structure's field individually using dot notation, as in this example:

```

...
y = X(1).a % Extracts the value of field a
           % of the first structure in array X
...

```

To reference all the values of a particular field for each structure in an array, use this notation in a for loop, as in this example:

```

...
s.a = 0;
s.b = 0;
X = repmat(s,1,5);
for i = 1:5
    X(i).a = i;
    X(i).b = i+1;
end
...

```

This example uses the `repmat` function to define an array of structures, each with two fields `a` and `b` as defined by `s`. See “Define Arrays of Structures for Code Generation” on page 25-6 for more information.

Do not reference fields dynamically

You cannot reference fields in a structure by using dynamic names, which express the field as a variable expression that MATLAB evaluates at run time (see “Generate Field Names from Variables”).

Assign Values to Structures and Fields

When assigning values to a structure, substructure, or field for code generation, use these guidelines:

Field properties must be consistent across structure-to-structure assignments

| If: | Then: |
|--|--|
| Assigning one structure to another structure. | Define each structure with the same number, type, and size of fields. |
| Assigning one structure to a substructure of a different structure and vice versa. | Define the structure with the same number, type, and size of fields as the substructure. |
| Assigning an element of one structure to an element of another structure. | The elements must have the same type and size. |

For structures with constant fields, do not assign field values inside control flow constructs

In the following code, the code generator recognizes that the structure fields `s.a` and `s.b` are constants.

```
function y = mystruct()
s.a = 3;
s.b = 5;
y = zeros(s.a,s.b);
```

If a field of a structure is assigned inside a control flow construct, the code generator does not recognize that `s.a` and `s.b` are constant. Consider the following code:

```
function y = mystruct(x)
s.a = 3;
if x > 1
    s.b = 4;
else
    s.b = 5;
end
y = zeros(s.a,s.b);
```

If variable-sizing is enabled, `y` is treated as a variable-size array. If variable-sizing is disabled, `y`, the code generator reports an error.

Do not assign mxArray to structures

You cannot assign `mxArrays` to structure elements; convert `mxArrays` to known types before code generation (see “Working with `mxArrays`” on page 14-9).

Do not assign handle classes or sparse arrays to global structure variables

Global structure variables cannot contain handle objects or sparse arrays.

Pass Large Structures as Input Parameters

If you generate a MEX function for a MATLAB function that takes a large structure as an input parameter, for example, a structure containing fields that are matrices, the MEX function might fail to load. This load failure occurs because, when you generate a MEX function from a MATLAB function that has input parameters, the code generator allocates memory for these input parameters on the stack. To avoid this issue, pass the structure by reference to the MATLAB function. For example, if the original function signature is:

```
y = foo(a, S)
```


where S is the structure input, rewrite the function to:

```
[y, S] = foo(a, S)
```


Functions, Classes, and System Objects Supported for Code Generation

Functions and Objects Supported for C/C++ Code Generation

You can generate efficient C/C++ code for a subset of MATLAB built-in functions and toolbox functions and System objects that you call from MATLAB code.

These functions and System objects are listed in the following tables. In these tables, a  icon before the name of a function or a System object indicates that there are specific usage notes and limitations related to C/C++ code generation for that function or System object. To view these usage notes and limitations, in the corresponding reference page, scroll down to the **Extended Capabilities** section at the bottom and expand the **C/C++ Code Generation** section.

- Functions and Objects Supported for C/C++ Code Generation (Category List)
- Functions and Objects Supported for C/C++ Code Generation (Alphabetical List)

See Also

Related Examples

- “MATLAB Language Features Supported for C/C++ Code Generation” on page 19-29

Code Generation for Tables

- “Code Generation for Tables” on page 27-2
- “Define Table Inputs” on page 27-5
- “Table Limitations for Code Generation” on page 27-8

Code Generation for Tables

In this section...

“Define Tables for Code Generation” on page 27-2

“Allowed Operations on Tables” on page 27-2

“MATLAB Toolbox Functions That Support Tables” on page 27-3

The `table` data type is a data type suitable for column-oriented or tabular data that is often stored as columns in a text file or in a spreadsheet. Tables consist of rows and column-oriented variables. Each variable in a table can have a different data type and a different size with one restriction: each variable must have the same number of rows. For more information, see “Tables”.

When you use tables with code generation, adhere to these restrictions.

Define Tables for Code Generation

For code generation, use the `table` function. For example, suppose the input arguments to your MATLAB function are three arrays that have the same number of rows and a cell array that has variable names. You can create a table that contains these arrays as table variables.

```
function T = foo(A,B,C,vnames) %#codegen
    T = table(A,B,C,'VariableNames',vnames);
end
```

You can use the `array2table`, `cell2table`, and `struct2table` functions to convert arrays, cell arrays, and structures to tables. For example, you can convert an input cell array to a table.

```
function T = foo(C,vnames) %#codegen
    T = cell2table(C,'VariableNames',vnames);
end
```

For code generation, you must supply table variable names when you create a table. Table variable names do not have to be valid MATLAB identifiers. The names must be composed of ASCII characters, but can include any ASCII characters (such as commas, dashes, and space characters).

Allowed Operations on Tables

For code generation, you are restricted to the operations on tables listed below.

| Operation | Example | Notes |
|------------------------|--|--|
| assignment operator: = | <code>T = table(A,B,C,'VariableNames',vnames)</code> <code>T(:,1) = D;</code> | Code generation does not support using the assignment operator = to: <ul style="list-style-type: none"> Delete a variable or a row. Add a variable or a row. |

| Operation | Example | Notes |
|--------------------|---|--|
| indexing operation | <pre>T = table(A,B,C, 'VariableNames', vnames); T(1:5,1:3);</pre> | <p>Code generation supports indexing by position, variable or row name, and logical indexing.</p> <p>Code generation supports:</p> <ul style="list-style-type: none"> • Table indexing with smooth parentheses, (). • Content indexing with curly braces, {}. • Dot notation to access a table variable. |
| concatenation | <pre>T1 = table(A,B,C, 'VariableNames', vnames); T2 = table(D,E,F, 'VariableNames', vnames); T = [T1 ; T2];</pre> | <p>Code generation supports table concatenation.</p> <ul style="list-style-type: none"> • For vertical concatenation, tables must have variables that have the same names in the same order. • For horizontal concatenation, tables must have the same number of rows. If the tables have row names, then they must have the same row names in the same order. |

MATLAB Toolbox Functions That Support Tables

For code generation, you can use tables with these MATLAB toolbox functions:

- `addvars`
- `array2table`
- `cat`
- `cell2table`
- `convertvars`
- `height`
- `horzcat`
- `innerjoin`
- `intersect`
- `isempty`
- `ismember`
- `issortedrows`
- `join`
- `mergevars`

- `movevars`
- `ndims`
- `numel`
- `outerjoin`
- `removevars`
- `renamevars`
- `rows2vars`
- `setdiff`
- `setxor`
- `size`
- `sortrows`
- `splitvars`
- `stack`
- `struct2table`
- `table`
- `table2array`
- `table2cell`
- `table2struct`
- `union`
- `unique`
- `unstack`
- `varfun`
- `vertcat`
- `width`

See Also

More About

- “Define Table Inputs” on page 27-5
- “Table Limitations for Code Generation” on page 27-8

Define Table Inputs

You can define table inputs at the command line. Programmatic specification of table input types by using preconditioning (assert statements) is not supported.

Define Table Inputs at the Command Line

Use one of these procedures:

- “Provide an Example Table Input” on page 27-5
- “Provide a Table Type” on page 27-5
- “Provide a Constant Table Input” on page 27-5

Alternatively, if you have a test file that calls your entry-point function with example inputs, you can determine the input types by using `coder.getArgTypes`.

Provide an Example Table Input

Use the `-args` option:

```
T = table(A,B,C,'VariableNames',vnames);
fiaccl myFunction -args {T}
```

Provide a Table Type

To provide a type for a table to `fiaccl`:

- 1 Define a table. For example:


```
T = table(A,B,C,'VariableNames',vnames);
```
- 2 Create a type from T.


```
t = coder.typeof(T);
```
- 3 Pass the type to `fiaccl` by using the `-args` option.


```
fiaccl myFunction -args {t}
```

Provide a Constant Table Input

To specify that a table input is constant, use `coder.Constant` with the `-args` option:

```
T = table(A,B,C,'VariableNames',vnames);
fiaccl myFunction -args {coder.Constant(T)}
```

Representation of Tables

A coder type object for a table describes the object and its properties. Use `coder.typeof` or pass `table` as a string scalar to `coder.newtype`.

The coder type object displays a succinct description of the object properties while excluding internal state values. Nonconstant properties display their type and size, while constant properties display only their values. For example:

```
A = [1 2 3]';
B = [4 5 6]';
```

```
C = [7 8 9]';
t = table(A,B,C);
tType = coder.typeof(t)
```

The representation of variable `t` is stored in coder type object `tType`.

```
tType =
    matlab.coder.type.TableType
    3x3 table
           Data : 1x3 homogeneous cell
           Description : 1x0 char
           UserData : 0x0 double
           DimensionNames : {'Row'} {'Variables'}
           VariableNames : {'A'} {'B'} {'C'}
VariableDescriptions : 1x3 homogeneous cell
           VariableUnits : 1x3 homogeneous cell
           VariableContinuity : 1x3 matlab.internal.coder.tabular.Continuity
           RowNames : 0x0 homogeneous cell
```

If your workflow requires the legacy representation of coder type objects, use the `getCoderType` function on the variable that has the new representation of your class or object. See “Legacy Representation of Coder Type Objects” (MATLAB Coder).

Resize Object Properties by Using `coder.resize`

You can resize most objects by using `coder.resize`. You can resize objects, its properties and create arrays within the properties.

For a table coder object, you can resize the object properties:

```
A = [1 2 3]';
B = [4 5 6]';
C = [7 8 9]';
t = table(A,B,C);
tType = coder.typeof(t)
tType.Description = coder.resize(tType.Description,[1 12],[0 1])
```

This code resizes the `Description` property to be a `1x:12 char` property which has an upper bound of 12.

```
tType =
    matlab.coder.type.TableType
    3x3 table
           Data : 1x3 homogeneous cell
           Description : 1x:12 char
           UserData : 0x0 double
           DimensionNames : {'Row'} {'Variables'}
           VariableNames : {'A'} {'B'} {'C'}
VariableDescriptions : 1x3 homogeneous cell
           VariableUnits : 1x3 homogeneous cell
           VariableContinuity : 1x3 matlab.internal.coder.tabular.Continuity
           RowNames : 0x0 homogeneous cell
```

You can also resize the object by using `coder.resize`. See “Edit and Represent Coder Type Objects and Properties” (MATLAB Coder).

See Also

`table` | `coder.Constant` | `coder.typeof`

More About

- “Code Generation for Tables” on page 27-2
- “Table Limitations for Code Generation” on page 27-8

Table Limitations for Code Generation

If you create tables, modify them, or use table functions in MATLAB code that you intend for code generation, then code generation has limitations described in the next sections. Limitations that apply to classes also apply to tables. For more information on class limitations, see “MATLAB Classes Definition for Code Generation” (MATLAB Coder).

Creating Tables Limitations

If your MATLAB code creates tables, then code generation has these limitations.

| Inputs for Table Creation | Limitations |
|---|--|
| Any inputs | <ul style="list-style-type: none"> Table variable names do not have to be valid MATLAB identifiers. The names must be composed of ASCII characters, which can include commas, dashes, and space characters. |
| Table created from input arrays | <ul style="list-style-type: none"> You must specify variables names by using the 'VariableNames' name-value argument when creating tables from input arrays by using the <code>table</code>, <code>array2table</code>, or <code>cell2table</code> functions. |
| Table created with preallocated variables | <ul style="list-style-type: none"> You do not have to specify the 'VariableNames' argument when you preallocate a table by using the <code>table</code> function and the 'Size' name-value argument. You can specify only the following data types by using the 'VariableTypes' name-value argument: <ul style="list-style-type: none"> 'double' 'single' 'doublenan' or 'doubleNaN' 'singlenan' or 'singleNaN' 'int8', 'int16', 'int32', or 'int64' 'uint8', 'uint16', 'uint32', or 'uint64' 'logical' 'duration' 'cellstr' 'char' |

Modifying Tables Limitations

If your MATLAB code modifies data in a table or its properties, then code generation has these limitations.

| Table Operation or Property | Limitations |
|---|---|
| VariableNames, RowNames, DimensionNames, or UserData properties | <ul style="list-style-type: none"> You cannot change the VariableNames, RowNames, DimensionNames, or UserData properties of a table after you create it. <p>You can specify the 'VariableNames', 'RowNames', and 'DimensionNames' input arguments when you create a table. These input arguments specify the properties.</p> |
| Table indices that specify variables as input arguments to generated code | <ul style="list-style-type: none"> To pass table indices that specify variables as input arguments into generated code, first make the indices constant by using the <code>coder.Constant</code> function. If table indices are not constant, then indexing into variables produces an error. |
| Custom metadata | <ul style="list-style-type: none"> You cannot add custom metadata to a table. The <code>addprop</code> and <code>rmprop</code> functions are not supported. |
| Assignments that change size of table | <ul style="list-style-type: none"> You cannot change the size of a table by assignments. For example, adding a new row produces an error. <pre data-bbox="906 989 1581 1129">function T = foo() %#codegen T = table((1:3)',(1:3)', 'VariableNames', ... {'Var1', 'Var2'}); T{4,2} = 5; end</pre> <p>Deleting a row or a variable also produces an error.</p> |
| Vertical concatenation | <ul style="list-style-type: none"> When you vertically concatenate tables, they must have the same variable names in the same order. In MATLAB, the variable names must be the same but can be in different orders in the tables. |
| Horizontal concatenation | <ul style="list-style-type: none"> When you horizontally concatenate tables and the tables have row names, they must have the same row names in the same order. In MATLAB, the row names must be the same but can be in different orders in the tables. |
| Table variables that are N-D cell arrays | <ul style="list-style-type: none"> If two tables have variables that are N-D cell arrays, then the tables cannot be vertically concatenated. You cannot use curly braces to extract data from multiple table variables that are N-D cell arrays because this operation is horizontal concatenation. |

Using Table Functions Limitations

If your MATLAB code uses the functions listed in the table, then code generation has these limitations.

| Function | Limitations |
|---|--|
| convertvars | <ul style="list-style-type: none"> Function handles are not supported. The second and third input arguments (<code>vars</code> and <code>dataType</code>) must be constant. You cannot specify <code>dataType</code> as <code>'cell'</code>, <code>'cellstr'</code>, or <code>'char'</code>. |
| innerjoin | <ul style="list-style-type: none"> In general, the input tables cannot have any nonkey variables with the same names. However, you can join subsets of the input tables if you specify the <code>'LeftVariables'</code> and <code>'RightVariables'</code> name-value arguments. Specify these arguments so that no variable name appears in both <code>'LeftVariables'</code> and <code>'RightVariables'</code>. The values of these name-value arguments must be constant: <ul style="list-style-type: none"> <code>'Keys'</code> <code>'LeftKeys'</code> <code>'RightKeys'</code> <code>'LeftVariables'</code> <code>'RightVariables'</code> Nested tables are not supported. |
| intersect setdiff setxor union | <ul style="list-style-type: none"> These functions support unsorted tables in all cases. You do not have to specify the <code>'stable'</code> option. |
| issortedrows | <ul style="list-style-type: none"> The input argument <code>vars</code> must be constant. If any table variables have multiple columns, then those variables must have fixed widths. |

| Function | Limitations |
|----------|--|
| join | <ul style="list-style-type: none">• In general, input tables cannot have nonkey variables with the same names. However, you can join subsets of the input tables if you specify the name-value arguments:<ul style="list-style-type: none">• 'KeepOneCopy', where you list variables to take from the left input table only.• 'LeftVariables' and 'RightVariables', where you list variables to take from either the left input table or the right input table, but not both.• The values of these name-value arguments must be constant:<ul style="list-style-type: none">• 'Keys'• 'LeftKeys'• 'RightKeys'• 'LeftVariables'• 'RightVariables'• 'KeepOneCopy'• Nested tables are not supported. |
| movevars | <ul style="list-style-type: none">• The input argument <code>vars</code> cannot contain duplicate variable names. |

| Function | Limitations |
|-----------|---|
| outerjoin | <ul style="list-style-type: none"> • Input tables cannot have key variables with the same names unless the value of 'MergeKeys' is true (logical 1). • In general, the input tables cannot have any nonkey variables with the same names. However, you can join subsets of the input tables if you specify the 'LeftVariables' and 'RightVariables' name-value arguments. Specify these arguments so that no variable name appears in both 'LeftVariables' and 'RightVariables'. • The values of these name-value arguments must be constant: <ul style="list-style-type: none"> • 'Keys' • 'LeftKeys' • 'RightKeys' • 'MergeKeys' • 'LeftVariables' • 'RightVariables' • 'Type' • Nested tables are not supported. |
| rows2vars | <ul style="list-style-type: none"> • The input table cannot be variable-size. • The 'VariableNamesSource' name-value argument is not supported. • The value of the 'DataVariables' name-value argument must be constant. • The value of the 'VariableNamingRule' name-value argument must be constant. • If you assign row names to the input table, then the vector of row names must be constant. |
| sortrows | <ul style="list-style-type: none"> • The input argument vars must be constant. • If tblA has a variable that is a cell array of character vectors with multiple columns, then you cannot sort the table using the values in that variable. |
| splitvars | <ul style="list-style-type: none"> • The value of the 'NewVariableNames' name-value argument must be constant. • The variables that are split cannot have a variable number of columns. |

| Function | Limitations |
|----------|---|
| stack | <ul style="list-style-type: none"> • The second input argument, <code>vars</code>, must be constant. • The values of the <code>'ConstantVariables'</code>, <code>'NewDataVariableName'</code>, and <code>'IndexVariableName'</code> name-value arguments must be constant. |
| unstack | <ul style="list-style-type: none"> • The <code>'NewDataVariableNames'</code> name-value argument must be specified. Its value must be constant. • The <code>vars</code> and <code>ivars</code> input arguments (data variables and indicator variables) must be constant. • If you specify grouping variables and constant variables, then they must be constant. • If you specify an aggregation function, then it must be constant. • If a variable of the input table is a cell array of character vectors, then <code>unstack</code> fills empty cells in the corresponding output variable with 1-by-0 character arrays in the generated code. In MATLAB, <code>unstack</code> fills such gaps with 0-by-0 character arrays. • The <code>unstack</code> function does not support code generation when the input table has a variable that is a heterogeneous cell array that cannot be converted to a homogeneous cell array. <ul style="list-style-type: none"> • If the input has a variable that is a homogeneous cell array, or that can be converted to one, then the <code>'AggregationFunction'</code> name-value argument must be specified. The default value of <code>'AggregationFunction'</code> is <code>'unique'</code>. But the <code>unique</code> function does not support cell arrays. |

| Function | Limitations |
|----------|--|
| varfun | <ul style="list-style-type: none"> • The function handle input, <code>func</code>, must be constant. • While function handles can be inputs to <code>varfun</code> itself, they cannot be inputs to your entry point functions. Specify <code>func</code> within the code meant for code generation. For more information, see “Function Handle Limitations for Code Generation” (MATLAB Coder). • The values for all name-value arguments must be constant. • The <code>'ErrorHandler'</code> name-value argument is not supported for code generation. • Variable-size input arguments are not supported. • Grouping variables cannot have duplicate values in generated code. • You cannot specify the value of <code>'OutputFormat'</code> as <code>'cell'</code> if you specify the <code>'GroupingVariables'</code> name-value argument and the function returns a different data type for each variable specified by <code>'InputVariables'</code>. • If you specify groups and the number of groups is not known at compile time, and that number is zero, then empty double variables in the output might have sizes of 1-by-0 in generated code. In MATLAB, such variables have sizes of 0-by-0. |

See Also

`array2table` | `cell2table` | `struct2table` | `table`

More About

- “Code Generation for Tables” on page 27-2
- “Define Table Inputs” on page 27-5

Code Generation for Timetables

- “Code Generation for Timetables” on page 28-2
- “Define Timetable Inputs” on page 28-6
- “Timetable Limitations for Code Generation” on page 28-9

Code Generation for Timetables

In this section...

“Define Timetables for Code Generation” on page 28-2

“Allowed Operations on Timetables” on page 28-2

“MATLAB Toolbox Functions That Support Timetables” on page 28-3

The `timetable` data type is a data type suitable for tabular data with time-stamped rows. Like tables, timetables consist of rows and column-oriented variables. Each variable in a timetable can have a different data type and a different size with one restriction: each variable must have the same number of rows.

The *row times* of a timetable are time values that label the rows. You can index into a timetable by row time and variable. To index into a timetable, use smooth parentheses () to return a subtable or curly braces { } to extract the contents. You can refer to variables and to the vector of row times by their names. For more information, see “Timetables”.

When you use timetables with code generation, adhere to these restrictions.

Define Timetables for Code Generation

For code generation, use the `timetable` function. For example, suppose the input arguments to your MATLAB function are three arrays that have the same number of rows (A, B, and C), a `datetime` or `duration` vector containing row times (D), and a cell array that has variable names (vnames). You can create a timetable that contains these arrays as timetable variables.

```
function TT = foo(A,B,C,D,vnames) %#codegen
    TT = table(A,B,C,'RowTimes',D,'VariableNames',vnames);
end
```

To convert arrays and tables to timetables, use the `array2timetable` and `table2timetable` functions. For example, you can convert an input M-by-N matrix to a timetable, where each column of the matrix becomes a variable in the timetable. Assign row times by using a `duration` vector.

```
function TT = foo(A,D,vnames) %#codegen
    TT = array2timetable(A,'RowTimes',D,'VariableNames',vnames);
end
```

For code generation, you must supply timetable variable names when you create a timetable. Timetable variable names do not have to be valid MATLAB identifiers. The names must be composed of ASCII characters, but can include any ASCII characters (such as commas, dashes, and space characters).

The row times can have either the `datetime` or `duration` data type.

Allowed Operations on Timetables

For code generation, you are restricted to the operations on timetables listed in this table.

| Operation | Example | Notes |
|------------------------|--|--|
| Assignment operator: = | <pre>TT = timetable(A,B,C,'RowTimes',D,... 'VariableNames',vnames); TT{: ,1} = X;</pre> | <p>Code generation does not support using the assignment operator = to:</p> <ul style="list-style-type: none"> • Delete a variable or a row. • Add a variable or a row. |
| Indexing operation | <pre>D = seconds(1:10); TT = timetable(A,B,C,'RowTimes',D,... 'VariableNames',vnames); TT(seconds(3:7),1:3);</pre> | <p>Code generation supports indexing by position, variable or row time, and logical indexing. Also, you can index using objects created by using the <code>timerange</code> or <code>withtol</code> functions.</p> <p>Code generation supports:</p> <ul style="list-style-type: none"> • Timetable indexing with smooth parentheses, <code>()</code>. • Content indexing with curly braces, <code>{}</code>. • Dot notation to access a timetable variable. |
| Concatenation | <pre>TT1 = timetable(A,B,C,'RowTimes',D1,... 'VariableNames',vnames); TT2 = timetable(D,E,F,'RowTimes',D2,... 'VariableNames',vnames); TT = [TT1 ; TT2];</pre> | <p>Code generation supports timetable concatenation.</p> <ul style="list-style-type: none"> • For vertical concatenation, timetables must have variables that have the same names in the same order. • For horizontal concatenation, timetables must have the same number of rows. They also must have the same row times in the same order. |

MATLAB Toolbox Functions That Support Timetables

For code generation, you can use timetables with these MATLAB toolbox functions:

- `addvars`
- `array2timetable`
- `cat`
- `convertvars`
- `height`
- `horzcat`
- `innerjoin`

- intersect
- isempty
- ismember
- isregular
- issorted
- issortedrows
- join
- mergevars
- movevars
- ndims
- numel
- outerjoin
- removevars
- renamevars
- rows2vars
- retime
- setdiff
- setxor
- size
- sortrows
- splitvars
- stack
- synchronize
- table2timetable
- timerange
- timetable
- timetable2table
- union
- unique
- unstack
- varfun
- vertcat
- width
- withtol

See Also

More About

- “Define Timetable Inputs” on page 28-6

- “Timetable Limitations for Code Generation” on page 28-9

Define Timetable Inputs

You can define timetable inputs at the command line. Programmatic specification of timetable input types by using preconditioning (`assert` statements) is not supported.

Define Timetable Inputs at the Command Line

Use one of these procedures:

- “Provide an Example Timetable Input” on page 28-6
- “Provide a Timetable Type” on page 28-6
- “Provide a Constant Timetable Input” on page 28-6

Alternatively, if you have a test file that calls your entry-point function with example inputs, you can determine the input types by using `coder.getArgTypes`.

Provide an Example Timetable Input

Use the `-args` option:

```
TT = timetable(A,B,C,'RowTimes',D,'VariableNames',vnames);  
fiaccel myFunction -args {TT}
```

Provide a Timetable Type

To provide a type for a timetable to `fiaccel`:

- 1 Define a timetable. For example:

```
TT = timetable(A,B,C,'RowTimes',D,'VariableNames',vnames);
```

- 2 Create a type from `T`.

```
t = coder.typeof(TT);
```

- 3 Pass the type to `fiaccel` by using the `-args` option.

```
fiaccel myFunction -args {t}
```

Provide a Constant Timetable Input

To specify that a timetable input is constant, use `coder.Constant` with the `-args` option:

```
TT = timetable(A,B,C,'RowTimes',D,'VariableNames',vnames);  
fiaccel myFunction -args {coder.Constant(TT)}
```

Representation of Timetables

A coder type object for a timetable describes the object and its properties. Use `coder.typeof` or pass `timetable` as a string scalar to `coder.newtype`.

The coder type object displays a succinct description of the object properties while excluding internal state values. Nonconstant properties display their type and size, while constant properties display only their values. For example:

```
t = timetable((1:5)',(11:15)', 'SampleRate',1);
tType = coder.typeof(t)
```

The representation of variable `t` is stored in coder type object `tType`.

```
tType =
  matlab.coder.type.RegularTimetableType
  5x2 timetable
      Data : 1x2 homogeneous cell
      Description : 1x0 char
      UserData : 0x0 double
      DimensionNames : {'Time'} {'Variables'}
      VariableNames : {'Var1'} {'Var2'}
  VariableDescriptions : 1x2 homogeneous cell
      VariableUnits : 1x2 homogeneous cell
  VariableContinuity : 1x2 matlab.internal.coder.tabular.Continuity
      StartTime : 1x1 matlab.coder.type.DurationType
      SampleRate : 1x1 double
      TimeStep : 1x1 matlab.coder.type.DurationType
```

Define a regular `timetable` by specifying the `SampleRate` or `TimeStep`. You can also define an irregular `timetable` by specifying the `RowTimes`. For example:

```
t1 = timetable((1:3)', 'RowTimes',seconds(1:3));
t1Type = coder.typeof(t)
```

The representation of irregular table `t1` is stored in coder type object `t1Type`.

```
t1Type =
  matlab.coder.type.TimetableType
  3x1 timetable
      Data : 1x1 homogeneous cell
      Description : 1x0 char
      UserData : 0x0 double
      DimensionNames : {'Time'} {'Variables'}
      VariableNames : {'Var1'}
  VariableDescriptions : 1x1 homogeneous cell
      VariableUnits : 1x1 homogeneous cell
  VariableContinuity : 1x1 matlab.internal.coder.tabular.Continuity
      RowTimes : 3x1 matlab.coder.type.DurationType
```

If your workflow requires the legacy representation of coder type objects, use the `getCoderType` function on the variable that has the new representation of your class or object. See “Legacy Representation of Coder Type Objects” (MATLAB Coder).

Resize Object Properties by Using `coder.resize`

You can resize most objects by using `coder.resize`. You can resize objects, its properties and create arrays within the properties.

For a `timetable` coder object, you can resize the object properties:

```
t = timetable((1:5)',(11:15)', 'SampleRate',1);
tType = coder.typeof(t);
tType.UserData = coder.resize(tType.UserData,[10 1],[1 0])
```

This code resizes the `UserData` property to be a `:10x1 double` property. The first dimension is upper-bound at 10.

```
tType =
```

```
matlab.coder.type.RegularTimetableType
5x2 timetable
      Data : 1x2 homogeneous cell
      Description : 1x0 char
      UserData : :10x1 double
      DimensionNames : {'Time'}    {'Variables'}
      VariableNames : {'Var1'}    {'Var2'}
      VariableDescriptions : 1x2 homogeneous cell
      VariableUnits : 1x2 homogeneous cell
      VariableContinuity : 1x2 matlab.internal.coder.tabular.Continuity
      StartTime : 1x1 matlab.coder.type.DurationType
      SampleRate : 1x1 double
      TimeStep : 1x1 matlab.coder.type.DurationType
```

You can also resize the object by using `coder.resize`. See “Edit and Represent Coder Type Objects and Properties” (MATLAB Coder).

See Also

`timetable` | `coder.Constant` | `coder.typeof`

More About

- “Code Generation for Timetables” on page 28-2
- “Timetable Limitations for Code Generation” on page 28-9

Timetable Limitations for Code Generation

If you create timetables, modify them, or use timetable functions in MATLAB code that you intend for code generation, then code generation has limitations described in the next sections. Limitations that apply to classes also apply to timetables. For more information on class limitations, see “MATLAB Classes Definition for Code Generation” (MATLAB Coder).

Creating Timetables Limitations

If your MATLAB code creates timetables, then code generation has these limitations.

| Inputs for Timetable Creation | Limitations |
|-------------------------------------|--|
| Any inputs | <ul style="list-style-type: none"> • The name of the first dimension of a timetable is 'Time' unless you specify it by using the 'DimensionNames' name-value argument. The name of the first dimension is also the name of the vector of row times, which you can refer to by using dot notation. • To create a regular timetable when the 'SampleRate', 'StartTime', or 'TimeStep' name-value arguments are passed in by an entry point function, first use the <code>coder.Constant</code> function to make the values constant. If you do not make them constant, then the row times are considered to be irregular. • If you create a regular timetable, and you attempt to set irregular row times, then an error is produced. • If you create an irregular timetable, then it remains irregular even if you set its sample rate or time step. • Timetable variable names do not have to be valid MATLAB identifiers. The names must be composed of ASCII characters, which can include commas, dashes, and space characters. |
| Timetable created from input arrays | <ul style="list-style-type: none"> • You must specify variables names by using the 'VariableNames' name-value argument when creating timetables from input arrays by using the <code>timetable</code> or <code>array2timetable</code> functions. |

| Inputs for Timetable Creation | Limitations |
|---|--|
| Timetable created with preallocated variables | <ul style="list-style-type: none"> • You do not have to specify the 'VariableNames' argument when you preallocate a timetable by using the <code>timetable</code> function and the 'Size' name-value argument. • You can specify only the following data types by using the 'VariableTypes' name-value argument: <ul style="list-style-type: none"> • 'double' • 'single' • 'doublenan' or 'doubleNaN' • 'singlenan' or 'singleNaN' • 'int8', 'int16', 'int32', or 'int64' • 'uint8', 'uint16', 'uint32', or 'uint64' • 'logical' • 'datetime' • 'duration' • 'cellstr' • 'char' |

Modifying Timetables Limitations

If your MATLAB code modifies data in a timetable, its row times, or its properties, then code generation has these limitations.

| Timetable Operation or Property | Limitations |
|---|--|
| VariableNames, DimensionNames, or UserData properties | <ul style="list-style-type: none"> • After you create a timetable, you cannot change the VariableNames, DimensionNames, or UserData properties. <p>When you create a timetable, you can specify the 'VariableNames', 'DimensionNames', and 'RowTimes' input arguments to set the properties having those names.</p> |

| Timetable Operation or Property | Limitations |
|---|---|
| <p>Timetable indices as input arguments to generated code</p> | <ul style="list-style-type: none"> • To pass timetable indices that specify variables into generated code as input arguments, first use the <code>coder.Constant</code> function to make the indices into the second dimension of the timetable constant. If indices into the second dimension are not constant, then indexing into variables produces an error. • If a timetable has row times that are <code>duration</code> values, and you index into it by using either <code>duration</code> values or an object produced by the <code>timerange</code> or <code>withtol</code> functions, then the output is nonconstant with a variable number of rows. • If a regular timetable has row times that are <code>duration</code> values, and you index into it by using either <code>duration</code> values or an object produced by the <code>timerange</code> or <code>withtol</code> functions, then the output is considered to be irregular. |
| <p>Custom metadata</p> | <ul style="list-style-type: none"> • You cannot add custom metadata to a timetable. The <code>addprop</code> and <code>rmprop</code> functions are not supported. |
| <p>Assignments that change size of timetable</p> | <ul style="list-style-type: none"> • You cannot change the size of a timetable by assignments. For example, this call to add a new row produces an error. <pre data-bbox="906 1192 1442 1360"> function TT = foo() %#codegen TT = timetable((1:3)',(1:3)',... 'RowTimes',seconds([0,5,10]),... 'VariableNames',{'Var1','Var2'}); TT{4,:} = [5,5]; end </pre> <p>Deleting a row or a variable by assignment also produces an error.</p> • You cannot add a new row by using a new row time in an assignment. For example, this call to add a new row by using a new row time instead of a numeric index does not produce an error, but also does not add the new row. <pre data-bbox="906 1654 1442 1822"> function TT = foo() %#codegen TT = timetable((1:3)',(1:3)',... 'RowTimes',seconds([0,5,10]),... 'VariableNames',{'Var1','Var2'}); TT{seconds(15),:} = [5,5]; end </pre> |

| Timetable Operation or Property | Limitations |
|--|---|
| Vertical concatenation | <ul style="list-style-type: none"> When you vertically concatenate timetables, they must have the same variable names in the same order. In MATLAB, the variable names must be the same but can be in different orders in the timetables. |
| Horizontal concatenation | <ul style="list-style-type: none"> When you horizontally concatenate timetables, they must have the same row times in the same order. In MATLAB, the row times must be the same but can be in different orders in the timetables. |
| Timetable variables that are N-D cell arrays | <ul style="list-style-type: none"> If two timetables have variables that are N-D cell arrays, then you cannot vertically concatenate the timetables. You cannot use curly braces to extract data from multiple timetable variables that are N-D cell arrays because this operation is horizontal concatenation. |

Using Timetable Functions Limitations

If your MATLAB code uses the functions listed in the table, then code generation has these limitations.

| Function | Limitations |
|-------------|---|
| convertvars | <ul style="list-style-type: none"> Function handles are not supported. The second and third input arguments (<code>vars</code> and <code>dataType</code>) must be constant. You cannot specify <code>dataType</code> as <code>'cell'</code>, <code>'cellstr'</code>, or <code>'char'</code>. |

| Function | Limitations |
|---|---|
| innerjoin | <ul style="list-style-type: none"> • In general, the input timetables cannot have any nonkey variables with the same names. However, you can join subsets of the input timetables if you specify the 'LeftVariables' and 'RightVariables' name-value arguments. Specify these arguments so that no variable name appears in both 'LeftVariables' and 'RightVariables'. • The values of these name-value arguments must be constant: <ul style="list-style-type: none"> • 'Keys' • 'LeftKeys' • 'RightKeys' • 'LeftVariables' • 'RightVariables' • Nested timetables are not supported. |
| intersect setdiff setxor union | <ul style="list-style-type: none"> • These functions support unsorted timetables in all cases. You do not have to specify the 'stable' option. |
| isregular | <ul style="list-style-type: none"> • Use <code>coder.Constant</code> to make the input argument <code>timeComponent</code> constant. • The input argument <code>timeComponent</code> cannot be a calendar unit. If you specify it, then its value must be 'time'. |
| issortedrows | <ul style="list-style-type: none"> • The input argument <code>vars</code> must be constant. • If any timetable variables have multiple columns, then those variables must have fixed widths. |

| Function | Limitations |
|----------|--|
| join | <ul style="list-style-type: none">• In general, input timetables cannot have nonkey variables with the same names. However, you can join subsets of the input timetables if you specify the name-value arguments:<ul style="list-style-type: none">• 'KeepOneCopy', where you list variables to take from the left input timetable only.• 'LeftVariables' and 'RightVariables', where you list variables to take from either the left input timetable or the right input timetable, but not both.• The values of these name-value arguments must be constant:<ul style="list-style-type: none">• 'Keys'• 'LeftKeys'• 'RightKeys'• 'LeftVariables'• 'RightVariables'• 'KeepOneCopy'• Nested timetables are not supported. |
| movevars | <ul style="list-style-type: none">• The input argument <code>vars</code> cannot contain duplicate variable names. |

| Function | Limitations |
|-----------------------|---|
| outerjoin | <ul style="list-style-type: none"> • Input timetables cannot have key variables with the same names unless the value of 'MergeKeys' is true (logical 1). • In general, the input timetables cannot have any nonkey variables with the same names. However, you can join subsets of the input timetables if you specify the 'LeftVariables' and 'RightVariables' name-value arguments. Specify these arguments so that no variable name appears in both 'LeftVariables' and 'RightVariables'. • The values of these name-value arguments must be constant: <ul style="list-style-type: none"> • 'Keys' • 'LeftKeys' • 'RightKeys' • 'MergeKeys' • 'LeftVariables' • 'RightVariables' • 'Type' • Nested timetables are not supported. |
| retime synchronize | <ul style="list-style-type: none"> • The row times of the output timetable are considered to be irregular, even when synchronized to row times that have a regular time step. • The 'makima' interpolation method is not supported. • If the VariableContinuity properties of the input timetables are not constant, then this function ignores them. • The 'weekly', 'monthly', and 'quarterly' time steps are not supported. <ul style="list-style-type: none"> • If the input timetables have row times that are datetime values, then the 'daily' and 'yearly' time steps also are not supported. |
| sortrows | <ul style="list-style-type: none"> • The input argument vars must be constant. • If tblA has a variable that is a cell array of character vectors with multiple columns, then you cannot sort the timetable using the values in that variable. |

| Function | Limitations |
|-----------|---|
| splitvars | <ul style="list-style-type: none">• The value of the 'NewVariableNames' name-value argument must be constant.• The variables that are split cannot have a variable number of columns. |
| stack | <ul style="list-style-type: none">• The second input argument, vars, must be constant.• The values of the 'ConstantVariables', 'NewDataVariableName', and 'IndexVariableName' name-value arguments must be constant. |
| timerange | <ul style="list-style-type: none">• The input argument unitOfTime is not supported. |

| Function | Limitations |
|----------|---|
| unstack | <ul style="list-style-type: none"> • The 'NewDataVariableNames' name-value argument must be specified. Its value must be constant. • The vars and ivars input arguments (data variables and indicator variables) must be constant. • If you specify grouping variables and constant variables, then they must be constant. • If you specify an aggregation function, then it must be constant. • If the input is a timetable with regular row times and you specify grouping variables that do not include the row times, then the output timetable might have irregular row times. Even though the intervals between output row times might look the same, the output timetable considers the vector of row times to be irregular. • If a variable of the input timetable is a cell array of character vectors, then unstack fills empty cells in the corresponding output variable with 1-by-0 character arrays in the generated code. In MATLAB, unstack fills such gaps with 0-by-0 character arrays. • The unstack function does not support code generation when the input timetable has a variable that is a heterogeneous cell array that cannot be converted to a homogeneous cell array. <ul style="list-style-type: none"> • If the input has a variable that is a homogeneous cell array, or that can be converted to one, then the 'AggregationFunction' name-value argument must be specified. The default value of 'AggregationFunction' is 'unique'. But the unique function does not support cell arrays. |

| Function | Limitations |
|----------|--|
| varfun | <ul style="list-style-type: none"> • The function handle input, <code>func</code>, must be constant. • While function handles can be inputs to <code>varfun</code> itself, they cannot be inputs to your entry point functions. Specify <code>func</code> within the code meant for code generation. For more information, see “Function Handle Limitations for Code Generation” (MATLAB Coder). • The values for all name-value arguments must be constant. • The <code>'ErrorHandler'</code> name-value argument is not supported for code generation. • Variable-size input arguments are not supported. • If you specify <code>'GroupingVariables'</code>, then the output is always an irregular timetable. • Grouping variables cannot have duplicate values in generated code. • You cannot specify the value of <code>'OutputFormat'</code> as <code>'cell'</code> if you specify the <code>'GroupingVariables'</code> name-value arguments and the function returns a different data type for each variable specified by <code>'InputVariables'</code>. • If you specify groups and the number of groups is not known at compile-time, and that number turns out to be zero, then empty double variables in the output might have sizes of 1-by-0 in generated code. In MATLAB, such variables have sizes of 0-by-0. |

See Also

`array2timetable` | `table2timetable` | `timetable`

More About

- “Code Generation for Timetables” on page 28-2
- “Define Timetable Inputs” on page 28-6

Code Generation for Variable-Size Data

- “Code Generation for Variable-Size Arrays” on page 29-2
- “Control Memory Allocation for Variable-Size Arrays” on page 29-4
- “Control Dynamic Memory Allocation for Fixed-Size Arrays” on page 29-6
- “Specify Upper Bounds for Variable-Size Arrays” on page 29-8
- “Define Variable-Size Data for Code Generation” on page 29-10
- “Diagnose and Fix Variable-Size Data Errors” on page 29-15
- “Incompatibilities with MATLAB in Variable-Size Support for Code Generation” on page 29-18
- “Variable-Sizing Restrictions for Code Generation of Toolbox Functions” on page 29-25
- “Generate Code With Implicit Expansion Enabled” on page 29-30
- “Optimize Implicit Expansion in Generated Code” on page 29-34
- “Representation of Arrays in Generated Code” on page 29-39
- “Control Memory Allocation for Fixed-Size Arrays” on page 29-43
- “Resolve Error: Size Mismatches” on page 29-45

Code Generation for Variable-Size Arrays

For code generation, an array dimension is fixed-size or variable-size. If the code generator can determine the size of the dimension and that the size of the dimension does not change, then the dimension is fixed-size. When all dimensions of an array are fixed-size, the array is a fixed-size array. In the following example, *Z* is a fixed-size array.

```
function Z = myfcn()
Z = zeros(1,4);
end
```

The size of the first dimension is 1 and the size of the second dimension is 4.

If the code generator cannot determine the size of a dimension or the code generator determines that the size changes, then the dimension is variable-size. When at least one of its dimensions is variable-size, an array is a variable-size array.

A variable-size dimension is either bounded or unbounded. A bounded dimension has a fixed upper size. An unbounded dimension does not have a fixed upper size.

In the following example, the second dimension of *Z* is bounded, variable-size. It has an upper bound of 16.

```
function s = myfcn(n)
if (n > 0)
    Z = zeros(1,4);
else
    Z = zeros(1,16);
end
s = length(Z);
```

In the following example, if the value of *n* is unknown at compile time, then the second dimension of *Z* is unbounded.

```
function s = myfcn(n)
Z = rand(1,n);
s = sum(Z);
end
```

You can define variable-size arrays by:

- Using constructors, such as `zeros`, with a nonconstant dimension
- Assigning multiple, constant sizes to the same variable before using it
- Declaring all instances of a variable to be variable-size by using `coder.varsize`

For more information, see “Define Variable-Size Data for Code Generation” on page 29-10.

You can control whether variable-size arrays are allowed for code generation. See “Enabling and Disabling Support for Variable-Size Arrays” on page 29-3.

Memory Allocation for Variable-Size Arrays

For fixed-size arrays and variable-size arrays whose size is less than a threshold, the code generator allocates memory statically on the stack. For unbounded, variable-size arrays and variable-size arrays

whose size is greater than or equal to a threshold, the code generator allocates memory dynamically on the heap.

You can control whether dynamic memory allocation is allowed or when it is used for code generation. See “Control Memory Allocation for Variable-Size Arrays” on page 29-4.

The code generator represents dynamically allocated data as a structure type called `emxArray`. The code generator generates utility functions that create and interact with `emxArrays`. If you use Embedded Coder, you can customize the generated identifiers for the `emxArray` types and utility functions. See “Identifier Format Control” (Embedded Coder).

Enabling and Disabling Support for Variable-Size Arrays

By default, support for variable-size arrays is enabled. To modify this support:

- In a code configuration object, set the `EnableVariableSizing` parameter to `true` or `false`.

Variable-Size Arrays in a Code Generation Report

You can tell whether an array is fixed-size or variable-size by looking at the **Size** column of the **Variables** tab in a code generation report.

| Name | Type | Size | Class |
|------|--------|---------|--------|
| y | Output | 1 × 1 | double |
| A | Input | 1 × :16 | char |
| n | Input | 1 × 1 | double |
| X | Local | 1 × :? | double |

A colon (:) indicates that a dimension is variable-size. A question mark (?) indicates that the size is unbounded. For example, a size of 1-by-:? indicates that the size of the first dimension is fixed-size 1 and the size of the second dimension is unbounded, variable-size. Italics indicates that the code generator produced a variable-size array, but the size of the array does not change during execution.

| Name | Type | Size | Class |
|------|--------|--------------|--------|
| y | Output | 1 × :5 | double |
| n | Input | 1 × 1 | double |
| Z | Local | <i>1 × 4</i> | double |

See Also

More About

- “Control Memory Allocation for Variable-Size Arrays” on page 29-4
- “Specify Upper Bounds for Variable-Size Arrays” on page 29-8
- “Define Variable-Size Data for Code Generation” on page 29-10

Control Memory Allocation for Variable-Size Arrays

Dynamic memory allocation allocates memory on the heap as needed at run-time, instead of allocating memory statically on the stack. Dynamic memory allocation is beneficial when:

- You do not know the upper bound of an array.
- You do not want to allocate memory on the stack for large arrays.

Dynamic memory allocation and the freeing of this memory can result in slower execution of the generated code. To control the use of dynamic memory allocation for variable-size arrays, you can:

- Provide upper bounds for variable-size arrays on page 29-4.
- Disable dynamic memory allocation on page 29-4.
- Configure the code generator to use dynamic memory allocation for arrays bigger than a threshold on page 29-4.

Provide Upper Bounds for Variable-Size Arrays

For an unbounded variable-size array, the code generator allocates memory dynamically on the heap. For a variable-size array with upper bound, whose size, in bytes, is less than the dynamic memory allocation threshold, the code generator allocates memory statically on the stack. To prevent dynamic memory allocation:

- 1 Specify upper bounds for a variable-size array. See “Specify Upper Bounds for Variable-Size Arrays” on page 29-8.
- 2 Make sure that the size of the array, in bytes, is less than the dynamic memory allocation threshold. See “Configure Code Generator to Use Dynamic Memory Allocation for Arrays Bigger Than a Threshold” on page 29-4.

Disable Dynamic Memory Allocation

By default, dynamic memory allocation for variable-size arrays is enabled. To disable it, in a configuration object for fixed-point acceleration, set the `EnableDynamicMemoryAllocation` parameter to `false`.

If you disable dynamic memory allocation, you must provide upper bounds for variable-size arrays.

Configure Code Generator to Use Dynamic Memory Allocation for Arrays Bigger Than a Threshold

Instead of disabling dynamic memory allocation for all variable-size arrays, you can specify for which size arrays the code generator uses dynamic memory allocation.

Use the dynamic memory allocation threshold to:

- Disable dynamic memory allocation for smaller arrays. For smaller arrays, static memory allocation can speed up generated code. However, static memory allocation can lead to unused storage space. You can decide that the unused storage space is not a significant consideration for smaller arrays.
- Enable dynamic memory allocation for larger arrays. For larger arrays, when you use dynamic memory allocation, you can significantly reduce storage requirements.

The default dynamic memory allocation threshold is 64 kilobytes. To change the threshold, in a configuration object for fixed-point acceleration, set the `EnableDynamicMemoryAllocation` to `true` and set a value for `DynamicMemoryAllocationThreshold`.

See Also

More About

- “Code Generation for Variable-Size Arrays” on page 29-2

Control Dynamic Memory Allocation for Fixed-Size Arrays

The code generated by MATLAB Coder allocates memory to the program stack for fixed-size arrays, whose size is less than a threshold. For example, in the following code, the array Z is a fixed-size array with a first dimension of size 1 and a second dimension of size 4.

```
function Z = myfcn()  
Z = zeros(1,4);  
end
```

Dynamic memory allocation allocates memory on the heap for fixed-size arrays, instead of the program stack. Consider dynamically allocating fixed-size arrays when you have large arrays that could exhaust stack memory.

Dynamic memory allocation might result in slower execution of the generated code.

Enable Dynamic Memory Allocation for Fixed-Size Arrays

By default, dynamic memory allocation for fixed-size arrays is disabled. To enable it:

- In a configuration object for code generation, set `EnableDynamicMemoryAllocation` and `DynamicMemoryAllocationForFixedSizeArrays` options to `true`.
- In the MATLAB Coder app, in the **Memory** settings, select **Enable dynamic memory allocation** and **Enable dynamic memory allocation for fixed-sized arrays**.

Dynamic Memory Allocation Threshold for Fixed-Size Arrays

When `DynamicMemoryAllocationForFixedSizeArrays` is enabled, the code generator allocates memory dynamically on the heap for fixed-size arrays whose size (in bytes) is

greater than or equal to `DynamicMemoryAllocationThreshold`.

The default dynamic memory allocation threshold is 64 kilobytes. To configure the threshold:

- In a configuration object for code generation, set a value for the `DynamicMemoryAllocationThreshold`.
- In the MATLAB Coder app, in the **Memory** settings, set a value for **Dynamic memory allocation threshold**.

Generating Code for Fixed-Size Arrays

Consider the following MATLAB function which calculates the product of two large fixed-size arrays a and b.

```
function y = tlargeSize(a,b)  
y = a*b;  
end
```

To generate C code, run these commands:

```
cfg = coder.config('lib');  
cfg.VerificationMode="SIL";  
cfg.DynamicMemoryAllocationForFixedSizeArrays = true;
```

```
t=coder.typeof(0,[1e4 1e4]);
codegen tlargeSize -args {t,t} -config cfg - report
```

By enabling the `DynamicMemoryAllocationForFixedSizeArrays` option, arrays `a` and `b` are dynamically allocated on the heap preventing stack overflow.

Generated Code

```
...
void tlargeSize(const mxArray_real_T *a, const mxArray_real_T *b,
               mxArray_real_T *y)
{
    const double *a_data;
    const double *b_data;
    double *y_data;
    ...
}
```

Usage Notes and Limitations

When enabling `DynamicMemoryAllocationForFixedSizeArrays`:

- Input and output variables of the function changes to:
 - `coder::array` for C++ code generation. See, “Use Dynamically Allocated C++ Arrays in Generated Function Interfaces” (MATLAB Coder).
 - `mxArray` for C code generation. See, “Use C Arrays in the Generated Function Interfaces” (MATLAB Coder)
- Ensure C/C++ function signatures are updated when passing struct with fixed-size arrays to `coder.ceval` or Code Replacement Library (CRL).

Note Enabling `DynamicMemoryAllocationForFixedSizeArrays` is not supported for GPU code generation.

See Also

`coder.EmbeddedCodeConfig` | `coder.MexCodeConfig` | `coder.CodeConfig`

Related Examples

- “Control Memory Allocation for Variable-Size Arrays” (MATLAB Coder)
- “Representation of Arrays in Generated Code” (MATLAB Coder)
- “Use C Arrays in the Generated Function Interfaces” (MATLAB Coder)
- “Use Dynamically Allocated C++ Arrays in Generated Function Interfaces” (MATLAB Coder)

Specify Upper Bounds for Variable-Size Arrays

Specify upper bounds for an array when:

- Dynamic memory allocation is disabled.

If dynamic memory allocation is disabled, you must specify upper bounds for all arrays.

- You do not want the code generator to use dynamic memory allocation for the array.

Specify upper bounds that result in an array size (in bytes) that is less than the dynamic memory allocation threshold.

Specify Upper Bounds for Variable-Size Inputs

This command specifies that the input to function `foo` is a matrix of real doubles with two variable dimensions. The upper bound for the first dimension is 3. The upper bound for the second dimension is 100.

To specify upper bounds for variable-size inputs, use the `coder.typeof` construct with the `fiaccel -args` option. For example:

```
fiaccel foo -args {coder.typeof(fi(0),[3 100],1)}
```

This command specifies that the input to function `foo` is a matrix of `fi` types with two variable dimensions. The upper bound for the first dimension is 3. The upper bound for the second dimension is 100.

Specify Upper Bounds for Local Variables

When using static allocation, the code generator uses a sophisticated analysis to calculate the upper bounds of local data. However, when the analysis fails to detect an upper bound or calculates an upper bound that is not precise enough for your application, you must specify upper bounds explicitly for local variables.

Constrain the Value of Variables That Specify the Dimensions of Variable-Size Arrays

To constrain the value of variables that specify the dimensions of variable-size arrays, use the `assert` function with relational operators. For example:

```
function y = dim_need_bound(n) %#codegen
assert (n <= 5);
L = ones(n,n);
M = zeros(n,n);
M = [L; M];
y = M;
```

This `assert` statement constrains input `n` to a maximum size of 5. `L` is variable-size with upper bounds of 5 in each dimension. `M` is variable-size with an upper bound of 10 in the first dimension and 5 in the second dimension.

Specify the Upper Bounds for All Instances of a Local Variable

To specify the upper bounds for all instances of a local variable in a function, use the `coder.varsize` function. For example:

```
function Y = example_bounds1(u) %#codegen
Y = [1 2 3 4 5];
coder.varsize('Y',[1 10]);
if (u > 0)
    Y = [Y Y+u];
else
    Y = [Y Y*u];
end
```

The second argument of `coder.varsize` specifies the upper bound for each instance of the variable specified in the first argument. In this example, the argument `[1 10]` indicates that for every instance of `Y`:

- The first dimension is fixed at size 1.
- The second dimension can grow to an upper bound of 10.

See Also

`coder.varsize` | `coder.typeof`

More About

- “Code Generation for Variable-Size Arrays” on page 29-2
- “Define Variable-Size Data for Code Generation” on page 29-10

Define Variable-Size Data for Code Generation

For code generation, before using variables in operations or returning them as outputs, you must assign them a specific class, size, and complexity. Generally, after the initial assignment, you cannot reassign variable properties. Therefore, after assigning a fixed size to a variable or structure field, attempts to grow the variable or structure field might cause a compilation error. In these cases, you must explicitly define the data as variable-size by using one of these methods.

| Method | See |
|--|---|
| Assign the data from a variable-size matrix constructor such as: <ul style="list-style-type: none"> • ones • zeros • repmat | “Use a Matrix Constructor with Nonconstant Dimensions” on page 29-10 |
| Assign multiple, constant sizes to the same variable before using (reading) the variable. | “Assign Multiple Sizes to the Same Variable” on page 29-10 |
| Grow an array by using (end + 1) indexing. | “Growing an Array by Using (end + 1)” on page 29-11 |
| Define all instances of a variable to be variable-size. | “Define Variable-Size Data Explicitly by Using coder.varsize” on page 29-12 |

Use a Matrix Constructor with Nonconstant Dimensions

You can define a variable-size matrix by using a constructor with nonconstant dimensions. For example:

```
function s = var_by_assign(u) %#codegen
y = ones(3,u);
s = numel(y);
```

If you are not using dynamic memory allocation, you must also add an `assert` statement to provide upper bounds for the dimensions. For example:

```
function s = var_by_assign(u) %#codegen
assert (u < 20);
y = ones(3,u);
s = numel(y);
```

Assign Multiple Sizes to the Same Variable

Before you use (read) a variable in your code, you can make it variable-size by assigning multiple, constant sizes to it. When the code generator uses static allocation on the stack, it infers the upper bounds from the largest size specified for each dimension. When you assign the same size to a given dimension across all assignments, the code generator assumes that the dimension is fixed at that size. The assignments can specify different shapes and sizes.

When the code generator uses dynamic memory allocation, it does not check for upper bounds. It assumes that the variable-size data is unbounded.

Inferring Upper Bounds from Multiple Definitions with Different Shapes

```
function s = var_by_multiassign(u) %#codegen
if (u > 0)
    y = ones(3,4,5);
else
    y = zeros(3,1);
end
s = numel(y);
```

When the code generator uses static allocation, it infers that `y` is a matrix with three dimensions:

- The first dimension is fixed at size 3
- The second dimension is variable-size with an upper bound of 4
- The third dimension is variable-size with an upper bound of 5

When the code generator uses dynamic allocation, it analyzes the dimensions of `y` differently:

- The first dimension is fixed at size 3.
- The second and third dimensions are unbounded.

Growing an Array by Using (`end + 1`)

To grow an array `X`, you can assign a value to `X(end + 1)`. If you make this assignment in your MATLAB code, the code generator treats the dimension you grow as variable-size.

For example, you can generate code for this code snippet:

```
...
a = [1 2 3 4 5 6];
a(end + 1) = 7;

b = [1 2];
for i = 3:10
    b(end + 1) = i;
end
...
```

When you use (`end + 1`) to grow an array, follow these restrictions:

- Use only (`end + 1`). Do not use (`end + 2`), (`end + 3`), and so on.
- Use (`end + 1`) with vectors only. For example, the following code is not allowed because `X` is a matrix, not a vector.

```
...
X = [1 2; 3 4];
X(end + 1) = 5;
...
```

- You can grow empty arrays of size `1x0` by using (`end + 1`). Growing arrays of size `0x1` is not supported. Growing an array of size `0x0` is supported only if you create that array by using `[]`.

Define Variable-Size Data Explicitly by Using `coder.varsize`

To explicitly define variable-size data, use the function `coder. varsize`. Optionally, you can also specify which dimensions vary along with their upper bounds. For example:

- Define B as a variable-size 2-dimensional array, where each dimension has an upper bound of 64.

```
coder. varsize('B', [64 64]);
```

- Define B as a variable-size array:

```
coder. varsize('B');
```

When you supply only the first argument, `coder. varsize` assumes that all dimensions of B can vary and that the upper bound is `size(B)`.

Specify Which Dimensions Vary

You can use the function `coder. varsize` to specify which dimensions vary. For example, the following statement defines B as an array whose first dimension is fixed at 2, but whose second dimension can grow to a size of 16:

```
coder. varsize('B',[2, 16],[0 1])
```

.

The third argument specifies which dimensions vary. This argument must be a logical vector or a double vector containing only zeros and ones. Dimensions that correspond to zeros or `false` have fixed size. Dimensions that correspond to ones or `true` vary in size. `coder. varsize` usually treats dimensions of size 1 as fixed. See “Define Variable-Size Matrices with Singleton Dimensions” on page 29-12.

Allow a Variable to Grow After Defining Fixed Dimensions

Function `var_by_if` defines matrix Y with fixed 2-by-2 dimensions before the first use (where the statement `Y = Y + u` reads from Y). However, `coder. varsize` defines Y as a variable-size matrix, allowing it to change size based on decision logic in the `else` clause:

```
function Y = var_by_if(u) %#codegen
if (u > 0)
    Y = zeros(2,2);
    coder. varsize('Y');
    if (u < 10)
        Y = Y + u;
    end
else
    Y = zeros(5,5);
end
```

Without `coder. varsize`, the code generator infers Y to be a fixed-size, 2-by-2 matrix. It generates a size mismatch error.

Define Variable-Size Matrices with Singleton Dimensions

A singleton dimension is a dimension for which `size(A, dim) = 1`. Singleton dimensions are fixed in size when:

- You specify a dimension with an upper bound of 1 in `coder. varsize` expressions.

For example, in this function, `Y` behaves like a vector with one variable-size dimension:

```
function Y = dim_singleton(u) %#codegen
Y = [1 2];
coder. varsize('Y', [1 10]);
if (u > 0)
    Y = [Y 3];
else
    Y = [Y u];
end
```

- You initialize variable-size data with singleton dimensions by using matrix constructor expressions or matrix functions.

For example, in this function, `X` and `Y` behave like vectors where only their second dimensions are variable-size.

```
function [X,Y] = dim_singleton_vects(u) %#codegen
Y = ones(1,3);
X = [1 4];
coder. varsize('Y','X');
if (u > 0)
    Y = [Y u];
else
    X = [X u];
end
```

You can override this behavior by using `coder. varsize` to specify explicitly that singleton dimensions vary. For example:

```
function Y = dim_singleton_vary(u) %#codegen
Y = [1 2];
coder. varsize('Y', [1 10], [1 1]);
if (u > 0)
    Y = [Y Y+u];
else
    Y = [Y Y*u];
end
```

In this example, the third argument of `coder. varsize` is a vector of ones, indicating that each dimension of `Y` varies in size.

Define Variable-Size Structure Fields

To define structure fields as variable-size arrays, use a colon (`:`) as the index expression. The colon (`:`) indicates that all elements of the array are variable-size. For example:

```
function y=struct_example() %#codegen

d = struct('values', zeros(1,0), 'color', 0);
data = repmat(d, [3 3]);
coder. varsize('data(:).values');

for i = 1:numel(data)
    data(i).color = rand-0.5;
    data(i).values = 1:i;
end
```

```
end

y = 0;
for i = 1:numel(data)
    if data(i).color > 0
        y = y + sum(data(i).values);
    end
end
```

The expression `coder. varsize('data(:).values')` defines the field `values` inside each element of matrix `data` to be variable-size.

Here are other examples:

- `coder. varsize('data.A(:).B')`

In this example, `data` is a scalar variable that contains matrix `A`. Each element of matrix `A` contains a variable-size field `B`.

- `coder. varsize('data(:).A(:).B')`

This expression defines field `B` inside each element of matrix `A` inside each element of matrix `data` to be variable-size.

See Also

`coder. varsize` | `coder. typeof`

More About

- “Code Generation for Variable-Size Arrays” on page 29-2
- “Specify Upper Bounds for Variable-Size Arrays” on page 29-8

Diagnose and Fix Variable-Size Data Errors

Diagnosing and Fixing Size Mismatch Errors

Issue: Assigning Variable-Size Matrices to Fixed-Size Matrices

You cannot assign variable-size matrices to fixed-size matrices in generated code. Consider this example:

```
function Y = example_mismatch1(n) %#codegen
assert(n < 10);
B = ones(n,n);
A = magic(3);
A(1) = mean(A(:));
if (n == 3)
    A = B;
end
Y = A;
```

Compiling this function produces this error:

```
??? Dimension 1 is fixed on the left-hand side
but varies on the right ...
```

There are several ways to fix this error:

- Allow matrix A to grow by adding the `coder.varsize` construct:

```
function Y = example_mismatch1_fix1(n) %#codegen
coder.varsize('A');
assert(n < 10);
B = ones(n,n);
A = magic(3);
A(1) = mean(A(:));
if (n == 3)
    A = B;
end
Y = A;
```

- Explicitly restrict the size of matrix B to 3-by-3 by modifying the `assert` statement:

```
function Y = example_mismatch1_fix2(n) %#codegen
coder.varsize('A');
assert(n == 3)
B = ones(n,n);
A = magic(3);
A(1) = mean(A(:));
if (n == 3)
    A = B;
end
Y = A;
```

- Use explicit indexing to make B the same size as A:

```
function Y = example_mismatch1_fix3(n) %#codegen
assert(n < 10);
B = ones(n,n);
A = magic(3);
```

```
A(1) = mean(A(:));
if (n == 3)
    A = B(1:3, 1:3);
end
Y = A;
```

Issue: Empty Matrix Reshaped to Match Variable-Size Specification

If you assign an empty matrix `[]` to variable-size data, MATLAB might silently reshape the data in generated code to match a `coder. varsize` specification. For example:

```
function Y = test(u) %#codegen
Y = [];
coder. varsize('Y', [1 10]);
if u < 0
    Y = [Y u];
end
```

In this example, `coder. varsize` defines `Y` as a column vector of up to 10 elements, so its first dimension is fixed at size 1. The statement `Y = []` designates the first dimension of `Y` as 0, creating a mismatch. The right hand side of the assignment is an empty matrix and the left hand side is a variable-size vector. In this case, MATLAB reshapes the empty matrix `Y = []` in generated code to `Y = zeros(1,0)` so it matches the `coder. varsize` specification.

Issue: Assigning Implicitly Expanded Outputs to Fixed-Size Variable

If you assign the implicitly expanded output of a binary operation or function to a variable of different size, the code generator might produce an error. For example:

```
function out = test(n) %#codegen
x = ones(n,1);
if mod(n,2) == 1
    y = ones(n,n);
    x = y + x;
end
out = out + x(2);
end
```

In this example, `x` is an unbounded vector. Due to implicit expansion, the plus operation on `x` and `y` results in an unbounded matrix (Inf-by-Inf). Assigning an unbounded matrix to `x`, which is an unbounded vector, results in an error.

If you want to use the implicitly expanded output, assign the output to a new variable with the same size as the output.

If you want `x` to retain its size and not apply implicit expansion in the generated code, use `coder. sameSizeBinaryOp` to apply the operation. You can also call `coder. noImplicitExpansionInFunction` in your function body to disable implicit expansion in the code generated for that function.

Implicit expansion automatically expands the operands to apply binary operations on arrays of compatible sizes. See “Generate Code With Implicit Expansion Enabled” (MATLAB Coder), “Optimize Implicit Expansion in Generated Code” (MATLAB Coder), and “Compatible Array Sizes for Basic Operations”.

Diagnosing and Fixing Errors in Detecting Upper Bounds

Issue: Using Nonconstant Dimensions in a Matrix Constructor

You can define variable-size data by assigning a variable to a matrix with nonconstant dimensions. For example:

```
function y = dims_vary(u) %#codegen
if (u > 0)
    y = ones(3,u);
else
    y = zeros(3,1);
end
```

However, compiling this function generates an error because you did not specify an upper bound for `u`.

There are several ways to fix the problem:

- Enable dynamic memory allocation and recompile. During code generation, MATLAB does not check for upper bounds when it uses dynamic memory allocation for variable-size data.
- If you do not want to use dynamic memory allocation, add an `assert` statement before the first use of `u`:

```
function y = dims_vary_fix(u) %#codegen
assert (u < 20);
if (u > 0)
    y = ones(3,u);
else
    y = zeros(3,1);
end
```

Incompatibilities with MATLAB in Variable-Size Support for Code Generation

In this section...

“Incompatibility with MATLAB for Scalar Expansion” on page 29-18

“Incompatibility with MATLAB in Determining Size of Variable-Size N-D Arrays” on page 29-19

“Incompatibility with MATLAB in Determining Size of Empty Arrays” on page 29-19

“Incompatibility with MATLAB in Determining Class of Empty Arrays” on page 29-21

“Incompatibility with MATLAB in Matrix-Matrix Indexing” on page 29-21

“Incompatibility with MATLAB in Vector-Vector Indexing” on page 29-22

“Incompatibility with MATLAB in Matrix Indexing Operations for Code Generation” on page 29-22

“Incompatibility with MATLAB in Concatenating Variable-Size Matrices” on page 29-23

“Differences When Curly-Brace Indexing of Variable-Size Cell Array Inside Concatenation Returns No Elements” on page 29-23

Incompatibility with MATLAB for Scalar Expansion

Scalar expansion is a method of converting scalar data to match the dimensions of vector or matrix data. If one operand is a scalar and the other is not, scalar expansion applies the scalar to every element of the other operand.

During code generation, scalar expansion rules apply except when operating on two variable-size expressions. In this case, both operands must be the same size. The generated code does not perform scalar expansion even if one of the variable-size expressions turns out to be scalar at run time. Therefore, when run-time error checks are enabled, a run-time error can occur.

Consider this function:

```
function y = scalar_exp_test_err1(u) %#codegen
y = ones(3);
switch u
    case 0
        z = 0;
    case 1
        z = 1;
    otherwise
        z = zeros(3);
end
y(:) = z;
```

When you generate code for this function, the code generator determines that `z` is variable size with an upper bound of 3.

| Summary | All Messages (0) | Build Logs | Code Insights (1) | Variables |
|---------|------------------|------------|-------------------|-----------|
| Name | | Type | Size | Class |
| y | | Output | 3 × 3 | double |
| u | | Input | 1 × 1 | double |
| z | | Local | :3 × :3 | double |

If you run the MEX function with `u` equal to 0 or 1, the generated code does not perform scalar expansion, even though `z` is scalar at run time. Therefore, when run-time error checks are enabled, a run-time error can occur.

```
scalar_exp_test_err1_mex(0)
Subscripted assignment dimension mismatch: [9] ~= [1].

Error in scalar_exp_test_err1 (line 11)
y(:) = z;
```

To avoid this issue, use indexing to force `z` to be a scalar value.

```
function y = scalar_exp_test_err1(u) %#codegen
y = ones(3);
switch u
    case 0
        z = 0;
    case 1
        z = 1;
    otherwise
        z = zeros(3);
end
y(:) = z(1);
```

Incompatibility with MATLAB in Determining Size of Variable-Size N-D Arrays

For variable-size N-D arrays, the `size` function can return a different result in generated code than in MATLAB. In generated code, `size(A)` returns a fixed-length output because it does not drop trailing singleton dimensions of variable-size N-D arrays. By contrast, `size(A)` in MATLAB returns a variable-length output because it drops trailing singleton dimensions.

For example, if the shape of array `A` is `?x:?x:?` and `size(A,3)==1`, `size(A)` returns:

- Three-element vector in generated code
- Two-element vector in MATLAB code

Workarounds

If your application requires generated code to return the same size of variable-size N-D arrays as MATLAB code, consider one of these workarounds:

- Use the two-argument form of `size`.

For example, `size(A,n)` returns the same answer in generated code and MATLAB code.

- Rewrite `size(A)`:

```
B = size(A);
X = B(1:ndims(A));
```

This version returns `X` with a variable-length output. However, you cannot pass a variable-size `X` to matrix constructors such as `zeros` that require a fixed-size argument.

Incompatibility with MATLAB in Determining Size of Empty Arrays

The size of an empty array in generated code might be different from its size in MATLAB source code. The size might be 1×0 or 0×1 in generated code, but 0×0 in MATLAB. Therefore, you should not write code that relies on the specific size of empty matrices.

For example, consider the following code:

```
function y = foo(n) %#codegen
x = [];
i = 0;
while (i < 10)
    x = [5 x];
    i = i + 1;
end
if n > 0
    x = [];
end
y = size(x);
end
```

Concatenation requires its operands to match on the size of the dimension that is not being concatenated. In the preceding concatenation, the scalar value has size 1×1 and x has size 0×0 . To support this use case, the code generator determines the size for x as $[1 \times ?]$. Because there is another assignment $x = []$ after the concatenation, the size of x in the generated code is 1×0 instead of 0×0 .

This behavior persists while determining the size of empty character vectors which are denoted as `''`. For example, consider the following code:

```
function out = string_size
out = size('');
end
```

Here, the value of `out` might be 1×0 or 0×1 in generated code, but 0×0 in MATLAB.

For incompatibilities with MATLAB in determining the size of an empty array that results from deleting elements of an array, see “Size of Empty Array That Results from Deleting Elements of an Array” on page 19-17.

Workaround

If your application checks whether a matrix is empty, use one of these workarounds:

- Rewrite your code to use the `isempty` function instead of the `size` function.
- Instead of using `x=[]` to create empty arrays, create empty arrays of a specific size using `zeros`. For example:

```
function y = test_empty(n) %#codegen
x = zeros(1,0);
i=0;
while (i < 10)
    x = [5 x];
    i = i + 1;
end
if n > 0
    x = zeros(1,0);
end
```

```
y=size(x);
end
```

Incompatibility with MATLAB in Determining Class of Empty Arrays

The class of an empty array in generated code can be different from its class in MATLAB source code. Therefore, do not write code that relies on the class of empty matrices.

For example, consider the following code:

```
function y = fun(n)
x = [];
if n > 1
    x = ['a' x];
end
y=class(x);
end
```

`fun(0)` returns `double` in MATLAB, but `char` in the generated code. When the statement `n > 1` is false, MATLAB does not execute `x = ['a' x]`. The class of `x` is `double`, the class of the empty array. However, the code generator considers all execution paths. It determines that based on the statement `x = ['a' x]`, the class of `x` is `char`.

Workaround

Instead of using `x=[]` to create an empty array, create an empty array of a specific class. For example, use `blanks(0)` to create an empty array of characters.

```
function y = fun(n)
x = blanks(0);
if n > 1
    x = ['a' x];
end
y=class(x);
end
```

Incompatibility with MATLAB in Matrix-Matrix Indexing

In matrix-matrix indexing, you use one matrix to index into another matrix. In MATLAB, the general rule for matrix-matrix indexing is that the size and orientation of the result match the size and orientation of the index matrix. For example, if `A` and `B` are matrices, `size(A(B))` equals `size(B)`. When `A` and `B` are vectors, MATLAB applies a special rule. The special vector-vector indexing rule is that the orientation of the result is the orientation of the data matrix. For example, if `A` is 1-by-5 and `B` is 3-by-1, then `A(B)` is 1-by-3.

The code generator applies the same matrix-matrix indexing rules as MATLAB. If `A` and `B` are variable-size matrices, to apply the matrix-matrix indexing rules, the code generator assumes that `size(A(B))` equals `size(B)`. If, at run time, `A` and `B` become vectors and have different orientations, then the assumption is incorrect. Therefore, when run-time error checks are enabled, an error can occur.

To avoid this issue, force your data to be a vector by using the colon operator for indexing. For example, suppose that your code intentionally toggles between vectors and regular matrices at run time. You can do an explicit check for vector-vector indexing.

```

...
if isvector(A) && isvector(B)
    C = A(:);
    D = C(B(:));
else
    D = A(B);
end
...

```

The indexing in the first branch specifies that `C` and `B(:)` are compile-time vectors. Therefore, the code generator applies the indexing rule for indexing one vector with another vector. The orientation of the result is the orientation of the data vector, `C`.

Incompatibility with MATLAB in Vector-Vector Indexing

In MATLAB, the special rule for vector-vector indexing is that the orientation of the result is the orientation of the data vector. For example, if `A` is 1-by-5 and `B` is 3-by-1, then `A(B)` is 1-by-3. If, however, the data vector `A` is a scalar, then the orientation of `A(B)` is the orientation of the index vector `B`.

The code generator applies the same vector-vector indexing rules as MATLAB. If `A` and `B` are variable-size vectors, to apply the indexing rules, the code generator assumes that the orientation of `B` matches the orientation of `A`. At run time, if `A` is scalar and the orientation of `A` and `B` do not match, then the assumption is incorrect. Therefore, when run-time error checks are enabled, a run-time error can occur.

To avoid this issue, make the orientations of the vectors match. Alternatively, index single elements by specifying the row and column. For example, `A(row, column)`.

Incompatibility with MATLAB in Matrix Indexing Operations for Code Generation

The following limitation applies to matrix indexing operations for code generation:

- Initialization of the following style:

```

for i = 1:10
    M(i) = 5;
end

```

In this case, the size of `M` changes as the loop is executed. Code generation does not support increasing the size of an array over time.

For code generation, preallocate `M`.

```

M = zeros(1,10);
for i = 1:10
    M(i) = 5;
end

```

The following limitation applies to matrix indexing operations for code generation when dynamic memory allocation is disabled:

- $M(i:j)$ where i and j change in a loop

During code generation, memory is not dynamically allocated for the size of the expressions that change as the program executes. To implement this behavior, use `for`-loops as shown:

```

...
M = ones(10,10);
for i=1:10
    for j = i:10
        M(i,j) = 2*M(i,j);
    end
end
...

```

Note The matrix M must be defined before entering the loop.

Incompatibility with MATLAB in Concatenating Variable-Size Matrices

For code generation, when you concatenate variable-size arrays, the dimensions that are not being concatenated must match exactly.

Differences When Curly-Brace Indexing of Variable-Size Cell Array Inside Concatenation Returns No Elements

Suppose that:

- c is a variable-size cell array.
- You access the contents of c by using curly braces. For example, $c\{2:4\}$.
- You include the results in concatenation. For example, $[a \ c\{2:4\} \ b]$.
- $c\{I\}$ returns no elements. Either c is empty or the indexing inside the curly braces produces an empty result.

For these conditions, MATLAB omits $c\{I\}$ from the concatenation. For example, $[a \ c\{I\} \ b]$ becomes $[a \ b]$. The code generator treats $c\{I\}$ as the empty array $[c\{I\}]$. The concatenation becomes $[\dots [c\{i\}] \dots]$. This concatenation then omits the array $[c\{I\}]$. So that the properties of $[c\{I\}]$ are compatible with the concatenation $[\dots [c\{i\}] \dots]$, the code generator assigns the class, size, and complexity of $[c\{I\}]$ according to these rules:

- The class and complexity are the same as the base type of the cell array.
- The size of the second dimension is always 0.
- For the rest of the dimensions, the size of N_i depends on whether the corresponding dimension in the base type is fixed or variable size.
 - If the corresponding dimension in the base type is variable size, the dimension has size 0 in the result.
 - If the corresponding dimension in the base type is fixed size, the dimension has that size in the result.

Suppose that c has a base type with class `int8` and size `10x7x8x?`. In the generated code, the class of $[c\{I\}]$ is `int8`. The size of $[c\{I\}]$ is `0x0x8x0`. The second dimension is 0. The first and

last dimensions are 0 because those dimensions are variable size in the base type. The third dimension is 8 because the size of the third dimension of the base type is a fixed size 8.

Inside concatenation, if curly-brace indexing of a variable-size cell array returns no elements, the generated code can have the following differences from MATLAB:

- The class of `[...c{i}...]` in the generated code can differ from the class in MATLAB.

When `c{I}` returns no elements, MATLAB removes `c{I}` from the concatenation. Therefore, `c{I}` does not affect the class of the result. MATLAB determines the class of the result based on the classes of the remaining arrays, according to a precedence of classes. See “Valid Combinations of Unlike Classes”. In the generated code, the class of `[c{I}]` affects the class of the result of the overall concatenation `[...[c{I}]...]` because the code generator treats `c{I}` as `[c{I}]`. The previously described rules determine the class of `[c{I}]`.

- In the generated code, the size of `[c{I}]` can differ from the size in MATLAB.

In MATLAB, the concatenation `[c{I}]` is a 0x0 double. In the generated code, the previously described rules determine the size of `[c{I}]`.

Variable-Sizing Restrictions for Code Generation of Toolbox Functions

In this section...

“Common Restrictions” on page 29-25

“Toolbox Functions with Restrictions for Variable-Size Data” on page 29-25

Common Restrictions

The following common restrictions apply to multiple toolbox functions, but only for code generation. To determine which of these restrictions apply to specific library functions, see the table in “Toolbox Functions with Restrictions for Variable-Size Data” on page 29-25.

Variable-length vector restriction

Inputs to the library function must be variable-length vectors or fixed-size vectors. A variable-length vector is a variable-size array that has the shape $1 \times n$ or $n \times 1$ (one dimension is variable sized and the other is fixed at size 1). Other shapes are not permitted, even if they are vectors at run time.

Automatic dimension restriction

This restriction applies to functions that take the working dimension (the dimension along which to operate) as input. In MATLAB and in code generation, if you do not supply the working dimension, the function selects it. In MATLAB, the function selects the first dimension whose size does not equal 1. For code generation, the function selects the first dimension that has a variable size or that has a fixed size that does not equal 1. If the working dimension has a variable size and it becomes 1 at run time, then the working dimension is different from the working dimension in MATLAB. Therefore, when run-time error checks are enabled, an error can occur.

For example, suppose that X is a variable-size matrix with dimensions $1 \times 3 \times 5$. In the generated code, `sum(X)` behaves like `sum(X,2)`. In MATLAB, `sum(X)` behaves like `sum(X,2)` unless `size(X,2)` is 1. In MATLAB, when `size(X,2)` is 1, `sum(X)` behaves like `sum(X,3)`.

To avoid this issue, specify the intended working dimension explicitly as a constant value. For example, `sum(X,2)`.

Array-to-vector restriction

The function issues an error when a variable-size array that is not a variable-length vector assumes the shape of a vector at run time. To avoid the issue, specify the input explicitly as a variable-length vector instead of a variable-size array.

Array-to-scalar restriction

The function issues an error if a variable-size array assumes a scalar value at run time. To avoid this issue, specify scalars as fixed size.

Toolbox Functions with Restrictions for Variable-Size Data

The following table lists functions that have code generation restrictions for variable-size data. For additional restrictions for these functions, and restrictions for all functions and objects supported for

code generation, see “Functions and Objects Supported for C/C++ Code Generation” (MATLAB Coder).

| Function | Restrictions for Variable-Size Data |
|----------|---|
| all | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. An error occurs if you pass the first argument a variable-size matrix that is 0-by-0 at run time. |
| any | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. An error occurs if you pass the first argument a variable-size matrix that is 0-by-0 at run time. |
| cat | <ul style="list-style-type: none"> Dimension argument must be a constant. |
| conv | <ul style="list-style-type: none"> See “Variable-length vector restriction” on page 29-25. Input vectors must have the same orientation, either both row vectors or both column vectors. |
| cov | <ul style="list-style-type: none"> For <code>cov(X)</code>, see “Array-to-vector restriction” on page 29-25. |
| cross | <ul style="list-style-type: none"> Variable-size array inputs that become vectors at run time must have the same orientation. |
| deconv | <ul style="list-style-type: none"> For both arguments, see “Variable-length vector restriction” on page 29-25. |
| detrend | <ul style="list-style-type: none"> For first argument for row vectors only, see “Array-to-vector restriction” on page 29-25. |
| diag | <ul style="list-style-type: none"> See “Array-to-vector restriction” on page 29-25. |
| diff | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. Length of the working dimension must be greater than the difference order input when the input is variable sized. For example, if the input is a variable-size matrix that is 3-by-5 at run time, <code>diff(x,2,1)</code> works but <code>diff(x,5,1)</code> generates a run-time error. |
| fft | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. |
| filter | <ul style="list-style-type: none"> For first and second arguments, see “Variable-length vector restriction” on page 29-25. See “Automatic dimension restriction” on page 29-25. |
| hist | <ul style="list-style-type: none"> For second argument, see “Variable-length vector restriction” on page 29-25. For second input argument, see “Array-to-scalar restriction” on page 29-25. |
| histc | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. |
| ifft | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. |
| ind2sub | <ul style="list-style-type: none"> First input (the size vector input) must be fixed size. |
| interp1 | <ul style="list-style-type: none"> For the <code>xq</code> input, see “Array-to-vector restriction” on page 29-25. If <code>v</code> becomes a row vector at run time, the array to vector restriction on page 29-25 applies. If <code>v</code> becomes a column vector at run time, this restriction does not apply. |

| Function | Restrictions for Variable-Size Data |
|----------|---|
| interpft | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. |
| ipermute | <ul style="list-style-type: none"> Order input must be fixed size. |
| issorted | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. |
| magic | <ul style="list-style-type: none"> Argument must be a constant. Output can be fixed-size matrices only. |
| max | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. |
| maxk | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. |
| mean | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. An error occurs if you pass as the first argument a variable-size matrix that is 0-by-0 at run time. |
| median | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. An error occurs if you pass as the first argument a variable-size matrix that is 0-by-0 at run time. |
| min | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. |
| mink | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. |
| mode | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. An error occurs if you pass as the first argument a variable-size matrix that is 0-by-0 at run time. |
| mtimes | <p>Consider the multiplication $A*B$. If the code generator is aware that A is scalar and B is a matrix, the code generator produces code for scalar-matrix multiplication. However, if the code generator is aware that A and B are variable-size matrices, it produces code for a general matrix multiplication. At run time, if A turns out to be scalar, the generated code does not change its behavior. Therefore, when run-time error checks are enabled, a size mismatch error can occur.</p> |
| nchoosek | <ul style="list-style-type: none"> The second input, k, must be a fixed-size scalar. The second input, k, must be a constant for static allocation. If you enable dynamic allocation, the second input can be a variable. You cannot create a variable-size array by passing in a variable, k, unless you enable dynamic allocation. |
| permute | <ul style="list-style-type: none"> Order input must be fixed-size. |
| planerot | <ul style="list-style-type: none"> Input must be a fixed-size, two-element column vector. It cannot be a variable-size array that takes on the size 2-by-1 at run time. |
| poly | <ul style="list-style-type: none"> See “Variable-length vector restriction” on page 29-25. |
| polyfit | <ul style="list-style-type: none"> For first and second arguments, see “Variable-length vector restriction” on page 29-25. |
| prod | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. An error occurs if you pass as the first argument a variable-size matrix that is 0-by-0 at run time. |

| Function | Restrictions for Variable-Size Data |
|----------|--|
| rand | <ul style="list-style-type: none"> For an upper-bounded variable N, <code>rand(1,N)</code> produces a variable-length vector of $1 \times M$ where M is the upper bound on N. For an upper-bounded variable N, <code>rand([1 N])</code> may produce a variable-length vector of $:1 \times M$ where M is the upper bound on N. |
| randi | <ul style="list-style-type: none"> For an upper-bounded variable N, <code>randi(imax,1,N)</code> produces a variable-length vector of $1 \times M$ where M is the upper bound on N. For an upper-bounded variable N, <code>randi(imax,[1 N])</code> may produce a variable-length vector of $:1 \times M$ where M is the upper bound on N. |
| randn | <ul style="list-style-type: none"> For an upper-bounded variable N, <code>randn(1,N)</code> produces a variable-length vector of $1 \times M$ where M is the upper bound on N. For an upper-bounded variable N, <code>randn([1 N])</code> may produce a variable-length vector of $:1 \times M$ where M is the upper bound on N. |
| reshape | <ul style="list-style-type: none"> If the input is a variable-size array and the output array has at least one fixed-length dimension, do not specify the output dimension sizes in a size vector <code>sz</code>. Instead, specify the output dimension sizes as scalar values, <code>sz1, . . . , szN</code>. Specify fixed-size dimensions as constants. When the input is a variable-size empty array, the maximum dimension size of the output array (also empty) cannot be larger than that of the input. |
| roots | <ul style="list-style-type: none"> See “Variable-length vector restriction” on page 29-25. |
| shiftdim | <ul style="list-style-type: none"> If you do not supply the second argument, the number of shifts is determined at compilation time by the upper bounds of the dimension sizes. Therefore, at run time the number of shifts is constant. An error occurs if the dimension that is shifted to the first dimension has length 1 at run time. To avoid the error, supply the number of shifts as the second input argument (must be a constant). First input argument must have the same number of dimensions when you supply a positive number of shifts. |
| sort | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. |
| std | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. An error occurs if you pass a variable-size matrix with 0-by-0 dimensions at run time. |
| sub2ind | <ul style="list-style-type: none"> First input (the size vector input) must be fixed size. |
| sum | <ul style="list-style-type: none"> See “Automatic dimension restriction” on page 29-25. An error occurs if you pass as the first argument a variable-size matrix that is 0-by-0 at run time. |

| Function | Restrictions for Variable-Size Data |
|-----------------|---|
| trapz | <ul style="list-style-type: none">• See “Automatic dimension restriction” on page 29-25.• An error occurs if you pass as the first argument a variable-size matrix that is 0-by-0 at run time. |
| typecast | <ul style="list-style-type: none">• See “Variable-length vector restriction” on page 29-25 on first argument. |
| var | <ul style="list-style-type: none">• See “Automatic dimension restriction” on page 29-25.• An error occurs if you pass a variable-size matrix with 0-by-0 dimensions at run time. |
| vecnorm | <ul style="list-style-type: none">• See “Automatic dimension restriction” on page 29-25. |

Generate Code With Implicit Expansion Enabled

Implicit expansion refers to the automatic size change of compatible operands to apply element-wise operations. Two dimensions have compatible sizes if, for every dimension, the dimension sizes of the arrays are either the same or one of them is singleton. See “Compatible Array Sizes for Basic Operations”.

Implicit expansion in the generated code is enabled by default. Code generated with implicit expansion enabled might differ from code generated with implicit expansion disabled in these ways:

- Output size
- Additional code generation
- Performance variation

For variable-size dynamic arrays, the generated code exhibits these changes to accomplish implicit expansion at run-time.

For fixed-size and constant arrays, because the values and sizes of the operands are known at compile time, the code generated to calculate the implicitly expanded output does not require additional code generation or cause performance variations.

To control implicit expansion in the generated code, see “Optimize Implicit Expansion in Generated Code” (MATLAB Coder).

Output Size

Implicit expansion automatically expands the operands to apply element-wise operations. For example, consider these input types of compatible size:

```
a_type = coder.typeof(1,[2 1]);  
b_type = coder.typeof(1,[2 inf]);
```

A binary operation on these two operands with implicit expansion enabled automatically expands the second dimension of `a_type` to result in an output size of 2-by-Inf. With implicit expansion disabled, the second dimension of `a_type` is not automatically expanded, and the output size is 2-by-1.

For existing workflows created with implicit expansion disabled in the generated code, generating code for the same MATLAB code with implicit expansion enabled might generate size mismatch errors or change the size of outputs from binary operations and functions. To troubleshoot size mismatch errors, see “Diagnose and Fix Variable-Size Data Errors” (MATLAB Coder).

Additional Code Generation

Implicit expansion enables the operands to be automatically expanded if the operand sizes are compatible. To perform this size change, the generated code introduces code that allows the operands to be expanded.

For example, consider the following code snippet. The function `vector_sum` finds the sum of two arrays.

```
function out = vector_sum(a,b)  
out = a + b;  
end
```

Consider the variable-size dynamic array defined here:

```
c_type = coder.typeof(1,[1 Inf]);
```

Generate code for `vector_sum` by using this command:

```
codegen vector_sum -args {c_type, c_type} -config:lib -report
```

The generated code for this function with implicit expansion:

```
static void plus(emxArray_real_T *out,
... const emxArray_real_T *b, const emxArray_real_T *a)
{
    int i;
    ....
    if (a->size[1] == 1) {
        out->size[1] = b->size[1];
    } else {
        out->size[1] = a->size[1];
    }
    ....
    if (a->size[1] == 1) {
        loop_ub = b->size[1];
    }
    else {
        loop_ub = a->size[1];
    }
    for (i = 0; i < loop_ub; i++) {
        out->data[i] = b->data[i * stride_0_1] + a->data[i * stride_1_1];
    }
}

void vector_sum(const emxArray_real_T *a, const emxArray_real_T *b, emxArray_real_T *out)
{
    int i;
    int loop_ub;
    if (b->size[1] == a->size[1]) {
        i = out->size[0] * out->size[1];
        out->size[0] = 1;
        out->size[1] = b->size[1];
        emxEnsureCapacity_real_T(out, i);
        loop_ub = b->size[1];
        for (i = 0; i < loop_ub; i++) {
            out->data[i] = b->data[i] + a->data[i];
        }
    } else {
        plus(out, b, a);
    }
}
```

The generated code for this function without implicit expansion:

```
void vector_sum(const emxArray_real_T *a,
... const emxArray_real_T *b, emxArray_real_T *out){
    int i;
    int loop_ub;
    i = out->size[0] * out->size[1];
    out->size[0] = 1;
    out->size[1] = b->size[1];
```

```
    emxEnsureCapacity_real_T(out, i);  
    loop_ub = b->size[1];  
    for (i = 0; i < loop_ub; i++) {  
        out->data[i] = b->data[i] + a->data[i];  
    }  
}
```

With implicit expansion enabled, the code generator creates a supporting function, in this case `plus`, to carry out the size change and to calculate the output.

In most cases, the supporting function carrying out implicit expansion is named after the binary operation it is assisting. In the previous example, if the expression `out = a + b` is changed to `out = a - b`, the name of the supporting function changes to `minus`.

Some supporting functions might also be named as `expand_op`, where `op` refers to the binary operation. In the previous example, if the expression `out = a + b` is replaced with `out = max(a, b)`, the name of the supporting function in the generated code changes to `expand_max`.

If multiple operations in an expression require implicit expansion, the generated code includes a supporting function that is named `binary_expand_op`. The supporting functions change the size of the operand and apply the binary operations.

If you want to apply specific binary operations and functions without implicit expansion, use `coder.sameSizeBinaryOp`. The code generated to apply this function does not include additional code to expand the operands. The output of this function does not expand the operands in MATLAB. This function does not support scalar expansion. Operands must be of the same size.

If you want to disable implicit expansion inside a function for all binary operations within that function in the generated code, call `coder.noImplicitExpansionInFunction` in the required function. Implicit expansion in MATLAB code is still enabled.

Performance Variation

Code generated with implicit expansion enabled might perform differently than when implicit expansion is disabled. Depending on the input to the generated code that uses implicit expansion, the code might take longer to evaluate the output.

If the generated code does not match the performance requirements of your workflow due to implicit expansion, generate code for your project by turning off implicit expansion for specific binary operations, specific function bodies, or for your whole project. See “Optimize Implicit Expansion in Generated Code” (MATLAB Coder).

Note Before disabling implicit expansion, ensure that the external code does not use implicit expansion. Disabling implicit expansion for an entire project might cause errors when generating code if your project includes MATLAB code from external sources.

See Also

`coder.noImplicitExpansionInFunction` | `coder.sameSizeBinaryOp`

Related Examples

- “Compatible Array Sizes for Basic Operations”
- “Diagnose and Fix Variable-Size Data Errors” (MATLAB Coder)
- “Optimize Implicit Expansion in Generated Code” (MATLAB Coder)

Optimize Implicit Expansion in Generated Code

Implicit expansion in the generated code is enabled by default. The code generator introduces modifications in the generated code to perform implicit expansion. The changes in the generated code might result in additional code to expand the operands. The expansion of the operands might affect the performance of the generated code. See “Generate Code With Implicit Expansion Enabled” (MATLAB Coder).

Implicit expansion might change the size of the outputs from the supported operators and functions causing size and type mismatch errors in your workflow.

For fine-grained control of where implicit expansion is enabled in the generated code, use the following functions in your MATLAB code:

- `coder.noImplicitExpansionInFunction`
- `coder.sameSizeBinaryOp`

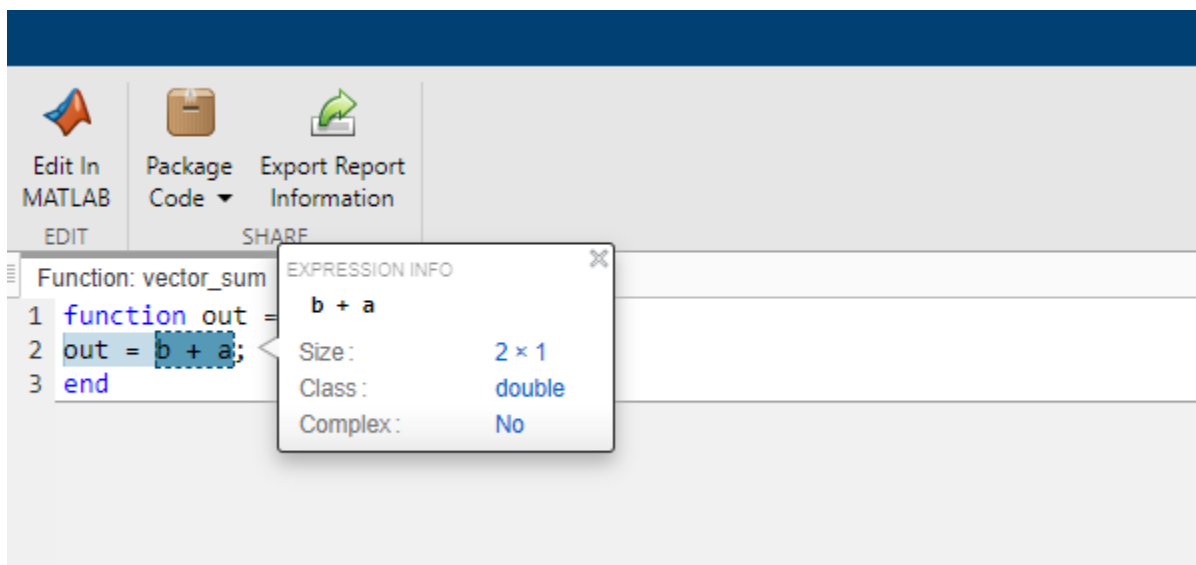
For example, consider this code snippet. The function `vector_sum` finds the sum of two arrays of compatible sizes.

```
function out = vector_sum(a,b)
out = b + a;
end
```

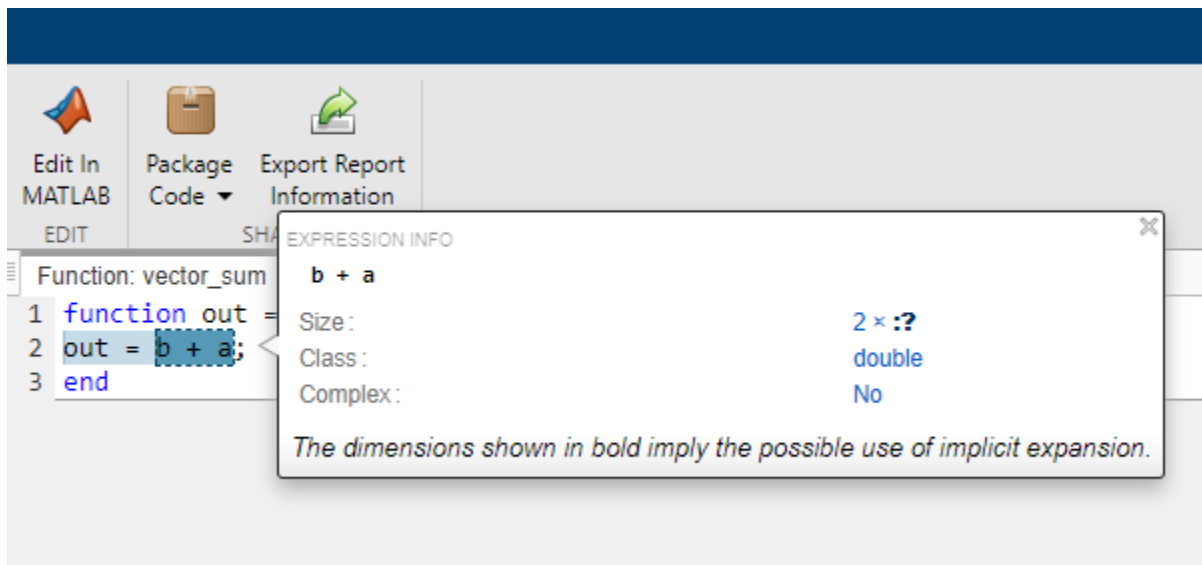
The types of operands `a` and `b` are defined as:

```
a_type = coder.typeof(1,[2 1])      %size: 2x1
b_type = coder.typeof(1,[2 Inf])    %size: 2x:inf
```

Without implicit expansion, the size of the `out` variable is calculated as `2x1`.



With implicit expansion, the size of the variable `out` is calculated as `2x:?`.



These code snippets outline the changes in the generated code for the function `vector_sum`, while implicit expansion is disabled and enabled. To generate the code, the types of operands `a` and `b` are defined as:

```
a_type = coder.typeof(1,[1 Inf])    %size: 1x:inf
b_type = coder.typeof(1,[1 Inf])    %size: 1x:inf
```

| Generated Code With Implicit Expansion Disabled | Generated Code With Implicit Expansion Enabled |
|--|--|
| <pre>void vector_sum(const emxArray_real_T *a, const emxArray_real_T *out) { int i; int loop_ub; i = out->size[0] * out->size[1]; out->size[0] = 1; out->size[1] = b->size[1]; emxEnsureCapacity_real_T(out, i); loop_ub = b->size[1]; for (i = 0; i < loop_ub; i++) { out->data[i] = b->data[i] + a->data[i]; } }</pre> | <pre>static void plus(const emxArray_real_T *out, const emxArray_real_T *a) { int i; if (a->size[1] == 1) { out->size[1] = b->size[1]; } else { out->size[1] = a->size[1]; } if (a->size[1] == 1) { loop_ub = b->size[1]; } else { loop_ub = a->size[1]; } for (i = 0; i < loop_ub; i++) { out->data[i] = b->data[i * stride_0_1] + a->data[i]; } } void vector_sum(const emxArray_real_T *a, const emxArray_real_T *out) { int i; int loop_ub; if (b->size[1] == a->size[1]) { i = out->size[0] * out->size[1]; out->size[0] = 1; out->size[1] = b->size[1]; emxEnsureCapacity_real_T(out, i); loop_ub = b->size[1]; for (i = 0; i < loop_ub; i++) { out->data[i] = b->data[i] + a->data[i]; } } else { plus(out, b, a); } }</pre> |

Disable Implicit Expansion in Specified Function by Using `coder.noImplicitExpansionInFunction`

If you require implicit expansion in your project but not in specific functions, disable implicit expansion for the generated code of that function by calling `coder.noImplicitExpansionInFunction` within the function.

For example, the code generated for `vector_sum` does not apply implicit expansion.

| MATLAB Code | Generated Code with <code>coder.sameSizeBinaryOp</code> |
|---|--|
| <pre>function out = vector_sum(a,b) coder.noImplicitExpansionInFunction(); out = a + b; end a = coder.typeof(1,[1 Inf]) %size: 1x:inf b = coder.typeof(1,[1 Inf]) %size: 1x:inf codegen vector_sum -launchreport ... -args {a,b} -config:lib</pre> | <pre>void vector_sum(const emxArray_real_T *a, const emxArray_real_T *b, emxArray_real_T *out) { int i; int loop_ub; i = out->size[0] * out->size[1]; out->size[0] = 1; out->size[1] = a->size[1]; emxEnsureCapacity_real_T(out, i); loop_ub = a->size[1]; for (i = 0; i < loop_ub; i++) { out->data[i] = a->data[i] + b->data[i]; } }</pre> |

Note `coder.noImplicitExpansionInFunction` does not disable implicit expansion in your MATLAB code. It disables implicit expansion only in the generated code.

Disable Implicit Expansion for Specific Binary Operation by Using `coder.sameSizeBinaryOp`

Use the function `coder.sameSizeBinaryOp` to perform an error check to ensure that the operands are the same size and prevent the code generator from generating implicitly expanded code for that function.

For example, this code snippet applies the plus operation by using `coder.sameSizeBinaryOp` without implicit expansion.

| MATLAB Code | Generated Code |
|--|--|
| <pre>function out = vector_sum(a,b) out = coder.sameSizeBinaryOp(@plus, a, b); end a = coder.typeof(1,[1 Inf]) %size: 1x:inf b = coder.typeof(1,[1 Inf]) %size: 1x:inf codegen vector_sum -launchreport ... -args {a,b} -config:lib</pre> | <pre>void vector_sum(const emxArray_real_T *a, const emxArray_real_T *b, emxArray_real_T *out) { int i; int loop_ub; i = out->size[0] * out->size[1]; out->size[0] = 1; out->size[1] = a->size[1]; emxEnsureCapacity_real_T(out, i); loop_ub = a->size[1]; for (i = 0; i < loop_ub; i++) { out->data[i] = a->data[i] + b->data[i]; } }</pre> |

`coder.sameSizeBinaryOp` does not support scalar expansion. Operands given to `coder.sameSizeBinaryOp` must be of the same size.

Disable Implicit Expansion in your Project

If you do not require implicit expansion in your generated code or do not want the modifications to affect your generated code, turn it off by setting the `EnableImplicitExpansion` flag in your `coder.config` object to `false`. This flag is set to `true` by default.

```
cfg = coder.config;  
cfg.EnableImplicitExpansion = false;
```

Disable implicit expansion in your Simulink model by setting the model-wide parameter **Enable Implicit Expansion in MATLAB functions** to `false`. Alternatively, use this command:

```
set_param(gcs, 'EnableImplicitExpansion', false);
```

Note Before turning off implicit expansion, ensure that the external code does not use implicit expansion. Disabling implicit expansion for an entire project might cause errors when generating code if your project includes MATLAB code from external sources.

See Also

`coder.noImplicitExpansionInFunction` | `coder.sameSizeBinaryOp`

Related Examples

- “Generate Code With Implicit Expansion Enabled” (MATLAB Coder)
- “Compatible Array Sizes for Basic Operations”
- “Diagnose and Fix Variable-Size Data Errors” (MATLAB Coder)

Representation of Arrays in Generated Code

The code generator produces C/C++ array definitions that depend on the array element type and whether the array uses static or dynamic memory allocation. Use the generated array implementations to interface your arrays with the generated code.

Memory allocation for arrays require different implementations:

- For a fixed-size array or a variable-size array whose size is bounded within a predefined memory threshold, the generated C/C++ definition consists of a fixed-size array of elements and a size vector that stores the total number of array elements. In some cases, the fixed-size element array and the size vector are stored within a structure. The memory for this array comes from the program stack and is statically allocated.
- For an array whose size is unbounded at compile time, or whose bounds exceed the predefined threshold, the generated C definition consists of a data structure called an `emxArray`. The generated C++ definition consists of a `coder::array` class template.

The predefined threshold size (in bytes) is specified in your configuration objects. The default value of the parameter is 65536. See `DynamicMemoryAllocationThreshold` in `coder.MexCodeConfig`, `coder.CodeConfig`, or `coder.EmbeddedCodeConfig`.

For dynamically allocated arrays, the run-time allocated size is set based on the current array size. During program execution, as run-time allocated size is exceeded, the generated code reallocates additional memory space from the heap and adds it to the dynamic array storage.

This table lists a few typical cases for array representation in the generated code.

| Algorithm Description and Array Size | MATLAB Function | Generated C/C++ Code |
|---|--|--|
| <p>Create a fixed-size 1-by-500 row vector. The array is the output of the MATLAB function</p> <p>The generated code allocates memory to a fixed-size vector on the program stack.</p> | <pre>function B = create_vec0 B = zeros(1,500); end</pre> | <pre>create_vec0(double B[500]) { memset(&B[0], 0, 500U * sizeof(double)); }</pre> <p>The array is the input to the function in the generated code.</p> |
| <p>Create a fixed-size 1-by-20 row vector. Declare the array as variable-size with bounds at 500 elements. Assign this variable-size array to the input array.</p> <p>This array is bound within the size threshold and is the input to the function in the generated code.</p> | <pre>function create_vec1(B) A = zeros(1,20); coder.varsize('A',[1 500],[0 1]); B = A; end</pre> | <pre>create_vec1(double B_data[], int B_size[]) { int i; B_size[0] = 1; B_size[1] = 20; for (i = 0; i < 20; i++) { B_data[i] = 1.0; } }</pre> <p>Note The generated code includes the inputs in the function parameters.</p> |

| Algorithm Description and Array Size | MATLAB Function | Generated C/C++ Code |
|---|---|---|
| <p>Create a local fixed-size 1-by-20000 row vector. Declare the array as variable-size with bounds at 30,000 elements.</p> <p>The variable-size array exceeds the predefined dynamic memory allocation threshold. This array is stored on heap memory.</p> <p>The generated code includes the output array in the function parameter.</p> | <pre>function B = create_vec2() %C A = ones(1,20000); coder.varsize("A",[1 30000], B = [1 A]; end</pre> | <pre>C: void create_vec2(emxArray_real_T *B) { double *B_data; int i; i = B->size[0] * B->size[1]; B->size[0] = 1; B->size[1] = 20001; emxEnsureCapacity_real_T(B, i); B_data = B->data; B_data[0] = 1.0; for (i = 0; i < 20000; i++) { B_data[i + 1] = 1.0; } } C++: void create_vec2(coder::array<double, 2 { B.set_size(1, 20001); B[0] = 1.0; for (int i{0}; i < 20000; i++) { B[i + 1] = 1.0; } }</pre> |

| Algorithm Description and Array Size | MATLAB Function | Generated C/C++ Code |
|--|--|---|
| <p>Create an array that has the size determined by an unbounded integer input.</p> <p>The generated array size is unknown and unbounded at compile time.</p> | <pre>function y = create_vec3(n) y = ones(1,n,'int8');</pre> | <pre>%codegen void create_vec3(double n, mxArray_int { int i; int loop_ub_tmp; signed char *y_data; i = y->size[0] * y->size[1]; y->size[0] = 1; loop_ub_tmp = (int)n; y->size[1] = (int)n; emxEnsureCapacity_int8_T(y, i); y_data = y->data; for (i = 0; i < loop_ub_tmp; i++) { y_data[i] = 1; } } C++: void create_vec3(double n, coder::array { int loop_ub_tmp; loop_ub_tmp = static_cast<int>(n); y.set_size(1, loop_ub_tmp); for (int i{0}; i < loop_ub_tmp; i++) y[i] = 1; } }</pre> |

To learn about the `emxArray` data structure, see “Use C Arrays in the Generated Function Interfaces” (MATLAB Coder).

To learn about the `coder::array` class template, see “Use Dynamically Allocated C++ Arrays in Generated Function Interfaces” (MATLAB Coder).

Customize Interface Generation

By default, the generated C++ code uses the `coder::array` template to implement dynamically allocated arrays. You can choose to generate C++ code that uses the C style `emxArray` data structure to implement dynamically allocated arrays. To generate C style `emxArray` data structures, do either of the following:

- In a code configuration object (`coder.MexCodeConfig`, `coder.CodeConfig`, or `coder.EmbeddedCodeConfig`), set the `DynamicMemoryAllocationInterface` parameter to 'C'.
- Alternatively, In the MATLAB Coder app, on the **Memory** tab, set **Dynamic memory allocation interface** to Use C style `EmxArray`.

To create dynamically allocated arrays for variable-size arrays in the generated code, do either of the following:

- Set the `EnabledDynamicMemoryAllocation` flag to `true`.
- Alternatively, in the MATLAB Coder App, on the **Memory** tab, select the **Enable dynamic memory allocation** option.

By default, arrays that are bounded within a threshold size do not use dynamic allocation in the generated code. Alternatively, you can disable dynamic memory allocation and change the dynamic memory allocation threshold. See “Control Memory Allocation for Variable-Size Arrays” (MATLAB Coder).

See Also

`coder.config` | `coder.MexCodeConfig` | `coder.CodeConfig` | `coder.EmbeddedCodeConfig`

Related Examples

- “Use C Arrays in the Generated Function Interfaces” (MATLAB Coder)
- “Use Dynamically Allocated C++ Arrays in Generated Function Interfaces” (MATLAB Coder)

Control Memory Allocation for Fixed-Size Arrays

Dynamic memory allocation for fixed-size arrays allocates memory for the array on the heap instead of allocating memory on the program stack. Consider using dynamic memory allocation when:

- The fixed-size arrays are large and you do not want to allocate memory on the stack.
- Your target hardware memory is limited and you do not want to allocate memory for the arrays on the program stack.

For larger arrays, you can significantly reduce storage requirements. Dynamic memory allocation might result in slower execution of the generated code.

Enable Dynamic Memory Allocation for All Fixed-Size Arrays

By default, dynamic memory allocation for fixed-size arrays is disabled. To enable it:

- In a configuration object for code generation, set the `DynamicMemoryAllocationForFixedSizeArrays` parameter to 'Always'.
- Alternatively, in the app, under **Memory** settings, set **Dynamic memory allocation fixed-sized arrays** to 'Always'.

The code generator dynamically allocates memory on the heap for all fixed-size arrays whose size is greater than 64 bytes.

Enable Dynamic Memory Allocation for Arrays Bigger Than a Threshold

Instead of allocating all fixed-size arrays dynamically on the heap, you can specify the threshold size above which memory is dynamically allocated. To instruct the code generator to use dynamic memory allocation for fixed-size arrays whose size is greater than or equal to the threshold:

- In the configuration object, set the `DynamicMemoryAllocationForFixedSizeArrays` to 'Threshold'.
- In the MATLAB Coder app, in the **Memory settings**, set **Dynamic memory allocation for fixed-sized arrays** to For arrays with max size at or above threshold.

To instruct the code generator to use dynamic memory allocation for fixed-size arrays whose size is greater than or equal to the threshold, in the configuration object, set the `DynamicMemoryAllocationForFixedSizeArrays` to 'Threshold'.

The default dynamic memory allocation threshold is 64 kilobytes. To change the threshold, in a configuration object for fixed-point acceleration, set the `DynamicMemoryAllocationThreshold`.

See Also

`coder.EmbeddedCodeConfig` | `coder.MexCodeConfig` | `coder.CodeConfig`

Related Examples

- “Control Memory Allocation for Variable-Size Arrays” (MATLAB Coder)
- “Representation of Arrays in Generated Code” (MATLAB Coder)
- “Use C Arrays in the Generated Function Interfaces” (MATLAB Coder)

- “Use Dynamically Allocated C++ Arrays in Generated Function Interfaces” (MATLAB Coder)

Resolve Error: Size Mismatches

Issue

The code generator produces size mismatch errors when array sizes are incompatible or implicit expansion is unavailable.

Most binary operators and functions in MATLAB and generated code support numeric arrays that have compatible sizes. Two inputs have compatible sizes if, for every dimension, the sizes of the inputs are either the same or one of them is 1. In the simplest cases, two array sizes are compatible if they are exactly the same or if one is a scalar. For example:

```
magic(4) + ones(4,1);
% where magic(4) =           ones(4,1) =
%  16     2     3     13           1
%   5    11    10     8           1
%   9     7     6    12           1
%   4    14    15     1           1

ans =

    17     3     4    14
     6    12    11     9
    10     8     7    13
     5    15    16     2
```

The second array implicitly expands to match the dimensions of the first matrix. For more information, see “Compatible Array Sizes for Basic Operations”.

Implicit expansion might be unavailable while performing binary operations on arrays of compatible size if any the following conditions are true :

- Your function scope includes the `coder.noImplicitExpansionInFunction` function.
- You use the `coder.sameSizeBinaryOp` function to carry out the binary operation.
- You turn off implicit expansion for your project.

Size mismatches or unavailability of implicit expansion generates the following error:

```
%Size mismatch between two arrays
Size mismatch (size [10][1] ~= size [1][10])
```

When the above conditions are true for structure fields and cell elements, the code generator produces the following errors respectively:

```
%Size mismatch in structure fields
Size mismatch (size [10][1] ~= size [1][10]) in field StructField

%Size mismatch in cell elements
Size mismatch (size [10][1] ~= size [1][10]) in element cellElement.
```

Possible Solutions

Verify that, in binary operations where you enable implicit expansion, the operations are in the scope of functions. Check for these conditions:

- Array size compatibility.
- Binary operations in the scope of functions that call `coder.noImplicitExpansionInFunction`.
- `coder.sameSizeBinaryOp` does not implicitly expand its operands or support scalar expansion.
- If you have turned off implicit expansion for the whole project, all operations that require implicit expansion generate an error.

Perform Binary Operations on Arrays of Compatible Sizes

If you must carry out a binary operation on arrays of differing sizes, make sure that sizes are compatible and implicit expansion is enabled in the function scope. See “Compatible Array Sizes for Basic Operations”.

Call Binary Operation Without `coder.noImplicitExpansionInFunction`

If you must include `coder.noImplicitExpansionInFunction` in your function, call the required binary operation in another function where implicit expansion is enabled.

Call Binary Operation Without `coder.sameSizeBinaryOp`

If you do not want implicit expansion for a specific operation, provide input arguments that are of same size to `coder.sameSizeBinaryOp`. `coder.sameSizeBinaryOp` does not allow scalar expansion and generates an error if the input arguments are not of the same size.

Enable Implicit Expansion for Project

If enabling implicit expansion does not affect your project, consider enabling it by setting the `EnableImplicitExpansion` property in your code configuration object to `true`.

If you need implicit expansion for specific operations, consider using `coder.sameSizeBinaryOp` or `coder.noImplicitExpansionInFunction` to prevent the other operations from implicitly expanding. See “Optimize Implicit Expansion in Generated Code” (MATLAB Coder).

See Also

`coder.noImplicitExpansionInFunction` | `coder.sameSizeBinaryOp`

Related Examples

- “Compatible Array Sizes for Basic Operations”
- “Generate Code With Implicit Expansion Enabled” (MATLAB Coder)
- “Optimize Implicit Expansion in Generated Code” (MATLAB Coder)

Code Generation for Cell Arrays

- “Code Generation for Cell Arrays” on page 30-2
- “Control Whether a Cell Array Is Variable-Size” on page 30-4
- “Define Cell Array Inputs” on page 30-6
- “Cell Array Limitations for Code Generation” on page 30-7

Code Generation for Cell Arrays

When you generate code from MATLAB code that contains cell arrays, the code generator classifies the cell arrays as homogeneous or heterogeneous. This classification determines how a cell array is represented in the generated code. It also determines how you can use the cell array in MATLAB code from which you generate code.

When you use cell arrays in MATLAB code that is intended for code generation, you must adhere to certain restrictions. See “Cell Array Limitations for Code Generation” on page 30-7.

Homogeneous vs. Heterogeneous Cell Arrays

A homogeneous cell array has these characteristics:

- The cell array is represented as an array in the generated code.
- All elements have the same properties. The type associated with the cell array specifies the properties of all elements rather than the properties of individual elements.
- The cell array can be variable-size.
- You can index into the cell array with an index whose value is determined at run time.

A heterogeneous cell array has these characteristics:

- The cell array is represented as a structure in the generated code. Each element is represented as a field of the structure.
- The elements can have different properties. The type associated with the cell array specifies the properties of each element individually.
- The cell array cannot be variable-size.
- You must index into the cell array with a constant index or with `for`-loops that have constant bounds.

The code generator uses heuristics to determine the classification of a cell array as homogeneous or heterogeneous. It considers the properties (class, size, complexity) of the elements and other factors, such as how you use the cell array in your program. Depending on how you use a cell array, the code generator can classify a cell array as homogeneous in one case and heterogeneous in another case. For example, consider the cell array `{1 [2 3]}`. The code generator can classify this cell array as a heterogeneous 1-by-2 cell array. The first element is double scalar. The second element is a 1-by-2 array of doubles. However, if you index into this cell array with an index whose value is determined at run time, the code generator classifies it as a homogeneous cell array. The elements are variable-size arrays of doubles with an upper bound of 2.

Controlling Whether a Cell Array Is Homogeneous or Heterogeneous

For cell arrays with certain characteristics, you cannot control the classification as homogeneous or heterogeneous:

- If the elements have different classes, the cell array must be heterogeneous.
- If the cell array is variable-size, it must be homogeneous.
- If you index into the cell array with an index whose value is determined at run time, the cell array must be homogeneous.

For other cell arrays, you can control the classification as homogeneous or heterogeneous.

To control the classification of cell arrays that are entry-point function inputs, use the `coder.CellType` methods `makeHomogeneous` and `makeHeterogeneous`.

To control the classification of cell arrays that are not entry-point function inputs:

- If the cell array elements have the same class, you can force a cell array to be homogeneous by using `coder.varsize`. See “Control Whether a Cell Array Is Variable-Size” on page 30-4.

Cell Arrays in Reports

To see whether a cell array is homogeneous or heterogeneous, view the variable in the code generation report.

For a homogeneous cell array, the report has one entry that specifies the properties of all elements. The notation `{:}` indicates that all elements of the cell array have the same properties.

| SUMMARY | ALL MESSAGES (0) | | CODE INSIGHTS (0) | |
|---------|------------------|-------|-------------------|--|
| Name | Type | Size | Class | |
| z | Output | 1 × 1 | double | |
| n | Input | 1 × 1 | embedded.fi | |
| ▲ c | Local | 1 × 3 | cell | |
| {:} | | 1 × 1 | double | |

For a heterogeneous cell array, the report has an entry for each element. For example, for a heterogeneous cell array `c` with two elements, the entry for `c{1}` shows the properties for the first element. The entry for `c{2}` shows the properties for the second element.

| SUMMARY | ALL MESSAGES (0) | | CODE INSIGHTS (0) | |
|---------|------------------|-------|-------------------|--|
| Name | Type | Size | Class | |
| n | I/O | 1 × 1 | embedded.fi | |
| z | Output | 1 × 1 | double | |
| ▲ c | Local | 1 × 2 | cell | |
| {1} | | 1 × 1 | double | |
| {2} | | 1 × 1 | char | |

See Also

`coder.CellType` | `coder.varsize`

More About

- “Control Whether a Cell Array Is Variable-Size” on page 30-4
- “Cell Array Limitations for Code Generation” on page 30-7
- “Code Generation Reports” on page 12-27

Control Whether a Cell Array Is Variable-Size

The code generator classifies a variable-size cell array as homogeneous. The cell array elements must have the same class. In the generated code, the cell array is represented as an array.

If a cell array is an entry-point function input, to make it variable-size, use `coder.typeof` or `coder.newtype` to create a type for a variable-size cell array. For example, to create a type for a cell array whose first dimension is fixed and whose second dimension has an upper bound of 10, use this code:

```
t = coder.typeof({1 2 3}, [1 10], [0 1])
```

See “Specify Variable-Size Cell Array Inputs” on page 12-44.

If a cell array is not an entry-point function input, to make it variable-size:

- Create the cell array by using the `cell` function. For example:

```
function z = mycell(n, j)
%#codegen
x = cell(1,n);
for i = 1:n
    x{i} = i;
end
z = x{j};
end
```

For code generation, when you create a variable-size cell array by using `cell`, you must adhere to certain restrictions. See “Definition of Variable-Size Cell Array by Using `cell`” on page 30-8.

- Grow the cell array. For example:

```
function z = mycell(n)
%#codegen
c = {1 2 3};
for i = 1:n
    c{end + 1} = 1;
end
z = c{n};
end
```

- Force the cell array to be variable-size by using `coder.varsize`. Consider this code:

```
function y = mycellfun()
%#codegen
c = {1 2 3};
coder.varsize('c', [1 10]);
y = c;
end
```

Without `coder.varsize`, `c` is fixed-size with dimensions 1-by-3. With `coder.varsize`, `c` is variable-size with an upper bound of 10.

Sometimes, using `coder.varsize` changes the classification of a cell array from heterogeneous to homogeneous. Consider this code:

```
function y = mycell()
%#codegen
```



```
c = {1 [2 3]};  
y = c{2};  
end
```

The code generator classifies `c` as heterogeneous because the elements have different sizes. `c` is fixed-size with dimensions 1-by-2. If you use `coder.varsize` with `c`, it becomes homogeneous. For example:

```
function y = mycell()  
%#codegen  
c = {1 [2 3]};  
coder.varsize('c', [1 10], [0 1]);  
y = c{2};  
end
```

`c` becomes a variable-size homogeneous cell array with dimensions 1-by-:10.

To force `c` to be homogeneous, but not variable-size, specify that none of the dimensions vary. For example:

```
function y = mycell()  
%#codegen  
c = {1 [2 3]};  
coder.varsize('c', [1 2], [0 0]);  
y = c{2};  
end
```

See Also

`coder.CellType` | `coder.varsize`

More About

- “Code Generation for Cell Arrays” on page 30-2
- “Cell Array Limitations for Code Generation” on page 30-7
- “Code Generation for Variable-Size Arrays” on page 29-2

Define Cell Array Inputs

To define types for cell arrays that are inputs to entry-point functions, use one of these approaches:

| Define Types | See |
|-------------------------------------|---|
| At the command line | “Specify Cell Array Inputs at the Command Line” on page 12-42 |
| Programmatically in the MATLAB file | “Define Input Properties Programmatically in MATLAB File” on page 12-35 |

See Also

`coder.CellType`

More About

- “Code Generation for Cell Arrays” on page 30-2

Cell Array Limitations for Code Generation

When you use cell arrays in MATLAB code that is intended for code generation, you must adhere to these restrictions:

- “Cell Array Element Assignment” on page 30-7
- “Variable-Size Cell Arrays” on page 30-8
- “Definition of Variable-Size Cell Array by Using `cell`” on page 30-8
- “Cell Array Indexing” on page 30-11
- “Growing a Cell Array by Using `{end + 1}`” on page 30-12
- “Cell Array Contents” on page 30-13
- “Passing Cell Arrays to External C/C++ Functions” on page 30-13

Cell Array Element Assignment

You must assign a cell array element on all execution paths before you use it. For example:

```
function z = foo(n)
%#codegen
c = cell(1,3);
if n < 1
    c{2} = 1;

else
    c{2} = n;
end
z = c{2};
end
```

The code generator considers passing a cell array to a function or returning it from a function as a use of all elements of the cell array. Therefore, before you pass a cell array to a function or return it from a function, you must assign all of its elements. For example, the following code is not allowed because it does not assign a value to `c{2}` and `c` is a function output.

```
function c = foo()
%#codegen
c = cell(1,3);
c{1} = 1;
c{3} = 3;
end
```

The assignment of values to elements must be consistent on all execution paths. The following code is not allowed because `y{2}` is double on one execution path and char on the other execution path.

```
function y = foo(n)
y = cell(1,3)
if n > 1;
    y{1} = 1;
    y{2} = 2;
    y{3} = 3;
else
    y{1} = 10;
    y{2} = 'a';
```

```

    y{3} = 30;
end

```

Variable-Size Cell Arrays

- `coder.versize` is not supported for heterogeneous cell arrays.
- If you use the `cell` function to define a fixed-size cell array, you cannot use `coder.versize` to specify that the cell array has a variable size. For example, this code causes a code generation error because `x = cell(1,3)` makes `x` a fixed-size, 1-by-3 cell array.

```

...
x = cell(1,3);
coder.versize('x',[1 5])
...

```

You can use `coder.versize` with a cell array that you define by using curly braces. For example:

```

...
x = {1 2 3};
coder.versize('x',[1 5])
...

```

- To create a variable-size cell array by using the `cell` function, use this code pattern:

```

function mycell(n)
    %#codegen
    x = cell(1,n);
    for i = 1:n
        x{i} = i;
    end
end

```

See “Definition of Variable-Size Cell Array by Using `cell`” on page 30-8.

To specify upper bounds for the cell array, use `coder.versize`.

```

function mycell(n)
    %#codegen
    x = cell(1,n);
    for i = 1:n
        x{i} = i;
    end
    coder.versize('x',[1,20]);
end

```

Definition of Variable-Size Cell Array by Using `cell`

For code generation, before you use a cell array element, you must assign a value to it. When you use `cell` to create a variable-size cell array, for example, `cell(1,n)`, MATLAB assigns an empty matrix to each element. However, for code generation, the elements are unassigned. For code generation, after you use `cell` to create a variable-size cell array, you must assign all elements of the cell array before any use of the cell array. For example:

```

function z = mycell(n, j)
    %#codegen

```

```
x = cell(1,n);
for i = 1:n
    x{i} = i;
end
z = x{j};
end
```

The code generator analyzes your code to determine whether all elements are assigned before the first use of the cell array. If the code generator detects that some elements are not assigned, code generation fails with an error message. For example, modify the upper bound of the `for`-loop to `j`.

```
function z = mycell(n, j)
%#codegen
x = cell(1,n);
for i = 1:j %<- Modified here
    x{i} = i;
end
z = x{j};
end
```

With this modification and with inputs `j` less than `n`, the function does not assign values to all of the cell array elements. Code generation produces the error:

```
Unable to determine that every element of 'x{:}' is assigned
before this line.
```

Sometimes, even though your code assigns all elements of the cell array, the code generator reports this message because the analysis does not detect that all elements are assigned. See “Unable to Determine That Every Element of Cell Array Is Assigned” on page 49-39.

To avoid this error, follow these guidelines:

- When you use `cell` to define a variable-size cell array, write code that follows this pattern:

```
function z = mycell(n, j)
%#codegen
x = cell(1,n);
for i = 1:n
    x{i} = i;
end
z = x{j};
end
```

Here is the pattern for a multidimensional cell array:

```
function z = mycell(m,n,p)
%#codegen
x = cell(m,n,p);
for i = 1:m
    for j =1:n
        for k = 1:p
            x{i,j,k} = i+j+k;
        end
    end
end
z = x{m,n,p};
end
```

- Increment or decrement the loop counter by 1.
- Define the cell array within one loop or one set of nested loops. For example, this code is not allowed:

```
function z = mycell(n, j)
x = cell(1,n);
for i = 1:5
    x{i} = 5;
end
for i = 6:n
    x{i} = 5;
end
z = x{j};
end
```

- Use the same variables for the cell dimensions and loop initial and end values. For example, code generation fails for the following code because the cell creation uses *n* and the loop end value uses *m*:

```
function z = mycell(n, j)
x = cell(1,n);
m = n;
for i = 1:m
    x{i} = 2;
end
z = x{j};
end
```

Rewrite the code to use *n* for the cell creation and the loop end value:

```
function z = mycell(n, j)
x = cell(1,n);
for i = 1:n
    x{i} = 2;
end
z = x{j};
end
```

- Create the cell array with this pattern:

```
x = cell(1,n)
```

Assign the cell array to a field of a structure or a property of an object by initializing a temporary variable with the required cell. For example:

```
t = cell(1,n)
for i = 1:n
    t{i} = i+1;
end
myObj.prop = t;
```

Do not assign a cell array to a field of a structure or a property of an object directly. For example, this code is not allowed:

```
myObj.prop = cell(1,n);
for i = 1:n
    myObj.prop{i} = i+1;
end
```

Do not use the `cell` function inside the cell array constructor `{}`. For example, this code is not allowed:

```
x = {cell(1,n)};
```

- The cell array creation and the loop that assigns values to the cell array elements must be together in a unique execution path. For example, the following code is not allowed.

```
function z = mycell(n)
if n > 3
    c = cell(1,n);
else
    c = cell(n,1);
end
for i = 1:n
    c{i} = i;
end
z = c{n};
end
```

To fix this code, move the assignment loop inside the code block that creates the cell array.

```
function z = cellerr(n)
if n > 3
    c = cell( 1,n);
    for i = 1:n
        c{i} = i;
    end
else
    c = cell(n,1);
    for i = 1:n
        c{i} = i;
    end
end
z = c{n};
end
```

Cell Array Indexing

- You cannot index cell arrays by using smooth parentheses `()`. Consider indexing cell arrays by using curly braces `{}` to access the contents of the cell.
- You must index into heterogeneous cell arrays by using constant indices or by using `for`-loops with constant bounds.

For example, the following code is not allowed.

```
x = {1, 'mytext'};
disp(x{randi});
```

You can index into a heterogeneous cell array in a `for`-loop with constant bounds because the code generator unrolls the loop. Unrolling creates a separate copy of the loop body for each loop iteration, which makes the index in each loop iteration constant. However, if the `for`-loop has a large body or it has many iterations, the unrolling can increase compile time and generate inefficient code.

If `A` and `B` are constant, the following code shows indexing into a heterogeneous cell array in a `for`-loop with constant bounds.

```
x = {1, 'mytext'};
for i = A:B
    disp(x{i});
end
```

Growing a Cell Array by Using {end + 1}

To grow a cell array X , you can use $X\{\text{end} + 1\}$. For example:

```
...
X = {1 2};
X{end + 1} = 'a';
...
```

When you use $\{\text{end} + 1\}$ to grow a cell array, follow these restrictions:

- Use only $\{\text{end} + 1\}$. Do not use $\{\text{end} + 2\}$, $\{\text{end} + 3\}$, and so on.
- Use $\{\text{end} + 1\}$ with vectors only. For example, the following code is not allowed because X is a matrix, not a vector:

```
...
X = {1 2; 3 4};
X{end + 1} = 5;
```

- Use $\{\text{end} + 1\}$ only with a variable. In the following code, $\{\text{end} + 1\}$ does not cause $\{1\ 2\ 3\}$ to grow. In this case, the code generator treats $\{\text{end} + 1\}$ as an out-of-bounds index into $X\{2\}$.

```
...
X = {'a' {1 2 3}};
X{2}{end + 1} = 4;
```

- When $\{\text{end} + 1\}$ grows a cell array in a loop, the cell array must be variable-size. Therefore, the cell array must be homogeneous on page 30-2.

This code is allowed because X is homogeneous.

```
...
X = {1 2};
for i=1:n
    X{end + 1} = 3;
end
...
```

This code is not allowed because X is heterogeneous.

```
...
X = {1 'a' 2 'b'};
for i=1:n
    X{end + 1} = 3;
end
...
```


Cell Array Contents

Cell arrays cannot contain `mxarrays`. In a cell array, you cannot store a value that an extrinsic function returns.

Passing Cell Arrays to External C/C++ Functions

You cannot pass a cell array to `coder.ceval`. If a variable is an input argument to `coder.ceval`, define the variable as an array or structure instead of as a cell array.

See Also

More About

- “Code Generation for Cell Arrays” on page 30-2
- “Differences Between Generated Code and MATLAB Code” on page 19-6

Primary Functions

Specify Properties of Entry-Point Function Inputs

| In this section... |
|---|
| “Why You Must Specify Input Properties” on page 31-2 |
| “Properties to Specify” on page 31-2 |
| “Rules for Specifying Properties of Primary Inputs” on page 31-3 |
| “Methods for Defining Properties of Primary Inputs” on page 31-4 |
| “Define Input Properties by Example at the Command Line” on page 31-4 |
| “Specify Constant Inputs at the Command Line” on page 31-6 |
| “Specify Variable-Size Inputs at the Command Line” on page 31-7 |
| “Input Type Specification and arguments blocks” on page 31-8 |

Why You Must Specify Input Properties

Fixed-Point Designer must determine the properties of all variables in the MATLAB files at compile time. To infer variable properties in MATLAB files, Fixed-Point Designer must be able to identify the properties of the inputs to the *primary* function, also known as the *top-level* or *entry-point* function. Therefore, if your primary function has inputs, you must specify the properties of these inputs, to Fixed-Point Designer. If your primary function has no input parameters, Fixed-Point Designer can compile your MATLAB file without modification. You do not need to specify properties of inputs to local functions or external functions called by the primary function.

Note Your primary function cannot be within a package. Create a wrapper function as the primary function outside the package. Call the desired function within the new function as the primary function.

Properties to Specify

If your primary function has inputs, you must specify the following properties for each input.

| For | Specify properties | | | | |
|---------------------------------|--|------|------------|-------------|--------|
| | Class | Size | Complexity | numerictype | fimath |
| Fixed-point inputs | ✓ | ✓ | ✓ | ✓ | ✓ |
| Each field in a structure input | <p>Specify properties for each field according to its class</p> <p>When a primary input is a structure, the code generator treats each field as a separate input. Therefore, you must specify properties for all fields of a primary structure input in the order that they appear in the structure definition:</p> <ul style="list-style-type: none"> • For each field of input structures, specify class, size, and complexity. • For each field that is fixed-point class, also specify <code>numerictype</code>, and <code>fimath</code>. | | | | |
| Other inputs | ✓ | ✓ | ✓ | | |

Default Property Values

Fixed-Point Designer assigns the following default values for properties of primary function inputs.

| Property | Default |
|------------|------------------------------|
| class | double |
| size | scalar |
| complexity | real |
| numericity | No default |
| fimath | MATLAB default fimath object |

Supported Classes

The following table presents the class names supported by Fixed-Point Designer.

| Class Name | Description |
|-------------|---|
| logical | Logical array of true and false values |
| char | Character array |
| int8 | 8-bit signed integer array |
| uint8 | 8-bit unsigned integer array |
| int16 | 16-bit signed integer array |
| uint16 | 16-bit unsigned integer array |
| int32 | 32-bit signed integer array |
| uint32 | 32-bit unsigned integer array |
| int64 | 64-bit signed integer array |
| uint64 | 64-bit unsigned integer array |
| single | Single-precision floating-point or fixed-point number array |
| double | Double-precision floating-point or fixed-point number array |
| struct | Structure array |
| embedded.fi | Fixed-point number array |

Rules for Specifying Properties of Primary Inputs

When specifying the properties of primary inputs, follow these rules:

- The order of elements in the cell array must correspond to the order in which inputs appear in the primary function signature. For example, the first element in the cell array defines the properties of the first primary function input.
- To generate fewer arguments than those arguments that occur in the MATLAB function, specify properties for only the number of arguments that you want in the generated function.
- If the MATLAB function has input arguments, to generate a function that has no input arguments, pass an empty cell array to `-args`.

- For each primary function input whose class is fixed point (`fi`), specify the input `numericType` and `fimath` properties.
- For each primary function input whose class is `struct`, specify the properties of each of its fields in the order that they appear in the structure definition.

Methods for Defining Properties of Primary Inputs

| Method | Advantages | Disadvantages |
|--|---|---|
| | <ul style="list-style-type: none"> • If you are working in a project, easy to use • Does not alter original MATLAB code • saves the definitions in the project file | <ul style="list-style-type: none"> • Not efficient for specifying memory-intensive inputs such as large structures and arrays |
| <p>“Define Input Properties by Example at the Command Line” on page 31-4</p> <hr/> <p>Note If you define input properties programmatically in the MATLAB file, you cannot use this method</p> | <ul style="list-style-type: none"> • Easy to use • Does not alter original MATLAB code • Designed for prototyping a function that has a few primary inputs | <ul style="list-style-type: none"> • Must be specified at the command line every time you invoke <code>fiaccl</code> (unless you use a script) • Not efficient for specifying memory-intensive inputs such as large structures and arrays |
| <p>“Define Input Properties Programmatically in the MATLAB File” (MATLAB Coder)</p> | <ul style="list-style-type: none"> • Integrated with MATLAB code; no need to redefine properties each time you invoke • Provides documentation of property specifications in the MATLAB code • Efficient for specifying memory-intensive inputs such as large structures | <ul style="list-style-type: none"> • Uses complex syntax • project files do not currently recognize properties defined programmatically. If you are using a project, you must reenter the input types in the project. |

Define Input Properties by Example at the Command Line

- “Command-Line Option `-args`” on page 31-4
- “Rules for Using the `-args` Option” on page 31-5
- “Specifying Properties of Primary Inputs by Example” on page 31-5
- “Specifying Properties of Primary Fixed-Point Inputs by Example” on page 31-5

Command-Line Option `-args`

The `fiaccl` function provides a command-line option `-args` for specifying the properties of primary (entry-point) function inputs as a cell array of example values or types. The cell array can be a variable or literal array of constant values. Using this option, you specify the properties of inputs at the same time as you generate code for the MATLAB function with `fiaccl`.

You can also create `coder.Type` objects interactively by using the Coder Type Editor. See “Create and Edit Input Types by Using the Coder Type Editor” (MATLAB Coder).

Rules for Using the `-args` Option

When using the `-args` command-line option to define properties by example, follow these rules:

- The order of elements in the cell array must correspond to the order in which inputs appear in the primary function signature. For example, the first element in the cell array defines the properties of the first primary function input.
- To generate fewer arguments than those arguments that occur in the MATLAB function, specify properties for only the number of arguments that you want in the generated function.
- If the MATLAB function has input arguments, to generate a function that has no input arguments, pass an empty cell array to `-args`.
- For each primary function input whose class is `fixed point (fi)`, specify the input `numeric type` and `fi math` properties.
- For each primary function input whose class is `struct`, specify the properties of each of its fields in the order that they appear in the structure definition.

Specifying Properties of Primary Inputs by Example

Consider a function that adds its two inputs:

```
function y = emcf(u,v) %#codegen
% The directive %#codegen indicates that you
% intend to generate code for this algorithm
y = u + v;
```

The following examples show how to specify different properties of the primary inputs `u` and `v` by example at the command line:

- Use a literal cell array of constants to specify that both inputs are real, scalar, fixed-point values:

```
fiaccel -o emcfx emcf ...
  -args {fi(0,1,16,15),fi(0,1,16,15)}
```

- Use a literal cell array of constants to specify that input `u` is an unsigned 16-bit, 1-by-4 vector and input `v` is a scalar, fixed-point value:

```
fiaccel -o emcfx emcf ...
  -args {zeros(1,4,'uint16'),fi(0,1,16,15)}
```

- Assign sample values to a cell array variable to specify that both inputs are real, unsigned 8-bit integer vectors:

```
a = fi([1;2;3;4],0,8,0)
b = fi([5;6;7;8],0,8,0)
ex = {a,b}
fiaccel -o emcfx emcf -args ex
```

Specifying Properties of Primary Fixed-Point Inputs by Example

Consider a function that calculates the square root of a fixed-point number:

```
function y = sqrtfi(x) %#codegen
y = sqrt(x);
```

To specify the properties of the primary fixed-point input `x` by example on the MATLAB command line, follow these steps:

- 1 Define the `numericType` properties for `x`, as in this example:

```
T = numericType('WordLength',32,...
    'FractionLength',23,'Signed',true);
```

- 2 Define the `fimath` properties for `x`, as in this example:

```
F = fimath('SumMode','SpecifyPrecision',...
    'SumWordLength',32,'SumFractionLength',23,...
    'ProductMode','SpecifyPrecision', ...
    'ProductWordLength',32,'ProductFractionLength',23);
```

- 3 Create a fixed-point variable with the `numericType` and `fimath` properties you defined, as in this example:

```
myeg = { fi(4.0,T,F) };
```

- 4 Compile the function `sqrtfi` using the `fiaccel` command, passing the variable `myeg` as the argument to the `-args` option, as in this example:

```
fiaccel sqrtfi -args myeg;
```

Specify Constant Inputs at the Command Line

If you know that your primary inputs do not change at run time, you can reduce overhead in the generated code by specifying that the primary inputs are constant values. Constant inputs are commonly used for flags that control how an algorithm executes and values that specify the sizes or types of data.

To specify that inputs are constants, use the `-args` command-line option with a `coder.Constant` object. To specify that an input is a constant with the size, class, complexity, and value of `constant_input`, use the following syntax:

```
-args {coder.Constant(constant_input)}
```

Calling Functions with Constant Inputs

`fiaccel` compiles constant function inputs into the generated code. As a result, the MEX function signature differs from the MATLAB function signature. At run time, you supply the constant argument to the MATLAB function, but not to the MEX function.

For example, consider the following function `identity` which copies its input to its output:

```
function y = identity(u) %#codegen
y = u;
```

To generate a MEX function `identity_mex` with a constant input, type the following command at the MATLAB prompt:

```
fiaccel -o identity_mex identity...
    -args {coder.Constant(fi(0.1,1,16,15))}
```

To run the MATLAB function, supply the constant argument as follows:

```
identity(fi(0.1,1,16,15))
```

You get the following result:


```
ans =
    0.1000
```

Now, try running the MEX function with this command:

```
identity_mex
```

You should get the same answer.

Specifying a Structure as a Constant Input

Suppose that you define a structure `tmp` in the MATLAB workspace to specify the dimensions of a matrix, as follows:

```
tmp = struct('rows', 2, 'cols', 3);
```

The following MATLAB function `rowcol` accepts a structure input `p` to define matrix `y`:

```
function y = rowcol(u,p) %#codegen
y = fi(zeros(p.rows,p.cols),1,16,15) + u;
```

The following example shows how to specify that primary input `u` is a double scalar variable and primary input `p` is a constant structure:

```
fiaccel rowcol ...
    -args {fi(0,1,16,15),coder.Constant(tmp)}
```

To run this code, use

```
u = fi(0.5,1,16,15)
y_m = rowcol(u,tmp)

y_mex = rowcol_mex(u)
```

Specify Variable-Size Inputs at the Command Line

Variable-size data is data whose size might change at run time. MATLAB supports bounded and unbounded variable-size data for code generation. Bounded variable-size data has fixed upper bounds. This data can be allocated statically on the stack or dynamically on the heap. Unbounded variable-size data does not have fixed upper bounds. This data must be allocated on the heap. You can define inputs to have one or more variable-size dimensions — and specify their upper bounds — using the `-args` option and `coder.typeof` function:

```
-args {coder.typeof(example_value, size_vector, variable_dims)}
```

Specifies a variable-size input with:

- Same class and complexity as *example_value*
- Same size and upper bounds as *size_vector*
- Variable dimensions specified by *variable_dims*

When you enable dynamic memory allocation, you can specify `Inf` in the size vector for dimensions with unknown upper bounds at compile time.

When *variable_dims* is a scalar, it is applied to all the dimensions, with the following exceptions:

- If the dimension is 1 or 0, which are fixed.
- If the dimension is unbounded, which is always variable size.

Specifying a Variable-Size Vector Input

- 1 Write a function that computes the sum of every n elements of a vector A and stores them in a vector B :

```
function B = nway(A,n) %#codegen
% Compute sum of every N elements of A and put them in B.

coder.extrinsic('error');
Tb = numerictype(1,32,24);
if ((mod(numel(A),n) == 0) && ...
    (n>=1 && n<=numel(A)))
    B = fi(zeros(1,numel(A)/n),Tb);
    k = 1;
    for i = 1 : numel(A)/n
        B(i) = sum(A(k + (0:n-1)));
        k = k + n;
    end
else
    B = fi(zeros(1,0),Tb);
    error('n<=0 or does not divide evenly');
end
```

- 2 Specify the first input A as a `fi` object. Its first dimension stays fixed in size and its second dimension can grow to an upper bound of 100. Specify the second input n as a double scalar.

```
fiaccel nway ...
-args {coder.typeof(fi(0,1,16,15,'SumMode','KeepLSB'),[1 100],1),0}...
-report
```

- 3 As an alternative, assign the `coder.typeof` expression to a MATLAB variable, then pass the variable as an argument to `-args`:

```
vareg = coder.typeof(fi(0,1,16,15,'SumMode','KeepLSB'),[1 100],1)
fiaccel nway -args {vareg, double(0)} -report
```

Input Type Specification and arguments blocks

Using function argument validation (`arguments` blocks) to specify input types of entry-point functions is not supported. Even if your entry-point function contains `arguments` blocks that validate the input arguments, you must specify the properties of these input arguments by using one of the three approaches listed in “Methods for Defining Properties of Primary Inputs” (MATLAB Coder).

Default Values for Entry-Point Function Inputs in Generated Code

The `arguments` block allows you to specify default values for one or more positional input arguments. Specifying a default value in the argument declaration makes a positional argument optional because MATLAB can use the default value when you do not pass a value in the function call. When you generate code by using the `codegen` command or accelerate fixed-point code by using the `fiaccel` command, you can choose to not specify the properties of one or more optional positional arguments that have constant default values. In such situations, the default values of these optional arguments are hard-coded in the generated code and these arguments do not appear in the generated code interface. For examples, see the following table.

| MATLAB Code | Generated Code |
|--|---|
| <pre>function out = useDefaults_1(a,b,c) arguments a (1,1) double = 3 b (1,1) double = 5 c (1,1) double = 7 end out = a + b + c; end</pre> | <pre>codegen command: codegen -config:lib -c useDefaults_1 -args {} -report Generated code: double useDefaults_1(void) { return 15.0; }</pre> |
| <pre>function out = useDefaults_2(a,b,c) arguments a (1,1) double b (1,1) double = 5 c (1,1) double = 7 end out = a + b + c; end</pre> | <pre>codegen command: codegen -config:lib -c useDefaults_2 -args 0 -report Generated code: double useDefaults_2(double a) { return (a + 5.0) + 7.0; }</pre> |
| | <pre>codegen command: codegen -config:lib -c useDefaults_2 -args {0,0} -rep Generated code: double useDefaults_2(double a, double b) { return (a + b) + 7.0; }</pre> |

See Also

More About

- “Specify Objects as Inputs” on page 15-26
- “Specify Cell Array Inputs at the Command Line” on page 12-42
- “Specify Number of Entry-Point Function Input or Output Arguments to Generate” on page 17-7

Define Input Properties Programmatically in the MATLAB File

For code generation, you can use the MATLAB `assert` function to define properties of primary function inputs directly in your MATLAB file.

How to Use `assert` with MATLAB Coder

Use the `assert` function to invoke standard MATLAB functions for specifying the class, size, and complexity of primary function inputs.

When specifying input properties using the `assert` function, use one of the following methods. Use the exact syntax that is provided; do not modify it.

- “Specify Any Class” on page 31-10
- “Specify `fi` Class” on page 31-10
- “Specify Structure Class” on page 31-11
- “Specify Cell Array Class” on page 31-11
- “Specify Fixed Size” on page 31-12
- “Specify Scalar Size” on page 31-12
- “Specify Upper Bounds for Variable-Size Inputs” on page 31-12
- “Specify Inputs with Fixed- and Variable-Size Dimensions” on page 31-12
- “Specify Size of Individual Dimensions” on page 31-13
- “Specify Real Input” on page 31-13
- “Specify Complex Input” on page 31-13
- “Specify `numerictype` of Fixed-Point Input” on page 31-13
- “Specify `fimath` of Fixed-Point Input” on page 31-14
- “Specify Multiple Properties of Input” on page 31-14

Specify Any Class

```
assert ( isa ( param, 'class_name' ) )
```

Sets the input parameter `param` to the MATLAB class `class_name`. For example, to set the class of input `U` to a 32-bit signed integer, call:

```
...  
assert(isa(U,'int32'));  
...
```

Specify `fi` Class

```
assert ( isfi ( param ) )  
assert ( isa ( param, 'embedded.fi' ) )
```

Sets the input parameter `param` to the MATLAB class `fi` (fixed-point numeric object). For example, to set the class of input `U` to `fi`, call:

```
...
assert(isfi(U));
...
```

or

```
...
assert(isa(U, 'embedded.fi'));
...
```

You must specify both the `fi` class and the `numericType`. See “Specify numericType of Fixed-Point Input” on page 31-13. You can also set the `fimath` properties, see “Specify fimath of Fixed-Point Input” on page 31-14. If you do not set the `fimath` properties, `codegen` uses the MATLAB default `fimath` value.

Specify Structure Class

```
assert ( isstruct ( param ) )
assert ( isa ( param, 'struct' ) )
```

Sets the input parameter *param* to the MATLAB class `struct` (structure). For example, to set the class of input `U` to a `struct`, call:

```
...
assert(isstruct(U));
...
```

or

```
...
assert(isa(U, 'struct'));
...
```

If you set the class of an input parameter to `struct`, you must specify the properties of all fields in the order that they appear in the structure definition.

Specify Cell Array Class

```
assert(iscell( param))
assert(isa(param, 'cell'))
```

Sets the input parameter *param* to the MATLAB class `cell` (cell array). For example, to set the class of input `C` to a `cell`, call:

```
...
assert(iscell(C));
...
```

or

```
...
assert(isa(C, 'cell'));
...
```

To specify the properties of cell array elements, see “Specifying Properties of Cell Arrays” on page 31-16.

Specify Fixed Size

```
assert ( all ( size (param) == [dims ] ) )
```

Sets the input parameter *param* to the size that dimensions *dims* specifies. For example, to set the size of input U to a 3-by-2 matrix, call:

```
...  
assert(all(size(U)== [3 2]));  
...
```

Specify Scalar Size

```
assert ( isscalar (param ) )  
assert ( all ( size (param) == [ 1 ] ) )
```

Sets the size of input parameter *param* to scalar. To set the size of input U to scalar, call:

```
...  
assert(isscalar(U));  
...
```

or

```
...  
assert(all(size(U)== [1]));  
...
```

Specify Upper Bounds for Variable-Size Inputs

```
assert ( all(size(param)<=[N0 N1 ...]));  
assert ( all(size(param)<[N0 N1 ...]));
```

Sets the upper-bound size of each dimension of input parameter *param*. To set the upper-bound size of input U to be less than or equal to a 3-by-2 matrix, call:

```
assert(all(size(U)<=[3 2]));
```

Note You can also specify upper bounds for variable-size inputs using `coder. varsize`.

Specify Inputs with Fixed- and Variable-Size Dimensions

```
assert ( all(size(param)>=[M0 M1 ...]));  
assert ( all(size(param)<=[N0 N1 ...]));
```

When you use `assert(all(size(param)>=[M0 M1 ...]))` to specify the lower-bound size of each dimension of an input parameter:

- You must also specify an upper-bound size for each dimension of the input parameter.
- For each dimension, *k*, the lower-bound *M_k* must be less than or equal to the upper-bound *N_k*.
- To specify a fixed-size dimension, set the lower and upper bound of a dimension to the same value.
- Bounds must be nonnegative.

To fix the size of the first dimension of input U to 3 and set the second dimension as variable size with upper bound of 2, call:

```
assert(all(size(U)>=[3 0]));
assert(all(size(U)<=[3 2]));
```

Specify Size of Individual Dimensions

```
assert (size(param, k)==Nk);
assert (size(param, k)<=Nk);
assert (size(param, k)<Nk);
```

You can specify individual dimensions and all dimensions simultaneously. You can also specify individual dimensions instead of specifying all dimensions simultaneously. The following rules apply:

- You must specify the size of each dimension at least once.
- The last dimension specification takes precedence over earlier specifications.

Sets the upper-bound size of dimension *k* of input parameter *param*. To set the upper-bound size of the first dimension of input *U* to 3, call:

```
assert(size(U,1)<=3)
```

To fix the size of the second dimension of input *U* to 2, call:

```
assert(size(U,2)==2)
```

Specify Real Input

```
assert ( isreal (param ) )
```

Specifies that the input parameter *param* is real. To specify that input *U* is real, call:

```
...
assert(isreal(U));
...
```

Specify Complex Input

```
assert ( ~isreal (param ) )
```

Specifies that the input parameter *param* is complex. To specify that input *U* is complex, call:

```
...
assert(~isreal(U));
...
```

Specify numerictype of Fixed-Point Input

```
assert ( isequal ( numerictype ( fiparam ), T ) )
```

Sets the *numerictype* properties of *fi* input parameter *fiparam* to the *numerictype* object *T*. For example, to specify the *numerictype* property of fixed-point input *U* as a signed *numerictype* object *T* with 32-bit word length and 30-bit fraction length, use the following code:

```
%#codegen
...
% Define the numerictype object.
T = numerictype(1, 32, 30);
```

```
% Set the numerictype property of input U to T.
assert(isequal(numerictype(U),T));
...
```

Specifying the `numerictype` for a variable does not automatically specify that the variable is fixed point. You must specify both the `fi` class and the `numerictype`.

Specify `fimath` of Fixed-Point Input

```
assert ( isequal ( fimath ( fiparam ), F ) )
```

Sets the `fimath` properties of `fi` input parameter `fiparam` to the `fimath` object `F`. For example, to specify the `fimath` property of fixed-point input `U` so that it saturates on integer overflow, use the following code:

```
 %#codegen
...
% Define the fimath object.
F = fimath('OverflowMode','saturate');

% Set the fimath property of input U to F.
assert(isequal(fimath(U),F));
...
```

If you do not specify the `fimath` properties using `assert`, `codegen` uses the MATLAB default `fimath` value.

Specify Multiple Properties of Input

```
assert ( function1 ( params ) &&
         function2 ( params ) &&
         function3 ( params ) && ... )
```

Specifies the class, size, and complexity of one or more inputs using a single `assert` function call. For example, the following code specifies that input `U` is a double, complex, 3-by-3 matrix, and input `V` is a 16-bit unsigned integer:

```
 %#codegen
...
assert(isa(U,'double') &&
       ~isreal(U) &&
       all(size(U) == [3 3]) &&
       isa(V,'uint16'));
...
```

Rules for Using `assert` Function

When using the `assert` function to specify the properties of primary function inputs, follow these rules:

- Call `assert` functions at the beginning of the primary function, before control-flow operations such as `if` statements or subroutine calls.
- Do not call `assert` functions inside conditional constructs, such as `if`, `for`, `while`, and `switch` statements.
- For a fixed-point input, you must specify both the `fi` class and the `numerictype`. See “Specify `numerictype` of Fixed-Point Input” on page 31-13. You can also set the `fimath` properties. See

“Specify fimath of Fixed-Point Input” on page 31-14. If you do not set the `fimath` properties, `codegen` uses the MATLAB default `fimath` value.

- If you set the class of an input parameter to `struct`, you must specify the class, size, and complexity of all fields in the order that they appear in the structure definition.
- When you use `assert(all(size(param)>=[M0 M1 ...]))` to specify the lower-bound size of each dimension of an input parameter:
 - You must also specify an upper-bound size for each dimension of the input parameter.
 - For each dimension, k , the lower-bound M_k must be less than or equal to the upper-bound N_k .
 - To specify a fixed-size dimension, set the lower and upper bound of a dimension to the same value.
 - Bounds must be nonnegative.
- If you specify individual dimensions, the following rules apply:
 - You must specify the size of each dimension at least once.
 - The last dimension specification takes precedence over earlier specifications.

Specifying General Properties of Primary Inputs

In the following code excerpt, a primary MATLAB function `mcspecgram` takes two inputs: `pennywhistle` and `win`. The code specifies the following properties for these inputs.

| Input | Property | Value |
|--------------|------------|--------------------|
| pennywhistle | class | int16 |
| | size | 220500-by-1 vector |
| | complexity | real (by default) |
| win | class | double |
| | size | 1024-by-1 vector |
| | complexity | real (by default) |

```

%#codegen
function y = mcspecgram(pennywhistle,win)
nx = 220500;
nfft = 1024;
assert(isa(pennywhistle,'int16'));
assert(all(size(pennywhistle) == [nx 1]));
assert(isa(win, 'double'));
assert(all(size(win) == [nfft 1]));
...

```

Alternatively, you can combine property specifications for one or more inputs inside `assert` commands:

```

%#codegen
function y = mcspecgram(pennywhistle,win)
nx = 220500;
nfft = 1024;
assert(isa(pennywhistle,'int16') && all(size(pennywhistle) == [nx 1]));
assert(isa(win, 'double') && all(size(win) == [nfft 1]));
...

```

Specifying Properties of Primary Fixed-Point Inputs

To specify fixed-point inputs, you must install Fixed-Point Designer software.

In the following example, the primary MATLAB function `mcsqrtfi` takes one fixed-point input `x`. The code specifies the following properties for this input.

| Property | Value |
|-------------------------|--|
| <code>class</code> | <code>fi</code> |
| <code>numericity</code> | numericity object <code>T</code> , as specified in the primary function |
| <code>fimath</code> | <code>fimath</code> object <code>F</code> , as specified in the primary function |
| <code>size</code> | scalar |
| <code>complexity</code> | real (by default) |

```
function y = mcsqrtfi(x) %#codegen
T = numericity('WordLength',32,'FractionLength',23,...
              'Signed',true);
F = fimath('SumMode','SpecifyPrecision',...
          'SumWordLength',32,'SumFractionLength',23,...
          'ProductMode','SpecifyPrecision',...
          'ProductWordLength',32,'ProductFractionLength',23);
assert(isfi(x));
assert(isequal(numericity(x),T));
assert(isequal(fimath(x),F));

y = sqrt(x);
```

You must specify both the `fi` class and the `numericity`.

Specifying Properties of Cell Arrays

To specify the class `cell` (cell array), use one of the following syntaxes:

```
assert(iscell(param))
assert(isa(param, 'cell'))
```

For example, to set the class of input `C` to `cell`, use:

```
...
assert(iscell(C));
...
```

or

```
...
assert(isa(C, 'cell'));
...
```

You can also specify the size of the cell array and the properties of the cell array elements. The number of elements that you specify determines whether the cell array is homogeneous or heterogeneous. See “Code Generation for Cell Arrays” (MATLAB Coder).

If you specify the properties of the first element only, the cell array is homogeneous. For example, the following code specifies that C is a 1x3 homogeneous cell array whose elements are 1x1 double.

```
...
assert(isa(C, 'cell'));
assert(all(size(C) == [1 3]));
assert(isa(C{1}, 'double'));
...
```

If you specify the properties of the first element only, but also assign a structure type name to the cell array, the cell array is heterogeneous. Each element has the properties of the first element. For example, the following code specifies that C is a 1x3 heterogeneous cell array. Each element is a 1x1 double.

```
...
assert(isa(C, 'cell'));
assert(all(size(C) == [1 3]));
assert(isa(C{1}, 'double'));
coder.cstructname(C, 'myname');
...
```

If you specify the properties of each element, the cell array is heterogeneous. For example, the following code specifies a 1x2 heterogeneous cell array whose first element is 1x1 char and whose second element is 1x3 double.

```
...
assert(isa(C, 'cell'));
assert(all(size(C) == [1 2]));
assert(isa(C{1}, 'char'));
assert(all(size(C{2}) == [1 3]));
assert(isa(C{2}, 'double'));
...
```

If you specify more than one element, you cannot specify that the cell array is variable size, even if all elements have the same properties. For example, the following code specifies a variable-size cell array. Because the code specifies the properties of the first and second elements, code generation fails.

```
...
assert(isa(C, 'cell'));
assert(all(size(C) <= [1 2]));
assert(isa(C{1}, 'double'));
assert(isa(C{2}, 'double'));
...
```

In the previous example, if you specify the first element only, you can specify that the cell array is variable-size. For example:

```
...
assert(isa(C, 'cell'));
assert(all(size(C) <= [1 2]));
assert(isa(C{1}, 'double'));
...
```

Specifying Class and Size of Scalar Structure

Suppose that you define `S` as the following scalar MATLAB structure:

```
S = struct('r',double(1),'i',int8(4));
```

The following code specifies the properties of the function input `S` and its fields:

```
function y = fcn(S) %#codegen

% Specify the class of the input as struct.
assert(isstruct(S));

% Specify the class and size of the fields r and i
% in the order in which you defined them.
assert(isa(S.r,'double'));
assert(isa(S.i,'int8'));
...
```

In most cases, when you do not explicitly specify values for properties, MATLAB Coder uses defaults—except for structure fields. The only way to name a field in a structure is to set at least one of its properties. At a minimum, you must specify the class of a structure field.

Specifying Class and Size of Structure Array

For structure arrays, you must choose a representative element of the array for specifying the properties of each field. For example, assume that you have defined `S` as the following 1-by-2 array of MATLAB structures:

```
S = struct('r',{double(1), double(2)},'i',{int8(4), int8(5)});
```

The following code specifies the class and size of each field of structure input `S` by using the first element of the array:

```
%#codegen
function y = fcn(S)

% Specify the class of the input S as struct.
assert(isstruct(S));

% Specify the size of the fields r and i
% based on the first element of the array.
assert(all(size(S) == [1 2]));
assert(isa(S(1).r,'double'));
assert(isa(S(1).i,'int8'));
```

The only way to name a field in a structure is to set at least one of its properties. At a minimum, you must specify the class of all fields.

Create and Edit Input Types by Using the Coder Type Editor

Fixed-Point Designer must determine the properties of all variables in the MATLAB files at compile time. To infer variable properties in MATLAB files, Fixed-Point Designer must be able to identify the properties of the inputs to the top-level MATLAB functions, also known as *entry-point functions*. Therefore, if your entry-point function has inputs, you must specify the properties of these inputs.

One of the ways to specify the properties of an input argument is by using a `coder.Type` object that contains information about class, size, and complexity (and sometimes other properties) of the argument. You can create and edit `coder.Type` objects programmatically at the command line, or interactively by using the Coder Type Editor.

For more information about creating `coder.Type` objects at the command line, see `coder.typeof` and `coder.newtype`.

Note To create and edit composite types such as structures and cell arrays, or types that have many customizable parameters such as `embedded.fi`, use the Coder Type Editor. Examples of such types are shown later in this topic.

Open the Coder Type Editor

To launch the Coder Type Editor, do one of the following:

- Launch an empty type editor by using the `coderTypeEditor` command:

```
coderTypeEditor
```

- Open the type editor pre-populated with `coder.Type` objects corresponding to the workspace variables `var1`, `var2`, and `var3` by typing:

```
coderTypeEditor var1 var2 var3
```


- Open a `coder.Type` object `myType` that already exists in your base MATLAB workspace:
 - Double click `myType` in the workspace.
 - Display `myType` at the command line and click the *Edit Type Object* link that appears at the end of the display.
 - Use this command at the MATLAB command line:

```
open myType
```

Common Editor Actions

By using the toolbar buttons in the type editor, you can perform these actions:

- Create a new type by clicking **New Type** and specifying the type, size, complexity, and other properties of the `coder.Type` object.
- Convert an existing variable to a type by clicking **From Variable** and specifying a variable that already exists in the base workspace.
- Create a new type from an example value by clicking **From Example** and entering MATLAB code that the software converts to a `coder.Type` object.


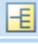





- Load all coder .Type objects from the base workspace to the **Type Browser** pane of the type editor by clicking **Load All**.
- Edit an existing type by selecting it in the **Type Browser** and modifying its properties.
- Save all coder .Type objects in the type editor by clicking **Save All**.
- Remove a selected type from **Type Browser** by clicking **Delete**. Alternatively, remove all types from the **Type Browser** by clicking **Delete > Delete all**. Deleting a coder .Type object from the **Type Browser** does not delete the object from the base MATLAB workspace.
- Export a MATLAB script that contains the code to recreate all the types by clicking **Share > MATLAB Script**. Or, create a MAT file that contains all the types by clicking **Share > MAT File**.
- Undo and redo your last action in the type editor by using the  buttons.

These are some additional actions that you can perform in the Coder Type Editor:

- In both the **Type Browser** pane and the **Type Properties** pane, copy a type object and paste it either as a new type or a field of an existing structure type. You can also copy the properties of one existing type into another existing type.
- Change the order of fields of a structure type. View the type in the properties pane and use drag-and-drop action.

Type Browser Pane

The **Type Browser** pane shows the name, class, and size of the coder .Type objects that are currently loaded in the type editor. For composite types such as structures, cell arrays, or classes, you can expand the display of the code .Type object in the **Type Browser** pane. The expanded view shows the name, class, and complexity of the individual fields or properties of the composite type.

| TYPE BROWSER | | |
|---|--------|-------|
| Name | Class | Size |
|  aType | double | 5×100 |
| ▼  bType | struct | 1×1 |
|  f1 | double | 1×1 |
|  f2 | char | 1×12 |
|  f3 | uint8 | 1×1 |
| ▼  cType | cell | 2×3 |
|  {} | double | 1×1 |

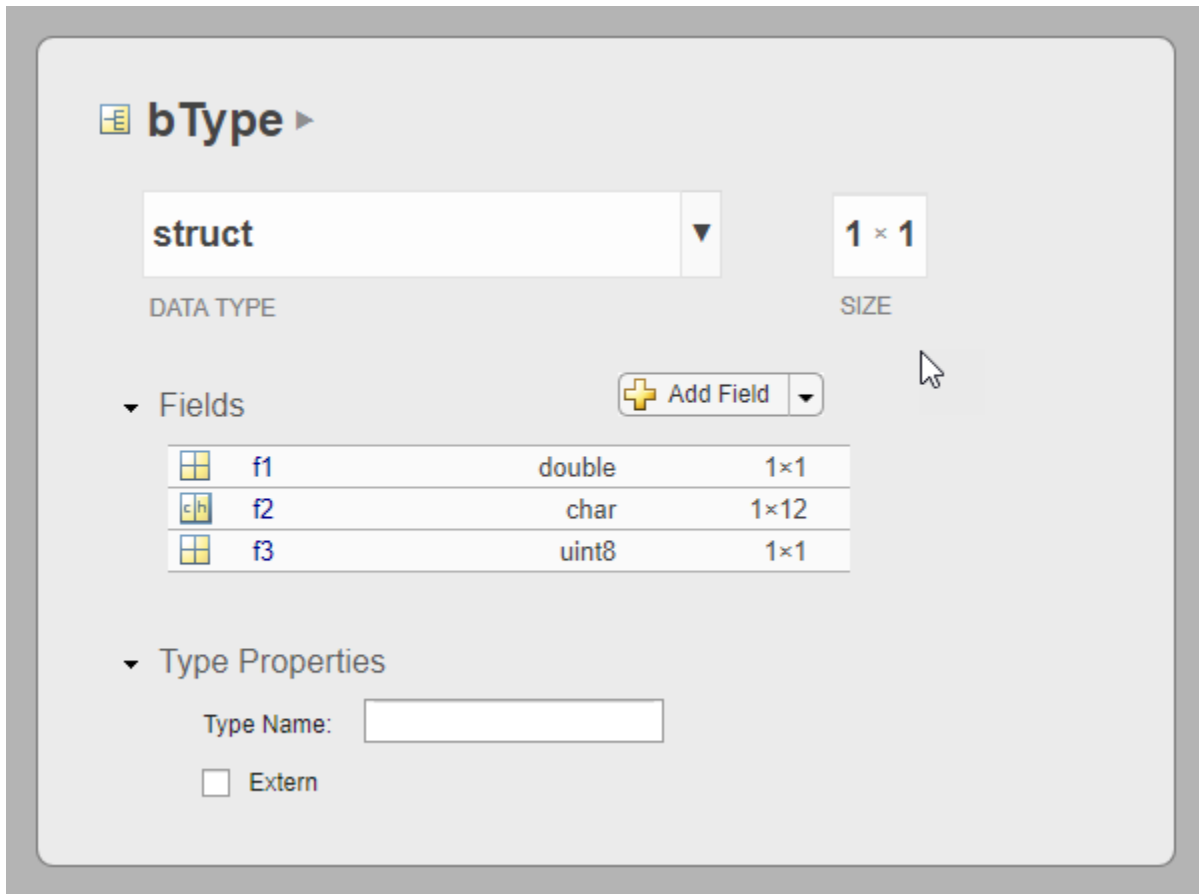
Visual Indicators on the Type Browser

| Indicator | Description |
|-----------|--|
| expander | The type has fields or properties that you can see by clicking the expander. |

| Indicator | Description |
|-----------|---|
| {:} | Homogeneous cell array (all elements have the same properties). |
| {n} | nth element of a heterogeneous cell array. |
| :n | Variable-size dimension with an upper bound of n. |
| :inf | Variable-size dimension that is unbounded. |

Type Properties Pane

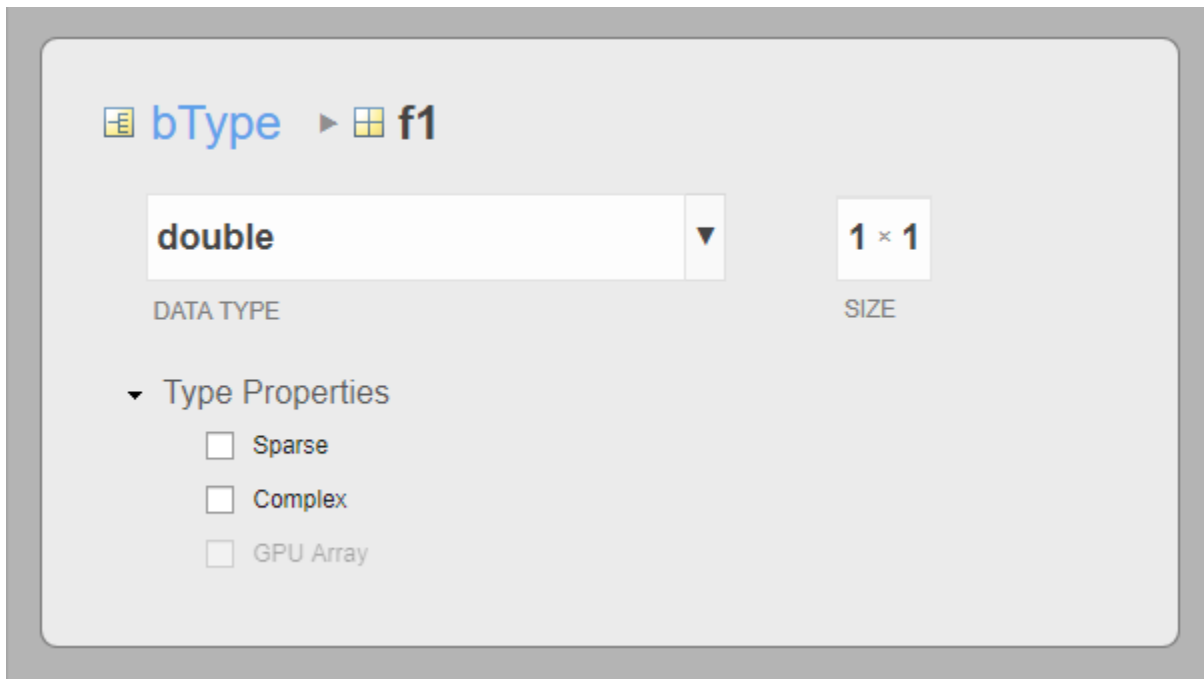
The type properties pane displays the class (data type), size, and other properties of the coder .Type object that is currently selected in the **Type Browser**. For composite types such as structures and classes, this pane also shows the name, class, and size of each constituent field or property.



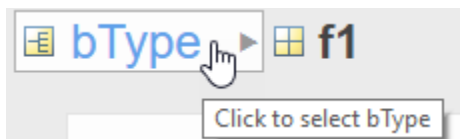
To edit the name, class, and size of a field in place, double-click the item.

| | | | |
|--|---------------------------------|--------|------|
| | <input type="text" value="f1"/> | double | 1x1 |
| | f2 | char | 1x12 |
| | f3 | uint8 | 1x1 |

Alternatively, click a field. The view in the type editor pane changes to display the properties of that field. Edit name, class(data type), size, or other properties in the pane.



The breadcrumb shows the nested path to the field that is currently open in the type properties pane. Click a field in the breadcrumb to display it in the pane. You can also edit the name of a type directly in the breadcrumb.



MATLAB Code Pane

The MATLAB Code pane displays the MATLAB script that creates the `coder.Type` object that is currently selected in the **Type Browser**. To automate the creation of this type, copy this script and include it in your build script.

```

MATLAB CODE
1 childTypes.f1 = coder.newtype('double', [1 1], [0 0]);
2 childTypes.f2 = coder.newtype('char', [1 12], [0 0]);
3 childTypes.f3 = coder.newtype('uint8', [1 1], [0 0]);
4 bType = coder.newtype('struct', childTypes, [1 1], [0 0]);
5
   clear childTypes;

```


See Also

`coderTypeEditor` | `coder.typeof` | `coder.newtype`

More About

- “Specify Properties of Entry-Point Function Inputs” on page 31-2

System Objects Supported for Code Generation

Code Generation for System Objects

You can generate C and C++ code for a subset of System objects provided by the following toolboxes.

| Toolbox Name | See |
|--------------------------------|---|
| Communications Toolbox™ | “System Objects in MATLAB Code Generation” (MATLAB Coder) |
| Computer Vision Toolbox™ | “System Objects in MATLAB Code Generation” (MATLAB Coder) |
| DSP System Toolbox | “System Objects in MATLAB Code Generation” (MATLAB Coder) |
| Image Acquisition Toolbox™ | <ul style="list-style-type: none"> • <code>imaq.VideoDevice</code> • “Code Generation with VideoDevice System Object” (Image Acquisition Toolbox) |
| Phased Array System Toolbox™ | “Code Generation” (Phased Array System Toolbox) |
| System Identification Toolbox™ | “Generate Code for Online Parameter Estimation in MATLAB” (System Identification Toolbox) |
| WLAN Toolbox™ | “System Objects in MATLAB Code Generation” (MATLAB Coder) |

To use these System objects, you need to install the requisite toolbox. For a list of System objects supported for C and C++ code generation, see “Functions and Objects Supported for C/C++ Code Generation” on page 26-2.

System objects are MATLAB object-oriented implementations of algorithms. They extend MATLAB by enabling you to model dynamic systems represented by time-varying algorithms. System objects are well integrated into the MATLAB language, regardless of whether you are writing simple functions, working interactively in the command window, or creating large applications.

In contrast to MATLAB functions, System objects automatically manage state information, data indexing, and buffering, which is particularly useful for iterative computations or stream data processing. This enables efficient processing of long data sets. For general information about MATLAB objects, see “Classes”.

Half Precision in MATLAB

- “Half Precision Code Generation Support” on page 33-2
- “Generate Native Half-Precision C Code Using MATLAB Coder” on page 33-13
- “What is Half Precision?” on page 33-19

Half Precision Code Generation Support

To assign a half-precision data type to a number or variable, use the `half` constructor. A half-precision data type occupies 16 bits of memory, but its floating-point representation enables it to handle wider dynamic ranges than integer or fixed-point data types of the same size. For more information, see “Floating-Point Numbers” on page 35-20.

A subset of MATLAB functions are supported for use with half-precision inputs. Additionally, some functions support code generation with half-precision data types. C and C++ code generation requires MATLAB Coder. CUDA code generation for NVIDIA® GPUs requires GPU Coder. Supported functions appear in alphabetical order in the following table. MATLAB System object supports half-precision data type and MATLAB System block supports half-precision data type with real values. For general information regarding code generation with half precision, see `half`.

| Function | MATLAB Simulation Support | C/C++ Code Generation Support | GPU Code Generation Support |
|--------------------------|---------------------------|---|--|
| <code>abs</code> | ✓ | ✓ | ✓ |
| <code>acos</code> | ✓ | ✓ | ✓ |
| <code>acosh</code> | ✓ | ✓ | ✓ |
| <code>activations</code> | ✓ | ✓ Half inputs are cast to single precision and computations are performed in single precision. | ✓ Half inputs are cast to single precision and computations are performed in single precision. To perform computations in half, set the library target to 'tensorrt' and set the data type to 'FP16' in <code>coder.DeepLearningConfig</code> . |
| <code>all</code> | ✓ | ✓ | ✓ |
| <code>allfinite</code> | ✓ | ✓ | ✓ |
| <code>and, &</code> | ✓ | ✓ | ✓ |
| Short-Circuit AND | ✓ | ✓ | ✓ |
| <code>any</code> | ✓ | ✓ | ✓ |
| <code>anynan</code> | ✓ | ✓ | ✓ |
| <code>area</code> | ✓ | | |
| <code>asin</code> | ✓ | ✓ | ✓ |
| <code>asinh</code> | ✓ | ✓ | ✓ |
| <code>atan</code> | ✓ | ✓ | ✓ |
| <code>atan2</code> | ✓ | ✓ | ✓ |
| <code>atanh</code> | ✓ | ✓ | ✓ |

| Function | MATLAB Simulation Support | C/C++ Code Generation Support | GPU Code Generation Support |
|-------------|--|---|--|
| bar | ✓ | | |
| barh | ✓ | | |
| cast | ✓ Supported syntax: cast(__, 'half') cast(__, 'like', p) | ✓ Supported syntax: cast(__, 'half') cast(__, 'like', p) | ✓ Supported syntax: cast(__, 'half') cast(__, 'like', p) |
| cat | ✓ | ✓ <ul style="list-style-type: none"> • Dimension argument must be a constant. • Dimension argument cannot be half precision. | ✓ <ul style="list-style-type: none"> • Dimension argument must be a constant. • Dimension argument cannot be half precision. |
| ceil | ✓ | ✓ | ✓ |
| cell | ✓ | ✓ | ✓ |
| chol | ✓ | | |
| circshift | ✓ | ✓ | ✓ |
| classify | ✓ | ✓ Half inputs are cast to single precision and computations are performed in single precision. | ✓ Half inputs are cast to single precision and computations are performed in single precision. To perform computations in half, set the library target to 'tensorrt' and set the data type to 'FP16' in coder.DeepLearningConfig. |
| coder.ceval | | ✓ | ✓ |
| colon, : | ✓ | ✓ | ✓ |
| complex | ✓ | ✓ | |
| conj | ✓ | ✓ | ✓ |
| conv | ✓ | ✓ | ✓ |
| conv2 | ✓ | ✓ | ✓ |
| cos | ✓ | ✓ | ✓ |
| cosh | ✓ | ✓ | ✓ |
| cospi | ✓ | ✓ | ✓ |

| Function | MATLAB Simulation Support | C/C++ Code Generation Support | GPU Code Generation Support |
|------------|--|---|---|
| ctranspose | ✓ | ✓ | ✓ |
| cumsum | ✓ | | |
| dot | ✓ | | |
| double | ✓ | ✓ | ✓ |
| empty | ✓ | | |
| eps | ✓ Supported syntax: eps('half') eps(half(1)) eps('like',half(1)) | ✓ eps(half(1)) | ✓ eps(half(1)) |
| eq, == | ✓ | ✓ | ✓ |
| exp | ✓ | ✓ | ✓ |
| expm1 | ✓ | ✓ | ✓ |
| eye | ✓ Supported syntax: eye(_, 'half') eye(_, 'like', p) | ✓ Supported syntax: eye(_, 'half') eye(_, 'like', p) where p is half precision. Other input arguments cannot be half precision. | ✓ Supported syntax: eye(_, 'half') eye(_, 'like', p) where p is half precision. Other input arguments cannot be half precision. |
| fft | ✓ | ✓ | |
| fft2 | ✓ | ✓ | |
| fftn | ✓ | ✓ | |
| fftshift | ✓ | ✓ | ✓ |
| fix | ✓ | ✓ | ✓ |
| flintmax | ✓ Supported syntax: flintmax('half') flintmax('like',half(1)) | | |

| Function | MATLAB Simulation Support | C/C++ Code Generation Support | GPU Code Generation Support |
|-----------|---|---|---|
| flip | ✓ | ✓ Dimension argument cannot be half precision. | ✓ Dimension argument cannot be half precision. |
| flipplr | ✓ | ✓ | ✓ |
| flipud | ✓ | ✓ | ✓ |
| floor | ✓ | ✓ | ✓ |
| fma | ✓ Complex half-precision inputs are not supported. | ✓ Complex half-precision inputs are not supported. | ✓ Complex half-precision inputs are not supported. |
| fplot | ✓ | | |
| ge, >= | ✓ | ✓ | ✓ |
| gt, > | ✓ | ✓ | ✓ |
| half | ✓ | ✓ | ✓ |
| horzcat | ✓ | ✓ | ✓ |
| hypot | ✓ | ✓ | ✓ |
| ifft | ✓ | ✓ | |
| ifft2 | ✓ | ✓ | |
| ifftn | ✓ | ✓ | |
| ifftshift | ✓ | ✓ | ✓ |
| imag | ✓ | ✓ | |
| Inf | ✓ Supported syntax: Inf(__, 'half') Inf(__, 'like', p) | ✓ Supported syntax: Inf(__, 'half') Inf(__, 'like', p) | ✓ Supported syntax: Inf(__, 'half') Inf(__, 'like', p) |
| int16 | ✓ | ✓ | ✓ |
| int32 | ✓ | ✓ | ✓ |
| int64 | ✓ | ✓ | ✓ |
| int8 | ✓ | ✓ | ✓ |
| isa | ✓ | ✓ | ✓ |
| iscolumn | ✓ | ✓ | ✓ |
| isempty | ✓ | ✓ | ✓ |
| isequal | ✓ | ✓ | ✓ |
| isequaln | ✓ | ✓ | ✓ |

| Function | MATLAB Simulation Support | C/C++ Code Generation Support | GPU Code Generation Support |
|-------------|---|--|--|
| isfinite | ✓ | ✓ | ✓ |
| isfloat | ✓ | ✓ | ✓ |
| isinf | ✓ | ✓ | ✓ |
| isinteger | ✓ | ✓ | ✓ |
| islogical | ✓ | ✓ | ✓ |
| ismatrix | ✓ | ✓ | ✓ |
| isnan | ✓ | ✓ | ✓ |
| isnumeric | ✓ | ✓ | ✓ |
| isobject | ✓ | ✓ | ✓ |
| | Returns true with half-precision input. | Returns false with half-precision input. | Returns false with half-precision input. |
| isreal | ✓ | ✓ | ✓ |
| isrow | ✓ | ✓ | ✓ |
| isscalar | ✓ | ✓ | ✓ |
| issorted | ✓ | | |
| isvector | ✓ | ✓ | ✓ |
| ldivide | ✓ | ✓ | ✓ |
| le, <= | ✓ | ✓ | ✓ |
| length | ✓ | ✓ | ✓ |
| line | ✓ | | |
| log | ✓ | ✓ | ✓ |
| log10 | ✓ | ✓ | ✓ |
| log1p | ✓ | ✓ | ✓ |
| log2 | ✓ | ✓ | ✓ |
| | | Two output syntax is not supported. | Two output syntax is not supported. |
| logical | ✓ | ✓ | ✓ |
| lt, < | ✓ | ✓ | ✓ |
| lu | ✓ | | |
| max | ✓ | ✓ | ✓ |
| mean | ✓ | ✓ | ✓ |
| min | ✓ | ✓ | ✓ |
| minus, - | ✓ | ✓ | ✓ |
| mldivide, \ | ✓ | | |
| | Left-hand side must be scalar | | |

| Function | MATLAB Simulation Support | C/C++ Code Generation Support | GPU Code Generation Support |
|------------------|---|---|--|
| mod | ✓ | ✓ | ✓ |
| mrdivide, / | ✓ Right-hand side must be scalar | ✓ Right-hand side must be scalar | ✓ Right-hand side must be scalar |
| mtimes, * | ✓ | ✓ | ✓ For GPU Code generation, you can perform half-precision matrix multiplication with real inputs. |
| NaN | ✓ Supported syntax: NaN(__, 'half') NaN(__, 'like', p) | ✓ Supported syntax: NaN(__, 'half') NaN(__, 'like', p) | ✓ Supported syntax: NaN(__, 'half') NaN(__, 'like', p) |
| ndims | ✓ | ✓ | ✓ |
| ne, ~= | ✓ | ✓ | ✓ |
| not | ✓ | ✓ | ✓ |
| numel | ✓ | ✓ | ✓ |
| ones | ✓ Supported syntax: ones(__, 'half') ones(__, 'like', p) | ✓ Supported syntax: ones(__, 'half') ones(__, 'like', p) | ✓ Supported syntax: ones(__, 'half') ones(__, 'like', p) |
| or, | ✓ | ✓ | ✓ |
| Short-Circuit OR | ✓ | ✓ | ✓ |
| permute | ✓ | ✓ | ✓ |
| plot | ✓ | | |
| plot3 | ✓ | | |
| plotmatrix | ✓ | | |
| plus, + | ✓ | ✓ | ✓ |
| pow10 | ✓ | ✓ | ✓ |
| pow2 | ✓ | ✓ | ✓ |
| power, .^ | ✓ | ✓ | ✓ |

| Function | MATLAB Simulation Support | C/C++ Code Generation Support | GPU Code Generation Support |
|-----------------------|---------------------------|---|--|
| predict | ✓ | ✓ Half inputs are cast to single precision and computations are performed in single precision. | ✓ Half inputs are cast to single precision and computations are performed in single precision. To perform computations in half, set the library target to 'tensorrt' and set the data type to 'FP16' in <code>coder.DeepLearningConfig</code> . |
| predictAndUpdateState | ✓ | ✓ Half inputs are cast to single precision and computations are performed in single precision. | ✓ Half inputs are cast to single precision and computations are performed in single precision. To perform computations in half, set the library target to 'tensorrt' and set the data type to 'FP16' in <code>coder.DeepLearningConfig</code> . |

| Function | MATLAB Simulation Support | C/C++ Code Generation Support | GPU Code Generation Support |
|----------|--|-------------------------------|-----------------------------|
| prod | <p>✓</p> <p>Half inputs are cast to single precision and computations are performed in single precision. As a result, saturation behavior differs between single and half inputs:</p> <pre>maxhalf = half.realmx; isequal(prod([maxhalf 2 0.5]), maxhalf) ans = logical 1 maxsingle = realmx('single'); isequal(prod([maxsingle 2 0.5]), maxsingle) ans = logical 0</pre> | ✓ | ✓ |
| rdivide | ✓ | ✓ | ✓ |
| real | ✓ | ✓ | ✓ |
| realmax | <p>✓</p> <p>Supported syntax:</p> <pre>realmax('half') realmax('like',half(1))</pre> | | |
| realmin | <p>✓</p> <p>Supported syntax:</p> <pre>realmin('half') realmin('like',half(1))</pre> | | |
| rem | ✓ | ✓ | ✓ |
| repelem | ✓ | ✓ | ✓ |

| Function | MATLAB Simulation Support | C/C++ Code Generation Support | GPU Code Generation Support |
|---------------|--|---|---|
| repmat | ✓ | ✓ Dimension argument cannot be half precision. | ✓ Dimension argument cannot be half precision. |
| reshape | ✓ | ✓ Dimension argument cannot be half precision. | ✓ Dimension argument cannot be half precision. |
| rgbplot | ✓ | | |
| round | ✓ Only one input supported | ✓ Only one input supported | ✓ Only one input supported |
| rsqrt | ✓ Complex half-precision inputs are not supported | | |
| scatter | ✓ | | |
| scatter3 | ✓ | | |
| sign | ✓ | ✓ | ✓ |
| sin | ✓ | ✓ | ✓ |
| single | ✓ | ✓ | ✓ |
| sinh | ✓ | ✓ | ✓ |
| sinpi | ✓ | ✓ | ✓ |
| size | ✓ | ✓ | ✓ |
| sort | ✓ | | |
| sqrt | ✓ | ✓ | ✓ |
| squeeze | ✓ | ✓ | ✓ |
| storedInteger | ✓ | | |

| Function | MATLAB Simulation Support | C/C++ Code Generation Support | GPU Code Generation Support |
|-----------|--|-------------------------------|-----------------------------|
| sum | <p>✓</p> <p>Half inputs are cast to single precision and computations are performed in single precision. As a result, saturation behavior differs between single and half inputs:</p> <pre>maxhalfint = half.flintmax; isequal(sum([maxhalfint, 1, -1]), maxhalfint) ans = logical 1 maxsingleint = flintmax('single'); isequal(sum([maxsingleint, 1, -1]), maxsingleint) ans = logical 0</pre> | ✓ | ✓ |
| tan | ✓ | ✓ | ✓ |
| tanh | ✓ | ✓ | ✓ |
| times, .* | ✓ | ✓ | ✓ |
| transpose | ✓ | ✓ | ✓ |
| typecast | ✓ | | |
| uint16 | ✓ | ✓ | ✓ |
| uint32 | ✓ | ✓ | ✓ |
| uint64 | ✓ | ✓ | ✓ |
| uint8 | ✓ | ✓ | ✓ |
| uminus | ✓ | ✓ | ✓ |
| uplus | ✓ | ✓ | ✓ |
| vertcat | ✓ | ✓ | ✓ |
| xlim | ✓ | | |
| ylim | ✓ | | |

| Function | MATLAB Simulation Support | C/C++ Code Generation Support | GPU Code Generation Support |
|----------|---|---|---|
| zeros | ✓ Supported syntax: zeros(__, 'half') zeros(__, 'like', p) | ✓ Supported syntax: zeros(__, 'half') zeros(__, 'like', p) | ✓ Supported syntax: zeros(__, 'half') zeros(__, 'like', p) |
| zlim | ✓ | | |

See Also

half

More About

- “Floating-Point Numbers” on page 35-20
- “What is Half Precision?” on page 33-19
- “Generate Code for Sobel Edge Detection That Uses Half-Precision Data Type” (MATLAB Coder)
- Edge Detection with Sobel Method in Half-Precision (GPU Coder)

Generate Native Half-Precision C Code Using MATLAB Coder

Some embedded hardware targets natively support special types for half precision, such as `_Float16` and `_fp16` data types for ARM[®] compilers. You can use MATLAB Coder to generate native half-precision C code for ARM Cortex[®]-A processors that natively support half precision floating-point data types.

The process to generate native half C code is as follows:

- Register a new hardware target device that natively supports half precision using the `target` package.
- Configure code generation configuration for half precision.
- Generate native half type code.

Fixed-Point Designer and MATLAB Coder include preconfigured language implementations for Armclang and GCC compilers. For other hardware targets, you can specify a custom language implementation based on your hardware specifications.

Generate Native Half-Precision C Code for ARM[®] Cortex[®]-A with GCC Compiler

In this example, an ARM Cortex[®]-A processor is used as the hardware target. The model is configured to use this ARM target and the GNU GCC compiler toolchain.

Register Target Hardware

Use the `target.create` function to create an ARM processor target that is compatible with half precision.

```
arm_half = target.create('Processor', 'Manufacturer', "Broadcom", 'Name', 'BCM2711 ARM Cortex A72');
```

Add the language implementation. Use the `target.get` function to retrieve the target object from the internal database.

```
li = target.get('LanguageImplementation', "GNU GCC ARM 32-bit");
```

Replace the default language implementation for ARM Cortex with Armclang.

```
arm_half.LanguageImplementations = li;
```

Use the `target.add` function to add the target object to the internal database.

```
target.add(arm_half);
```

```
target.add summary:
```

```
Objects added to internal database for current MATLAB session:
target.Processor          "Broadcom-BCM2711 ARM Cortex A72"
```

```
1 object not added because they already exist.
```

Open MATLAB Code that Uses Half-Precision Data Type

```
edit testNativeHalffp16.m
```

Configure Code Generation Configuration for Half Precision

Create the code generation configuration object.

```
cfg = coder.config('lib');
```

Specify the production hardware device type.

```
cfg.HardwareImplementation.ProdHWDeviceType = 'Broadcom->BCM2711 ARM Cortex A72';
```

Select the toolchain compatible with the selected hardware.

```
cfg.Toolchain = 'GNU Tools for ARM Embedded Processors';
```

Add the half-precision flags for compilation.

```
cfg.BuildConfiguration = 'Specify';  
cfg.CustomToolchainOptions{4} = '-c -MMD -MP -MF"$(@:%.o=%.dep)" -MT"$@" -O0 -mfp16-format=ieee';
```

Generate Code

```
codegen testNativeHalffp16 -args {half(3)} -launchreport -config cfg
```

You can inspect the code generation report to confirm that the custom half-precision type definitions are used.

The screenshot displays the MATLAB Coder interface. On the left, the 'MATLAB Source' pane shows a function list for 'testNativeHalffp16.m', including 'testNativeHalffp16' and 'ProcessHalf'. Below it, the 'Generated Code' pane shows a tree of source files, with 'rtwhalf.h' selected. The main editor window shows the C code for 'rtwhalf.h'. The code includes data type conversion functions (lines 101-106), math functions (lines 108-117), and a conditional compilation block (lines 119-120) for _Float16. A red box highlights the #else and typedef _Float16 real16_T; lines. The interface also shows a summary of successful code generation on the right, indicating that the code was generated on 22-Jun-2021 13:30:29 as a static library using GNU Tools for ARM Embedded Processors.

```

100
101 /* Data Type Conversion */
102 float halfToFloat(real16_T a);
103 double halfToDouble(real16_T a);
104
105 real16_T floatToHalf(float a);
106 real16_T doubleToHalf(double a);
107
108 /* Math functions */
109 real16_T sin_half(real16_T a);
110 real16_T cos_half(real16_T a);
111 real16_T ceil_half(real16_T a);
112 real16_T fix_half(real16_T a);
113 real16_T floor_half(real16_T a);
114 real16_T exp_half(real16_T a);
115 real16_T log_half(real16_T a);
116 real16_T log10_half(real16_T a);
117 real16_T sqrt_half(real16_T a);
118
119 #else
120 typedef _Float16 real16_T;
121
122 #endif
123 #endif
124 /*
125 * File trailer for rtwhalf.h
126 *
127 * [EOF]
128 */
129

```

Summary All Messages (0)

Code generation successful

Generated on: 22-Jun-2021 13:30:29
Build type: Static Library
Toolchain: GNU Tools for ARM Embedded Processors
Build Configuration: Specify

The half-precision constants use the f16 suffix.

The screenshot displays the MATLAB IDE interface during the generation of native half-precision C code. On the left, the 'MATLAB Source' pane shows the project structure, including the function `testNativeHalfp16.m` and its associated files. The 'Generated Code' pane lists the source files for the native code, such as `rtwhalf.c`, `rtwhalf.h`, and `testNativeHalfp16.c`.

The main editor shows the C code for `testNativeHalfp16.c`. A red box highlights the following code snippet:

```
57 static const real16_T hv[8] = {101.0f16, 103.0f16, 105.0f16, 107.0f16,
58                               102.0f16, 104.0f16, 106.0f16, 108.0f16};
```

The bottom pane shows a 'Summary' of the code generation process, indicating that the code generation was successful. The summary includes the following details:

- Generated on: 22-Jun-2021 13:30:29
- Build type: Static Library
- Toolchain: GNU Tools for ARM Embedded Processors
- Build Configuration: Specify

Generate Native Half-Precision C Code for ARM Cortex-A with Armclang Compiler

In this example, an ARM Cortex-A processor is used as the hardware target. The model is configured to use this ARM target and the Armclang compiler toolchain.

Register Target Hardware

Use the `target.create` function to create an ARM processor target that is compatible with half precision.

```
arm_half = target.create('Processor',...
    'Manufacturer','Broadcom',...
    'Name','ARM Cortex A75');
```

Add the language implementation. Use the `target.get` function to retrieve the target object from the internal database.

```
li = target.get('LanguageImplementation','Clang ARM 32-bit');
```

Replace the default language implementation for ARM Cortex with Armclang.

```
arm_half.LanguageImplementations = li;
```

Use the `target.add` function to add the target object to the internal database.

```
target.add(arm_half);
```

Configure Code Generation Configuration for Half Precision

Create the code generation configuration object.

```
cfg = coder.config('lib');
```

Specify the production hardware type.

```
cfg.HardwareImplementation.ProdHWDeviceType = 'Broadcom->ARM Cortex A75';
```

Select the toolchain compatible with the selected hardware.

```
cfg.Toolchain = 'Armclang Compiler';
```

Add the half-precision flags for compilation.

```
cfg.BuildConfiguration = 'Specify';
cfg.CustomToolchainOptions{4} = '-c -MMD -MP -MF"$(@:%.o=%%.dep)" -MT"$@" -O0 --target=arm-arm-non';
```

Generate Code

```
codegen testNativeHalffp16 -args {half(3)} -launchreport -config cfg
```

Register ARM Target Hardware with Custom Language Implementation

In this example, create a new custom language implementation with half precision for a compatible ARM target.

Register Target Hardware

Use the `target.create` function to copy the ARM Compatible-ARM Cortex language implementation.

```
languageImplementation = target.create('LanguageImplementation',...
    'Name','ARM with half',...
    'Copy','ARM Compatible-ARM Cortex');
```

Specify custom half information and target specific headers, as given by your target hardware documentation. For more information, see “Register New Hardware Devices”. For example,

```
customHalf = target.create('FloatingPointDataType',...
    'Name','BCM2711 Half Type', ...
    'TypeName','_Float16',...
    'LiteralSuffix','f16',...
    'Size',16, ...
    'SystemIncludes',["arm_fp16.h" "arm_neon.h"]);
languageImplementation.DataTypes.NonStandardDataTypes = customHalf;
```

Provide information about your target processor. For example,

```
% Broadcom BCM2711
% Quad core Cortex-A72 (ARM v8) 64-bit SoC
pi4a72 = target.create('Processor','Manufacturer',...
    'Broadcom','Name','BCM2711');
```

Add the custom half-precision language implementation.

```
pi4a72.LanguageImplementations = languageImplementation;
```

Use the `target.add` function to add the target object to the internal database.

```
target.add(pi4a72);
```

See Also

`half` | `target.FloatingPointDataType` | `target.add` | `target.create` | `target.get` | `target.remove`

Related Examples

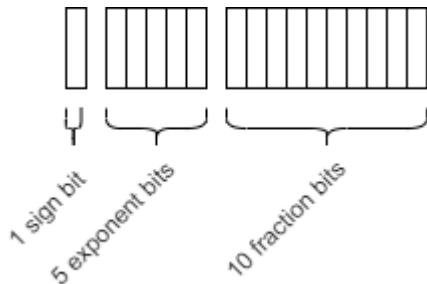
- “Half Precision Code Generation Support” on page 33-2
- “Register New Hardware Devices”

External Websites

- Clang Language Extensions for Half-Precision Floating Point
- Arm Compiler armclang Reference Guide: Half-precision floating-point data types
- GCC Half-Precision Floating Point
- Reduce the Program Data Size with Ease! Introducing Half-Precision Floating-Point Feature in Renesas Compiler Professional Edition

What is Half Precision?

The IEEE® 754 half-precision floating-point format is a 16-bit word divided into a 1-bit sign indicator s , a 5-bit biased exponent e , and a 10-bit fraction f .



Because numbers of type `half` are stored using 16 bits, they require less memory than numbers of type `single`, which uses 32 bits, or `double`, which uses 64 bits. However, because they are stored with fewer bits, numbers of type `half` are represented to less precision than numbers of type `single` or `double`.

The range, bias, and precision for supported floating-point data types are given in the table below.

| Data Type | Low Limit | High Limit | Exponent Bias | Precision |
|-----------|---------------------------------------|---|---------------|----------------------------|
| Half | $2^{-14} \approx 6.1 \cdot 10^{-5}$ | $(2-2^{-10}) \cdot 2^{15} \approx 6.5 \cdot 10^4$ | 15 | $2^{-10} \approx 10^{-3}$ |
| Single | $2^{-126} \approx 10^{-38}$ | $2^{128} \approx 3 \cdot 10^{38}$ | 127 | $2^{-23} \approx 10^{-7}$ |
| Double | $2^{-1022} \approx 2 \cdot 10^{-308}$ | $2^{1024} \approx 2 \cdot 10^{308}$ | 1023 | $2^{-52} \approx 10^{-16}$ |

For a video introduction to the half-precision data type, see [What Is Half Precision?](#) and [Half-Precision Math in Modeling and Code Generation](#).

Half Precision Applications

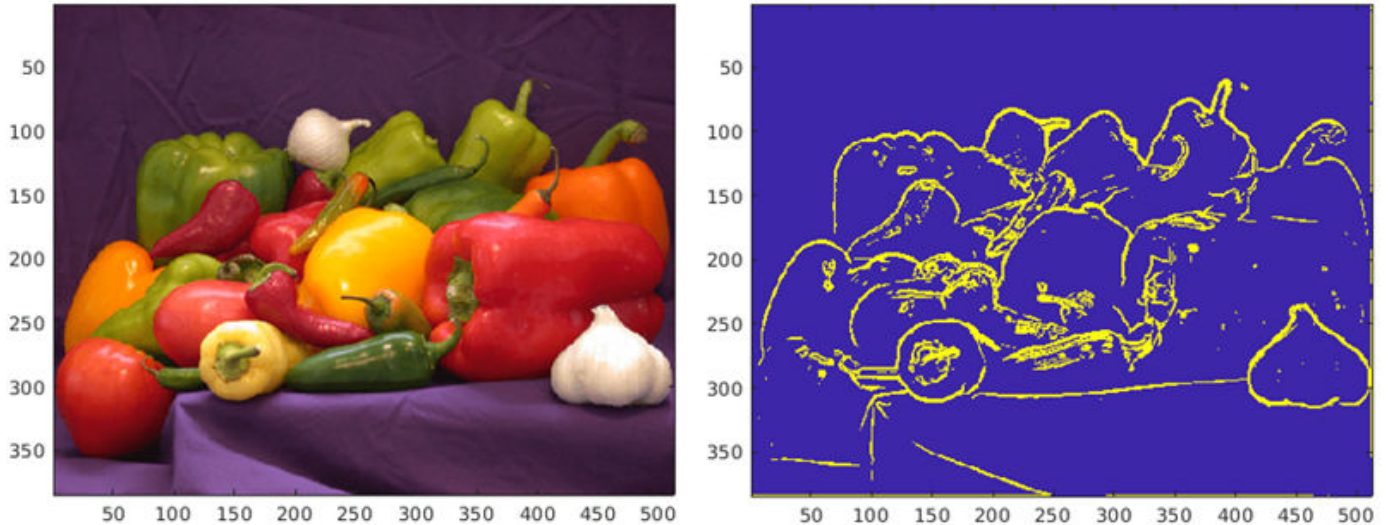
When an algorithm contains large or unknown dynamic ranges (for example integrators in feedback loops) or when the algorithm uses operations that are difficult to design in fixed-point (for example `atan2`), it can be advantageous to use floating-point representations. The half-precision data type occupies only 16 bits of memory, but its floating-point representation enables it to handle wider dynamic ranges than integer or fixed-point data types of the same size. This makes half precision particularly suitable for some image processing and graphics applications. When half-precision is used with deep neural networks, the time needed for training and inference can be reduced. By using half precision as a storage time for lookup tables, the memory footprint of the lookup table can be reduced.

MATLAB Examples

- “Fog Rectification” (GPU Coder) — The fog rectification image processing algorithm uses convolution, image color space conversion, and histogram-based contrast stretching to enhance the input image. This example shows how to generate and execute CUDA MEX with half-precision data types for these image processing operations.



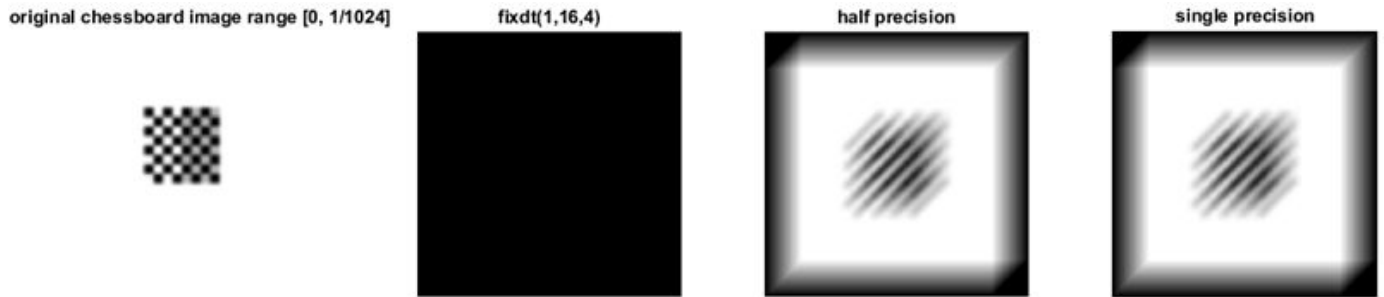
- “Edge Detection with Sobel Method in Half-Precision” (GPU Coder) — The sobel edge detection algorithm takes an input image and returns an output image that emphasizes high spatial frequency regions that correspond to edges in the input image. This example shows how to generate and execute CUDA MEX with the half-precision data type used for the input image and Sobel operator kernel values.



- “Generate Code for Sobel Edge Detection That Uses Half-Precision Data Type” (MATLAB Coder) — This example shows how to generate a standalone C++ library from a MATLAB function that performs Sobel edge detection of images by using half-precision floating point numbers.

Simulink Examples

- “Half-Precision Field-Oriented Control Algorithm” on page 51-11 — This example implements a Field-Oriented Control (FOC) algorithm using both single precision and half precision.
- “Image Quantization with Half-Precision Data Types” on page 51-14 — This example shows the effects of quantization on images. While the fixed-point data type does not always produce an acceptable results, the half-precision data type, which uses the same number of bits as the fixed-point data type, produces a result comparable to the single-precision result.



- “Digit Classification with Half-Precision Data Types” on page 51-20 — This example compares the results of a trained neural network classification model in double precision and half precision.
- “Convert Single Precision Lookup Table to Half Precision” on page 51-15 — This example demonstrates how to convert a single-precision lookup table to use half precision. Half precision is the storage type; the lookup table computations are performed using single precision. After converting to half precision, the memory size of the Lookup Table blocks are reduced by half while maintaining the desired system performance.

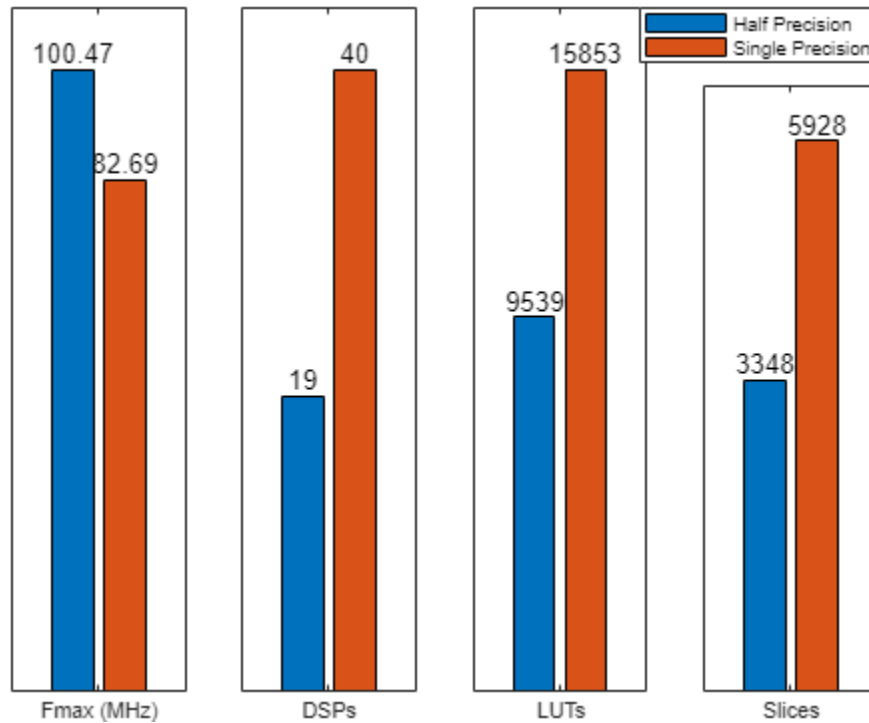
Benefits of Using Half Precision in Embedded Applications

The half precision data type uses less memory than other floating-point types like single and double. Though it occupies only 16 bits of memory, its floating-point representations enables it to handle wider dynamic ranges than integer or fixed-point data types of the same size.

FPGA

The half precision data type uses significantly less area and has low latency compared to the single precision data type when used on hardware. Half precision is particularly advantageous for low dynamic range applications.

The following plot shows the advantage of using half precision for an implementation of a field-oriented control algorithm in Xilinx® Virtex® 7 hardware.



GPU

In GPUs that support the half-precision data type, arithmetic operations are faster as compared to single or double precision.

In applications like deep learning, which require a large number of computations, using half precision can provide significant performance benefits without significant loss of precision. With GPU Coder, you can generate optimized code for prediction of a variety of trained deep learning networks from the Deep Learning Toolbox™. You can configure the code generator to take advantage of the NVIDIA TensorRT high performance inference library for NVIDIA GPUs. TensorRT provides improved latency, throughput, and memory efficiency by combining network layers and optimizing kernel selection. You can also configure the code generator to take advantage TensorRT's precision modes (FP32, FP16, or INT8) to further improve performance and reduce memory requirements.

CPU

In CPUs that support the half-precision data type, arithmetic operations are faster as compared to single or double precision. For ARM targets that natively support half-precision data types, you can generate native half C code from MATLAB or Simulink. See “Code Generation with Half Precision” on page 33-23.

Half Precision in MATLAB

Many functions in MATLAB support the half-precision data type. For a full list of supported functions, see `half`.

Half Precision in Simulink

Signals and block outputs in Simulink can specify a half-precision data type. The half-precision data type is supported for simulation and code generation for parameters and a subset of blocks. To view the blocks that support half precision, at the command line, type:

```
showblockdatatypetable
```

Blocks that support half precision display an X in the column labeled **Half**. For detailed information about half precision support in Simulink, see “The Half-Precision Data Type in Simulink” on page 51-2.

Code Generation with Half Precision

The half precision data type is supported for C/C++ code generation, CUDA code generation using GPU Coder, and HDL code generation using HDL Coder. For GPU targets, the half-precision data type uses the native half data type available in NVIDIA GPU for maximum performance.

For detailed code generation support for half precision in MATLAB and Simulink, see “Half Precision Code Generation Support” on page 33-2 and “The Half-Precision Data Type in Simulink” on page 51-2.

For embedded hardware targets that natively support special types for half precision, such as `_Float16` and `_fp16` data types for ARM compilers, you can generate native half precision C code using Embedded Coder or MATLAB Coder. For more information, see “Generate Native Half-Precision C Code from Simulink Models” on page 51-5 and “Generate Native Half-Precision C Code Using MATLAB Coder” on page 33-13.

See Also

`half` | “The Half-Precision Data Type in Simulink” on page 51-2 | “Half Precision” 16-bit Floating Point Arithmetic | “Floating-Point Numbers” on page 35-20

Related Examples

- “Fog Rectification” (GPU Coder)
- “Edge Detection with Sobel Method in Half-Precision” (GPU Coder)
- “Generate Code for Sobel Edge Detection That Uses Half-Precision Data Type” (MATLAB Coder)
- “Half-Precision Field-Oriented Control Algorithm” on page 51-11
- “Image Quantization with Half-Precision Data Types” on page 51-14
- “Digit Classification with Half-Precision Data Types” on page 51-20
- “Convert Single Precision Lookup Table to Half Precision” on page 51-15

Fixed-Point Designer for Simulink Models

Getting Started

- “Sharing Fixed-Point Models” on page 34-2
- “Physical Quantities and Measurement Scales” on page 34-3
- “Why Use Fixed-Point Hardware?” on page 34-9
- “Why Use the Fixed-Point Designer Software?” on page 34-10
- “Developing and Testing Fixed-Point Systems” on page 34-11
- “Supported Data Types” on page 34-13
- “Configure Blocks with Fixed-Point Output” on page 34-14
- “Configure Blocks with Fixed-Point Parameters” on page 34-20
- “Pass Fixed-Point Data Between Simulink Models and MATLAB” on page 34-22
- “Cast from Doubles to Fixed Point” on page 34-25

Sharing Fixed-Point Models

You can edit a model containing fixed-point blocks without the Fixed-Point Designer software. However, you must have a Fixed-Point Designer software license to:

- Update a Simulink diagram (**Ctrl+D**) containing fixed-point data types
- Simulate a model containing fixed-point data types
- Generate code from a model containing fixed-point data types
- Log the minimum and maximum values produced by a simulation
- Automatically scale the output of a model

If the Fixed-Point Designer product is not installed on your system, you can work with a model containing Simulink blocks with fixed-point settings as follows:

- 1 Instrumentation requires a Fixed-Point Designer license. To disable fixed-point instrumentation on a model, set the `MinMaxOverflowLogging` parameter to `ForceOff`. At the command line, enter:

```
set_param(gcs, 'MinMaxOverflowLogging', 'ForceOff')
```

- 2 If you do not have Fixed-Point Designer software, you can still configure data type override settings to simulate a model that specifies fixed-point data types. Using this setting, the software temporarily overrides data types with floating-point data types during simulation. To simulate a model without using Fixed-Point Designer, at the command line enter:

```
set_param(gcs, 'DataTypeOverride', 'Double', ...  
'DataTypeOverrideAppliesTo', 'AllNumericTypes')
```

- 3 If you use `fi` objects or embedded numeric data types in your model or workspace, you might introduce fixed-point data types into your model. To prevent the checkout of a Fixed-Point Designer license, set the `fipref` `DataTypeOverride` property to `TrueDoubles` and the `DataTypeOverrideAppliesTo` property to `AllNumericTypes`.

At the MATLAB command line, enter:

```
p = fipref('DataTypeOverride', 'TrueDoubles', ...  
'DataTypeOverrideAppliesTo', 'AllNumericTypes');
```

See Also

`fipref` | “Fixed-Point Instrumentation and Data Type Override” on page 42-61

Physical Quantities and Measurement Scales

| In this section... |
|---|
| “Introduction” on page 34-3 |
| “Selecting a Measurement Scale” on page 34-3 |
| “Select a Measurement Scale for Temperature” on page 34-5 |

Introduction

The decision to use fixed-point hardware is simply a choice to represent numbers in a particular form. This representation often offers advantages in terms of the power consumption, size, memory usage, speed, and cost of the final product.

A measurement of a physical quantity can take many numerical forms. For example, the boiling point of water is 100 degrees Celsius, 212 degrees Fahrenheit, 373 kelvin, or 671.4 degrees Rankine. No matter what number is given, the physical quantity is exactly the same. The numbers are different because four different scales are used.

Well known standard scales like Celsius are convenient for the exchange of information. However, there are situations where it makes sense to create and use unique nonstandard scales. These situations usually involve making the most of a limited resource.

For example, nonstandard scales allow map makers to get the maximum detail on a fixed size sheet of paper. A typical road atlas of the USA will show each state on a two-page display. The scale of inches to miles will be unique for most states. By using a large ratio of miles to inches, all of Texas can fit on two pages. Using the same scale for Rhode Island would make poor use of the page. Using a much smaller ratio of miles to inches would allow Rhode Island to be shown with the maximum possible detail.

Fitting measurements of a variable inside an embedded processor is similar to fitting a state map on a piece of paper. The map scale should allow all the boundaries of the state to fit on the page. Similarly, the binary scale for a measurement should allow the maximum and minimum possible values to fit. The map scale should also make the most of the paper in order to get maximum detail. Similarly, the binary scale for a measurement should make the most of the processor in order to get maximum precision.

Use of standard scales for measurements has definite compatibility advantages. However, there are times when it is worthwhile to break convention and use a unique nonstandard scale. There are also occasions when a mix of uniqueness and compatibility makes sense. See the sections that follow for more information.

Selecting a Measurement Scale

Suppose that you want to make measurements of the temperature of liquid water, and that you want to represent these measurements using 8-bit unsigned integers. Fortunately, the temperature range of liquid water is limited. No matter what scale you use, liquid water can only go from the freezing point to the boiling point. Therefore, this is the range of temperatures that you must capture using just the 256 possible 8-bit values: 0,1,2,...,255.

One approach to representing the temperatures is to use a standard scale. For example, the units for the integers could be Celsius. Hence, the integers 0 and 100 represent water at the freezing point

and at the boiling point, respectively. On the upside, this scale gives a trivial conversion from the integers to degrees Celsius. On the downside, the numbers 101-255 are unused. By using this standard scale, more than 60% of the number range has been wasted.

A second approach is to use a nonstandard scale. In this scale, the integers 0 and 255 represent water at the freezing point and at the boiling point, respectively. On the upside, this scale gives maximum precision since there are 254 values between freezing and boiling instead of just 99. On the downside, the units are roughly 0.3921568 degrees Celsius per bit so the conversion to Celsius requires division by 2.55, which is a relatively expensive operation on most fixed-point processors.

A third approach is to use a “semistandard” scale. For example, the integers 0 and 200 could represent water at the freezing point and at the boiling point, respectively. The units for this scale are 0.5 degrees Celsius per bit. On the downside, this scale does not use the numbers from 201-255, which represents a waste of more than 21%. On the upside, this scale permits relatively easy conversion to a standard scale. The conversion to Celsius involves division by 2, which is an easy shift operation on most processors.

Measurement Scales: Beyond Multiplication

One of the key operations in converting from one scale to another is multiplication. The preceding case study gave three examples of conversions from a quantized integer value Q to a real-world Celsius value V that involved only multiplication:

$$V = \begin{cases} \frac{100^{\circ}\text{C}}{100} Q_1 & \text{Conversion 1} \\ \frac{100^{\circ}\text{C}}{255} Q_2 & \text{Conversion 2} \\ \frac{100^{\circ}\text{C}}{200} Q_3 & \text{Conversion 3} \end{cases}$$

Graphically, the conversion is a line with slope S , which must pass through the origin. A line through the origin is called a purely linear conversion. Restricting yourself to a purely linear conversion can be wasteful and it is often better to use the general equation of a line:

$$V = SQ + B.$$

By adding a bias term B , you can obtain greater precision when quantizing to a limited number of bits.

The general equation of a line gives a useful conversion to a quantized scale. However, like all quantization methods, the precision is limited and errors can be introduced by the conversion. The general equation of a line with quantization error is given by

$$V = SQ + B \pm \text{Error}.$$

If the quantized value Q is rounded to the nearest representable number, then

$$-\frac{S}{2} \leq \text{Error} \leq \frac{S}{2}.$$

That is, the amount of quantization error is determined by both the number of bits and by the scale. This scenario represents the best-case error. For other rounding schemes, the error can be twice as large.

Select a Measurement Scale for Temperature

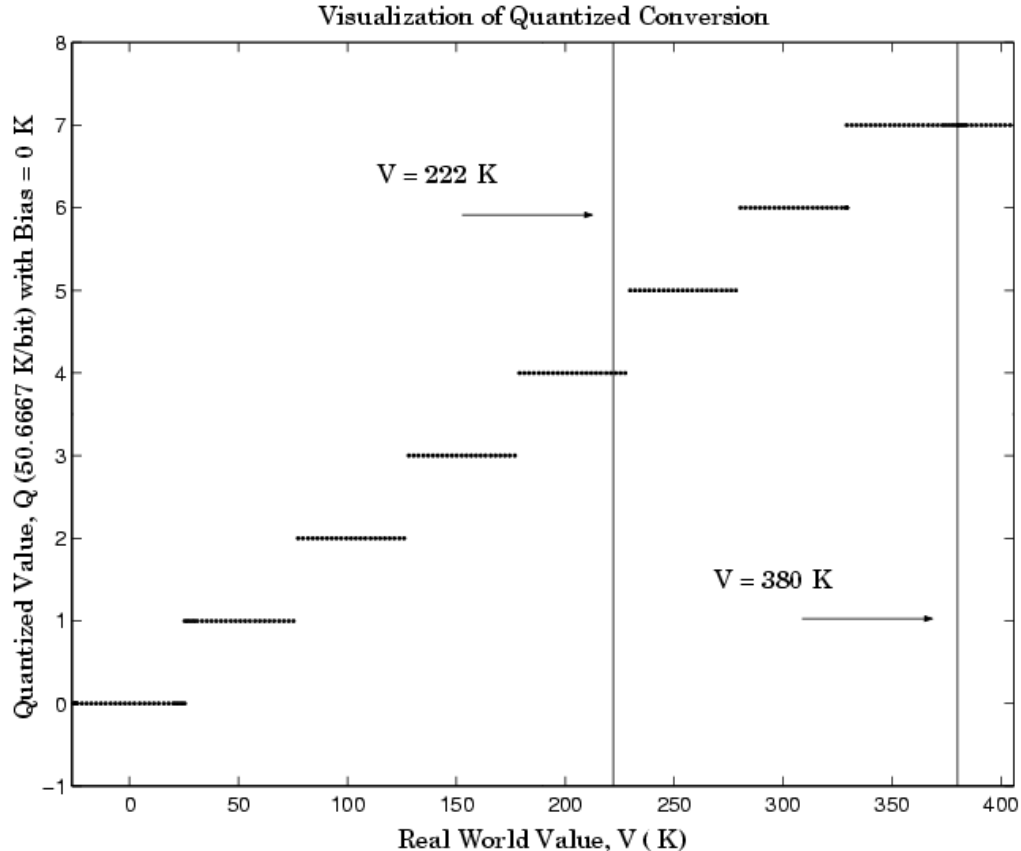
On typical electronically controlled internal combustion engines, the flow of fuel is regulated to obtain the desired ratio of air to fuel in the cylinders just prior to combustion. Therefore, knowledge of the current air flow rate is required. Some manufacturers use sensors that directly measure air flow, while other manufacturers calculate air flow from measurements of related signals. The relationship of these variables is derived from the ideal gas equation. The ideal gas equation involves division by air temperature. For proper results, an absolute temperature scale such as kelvin or Rankine must be used in the equation. However, quantization directly to an absolute temperature scale would cause needlessly large quantization errors.

The temperature of the air flowing into the engine has a limited range. On a typical engine, the radiator is designed to keep the block below the boiling point of the cooling fluid. Assume a maximum of 225°F (380 K). As the air flows through the intake manifold, it can be heated to this maximum temperature. For a cold start in an extreme climate, the temperature can be as low as -60°F (222 K). Therefore, using the absolute kelvin scale, the range of interest is 222-380 K.

The air temperature needs to be quantized for processing by the embedded control system. Assuming an unrealistic quantization to 3-bit unsigned numbers: 0,1,2,...,7, the purely linear conversion with maximum precision is

$$V = \frac{380 \text{ K}}{7.5 \text{ bit}} Q.$$

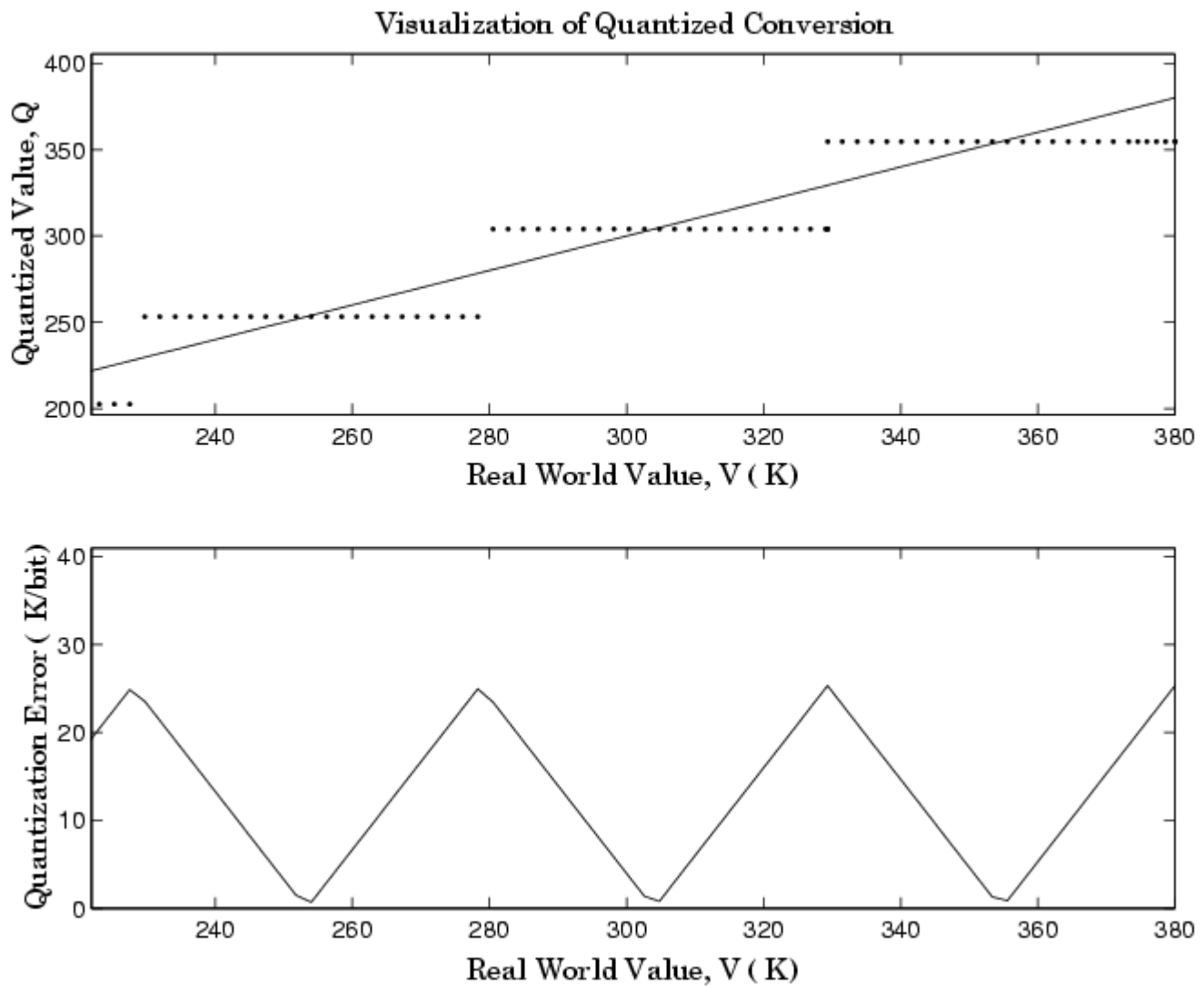
The quantized conversion and range of interest are shown in the following figure.



Notice that there are 7.5 possible quantization values. This is because only half of the first bit corresponds to temperatures (real-world values) greater than zero.

The quantization error is $-25.33 \text{ K/bit} \leq \text{Error} \leq 25.33 \text{ K/bit}$.

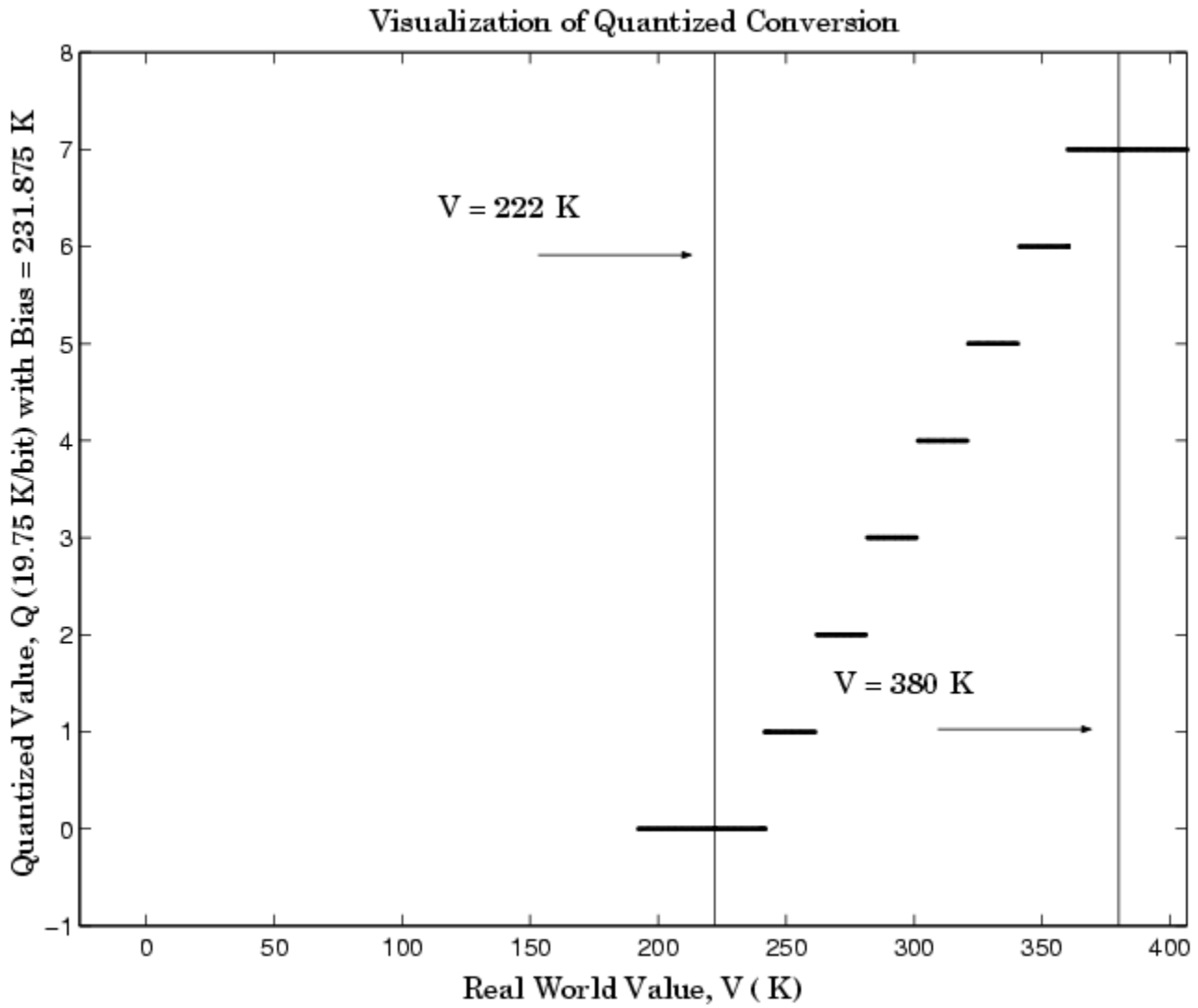
The range of interest of the quantized conversion and the absolute value of the quantized error are shown in the following figure.



As an alternative to the purely linear conversion, consider the general linear conversion with maximum precision:

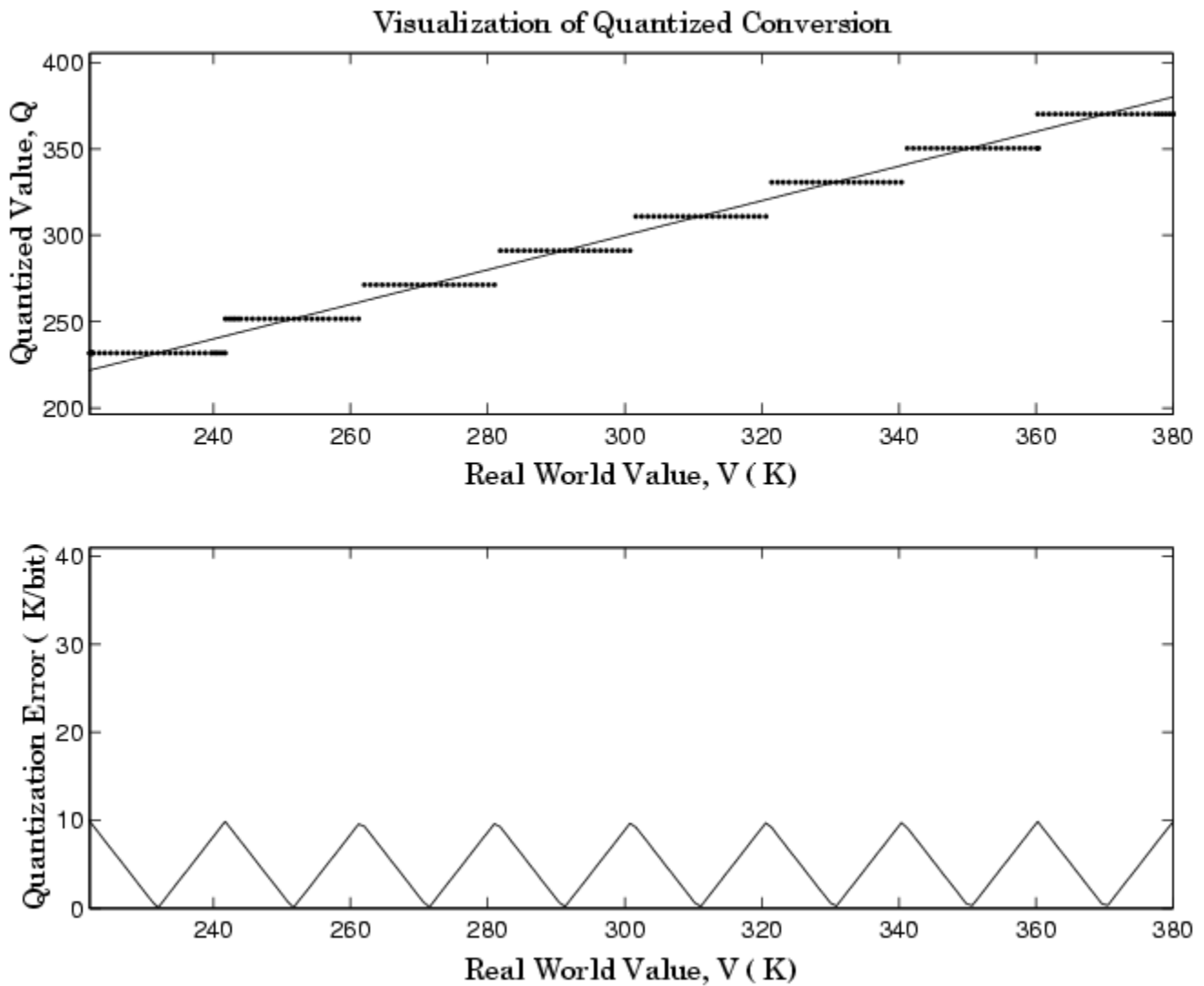
$$V = \left(\frac{380 \text{ K} - 222 \text{ K}}{8} \right) Q + 222 \text{ K} + 0.5 \left(\frac{380 \text{ K} - 222 \text{ K}}{8} \right)$$

The quantized conversion and range of interest are shown in the following figure.



The quantization error is $-9.875 \text{ K/bit} \leq \text{Error} \leq 9.875 \text{ K/bit}$, which is approximately 2.5 times smaller than the error associated with the purely linear conversion.

The range of interest of the quantized conversion and the absolute value of the quantized error are shown in the following figure.



Clearly, the general linear scale gives much better precision than the purely linear scale over the range of interest.

Why Use Fixed-Point Hardware?

Digital hardware is becoming the primary means by which control systems and signal processing filters are implemented. Digital hardware can be classified as either off-the-shelf hardware (for example, microcontrollers, microprocessors, general-purpose processors, and digital signal processors) or custom hardware. Within these two types of hardware, there are many architecture designs. These designs range from systems with a single instruction, single data stream processing unit to systems with multiple instruction, multiple data stream processing units.

Within digital hardware, numbers are represented as either fixed-point or floating-point data types. For both of these data types, word sizes are fixed at a set number of bits. However, the dynamic range of fixed-point values is much less than floating-point values with equivalent word sizes. Therefore, in order to avoid overflow or unreasonable quantization errors, fixed-point values must be scaled. Since floating-point processors can greatly simplify the real-time implementation of a control law or digital filter, and floating-point numbers can effectively approximate real-world numbers, then why use a microcontroller or processor with fixed-point hardware support?

- **Size and Power Consumption** — The logic circuits of fixed-point hardware are much less complicated than those of floating-point hardware. This means that the fixed-point chip size is smaller with less power consumption when compared with floating-point hardware. For example, consider a portable telephone where one of the product design goals is to make it as portable (small and light) as possible. If one of today's high-end, floating-point, general-purpose processors is used, a large heat sink and battery would also be needed, resulting in a costly, large, and heavy portable phone.
- **Memory Usage and Speed** — In general fixed-point calculations require less memory and less processor time to perform.
- **Cost** — Fixed-point hardware is more cost effective where price/cost is an important consideration. When digital hardware is used in a product, especially mass-produced products, fixed-point hardware costs much less than floating-point hardware and can result in significant savings.

After making the decision to use fixed-point hardware, the next step is to choose a method for implementing the dynamic system (for example, control system or digital filter). Floating-point software emulation libraries are generally ruled out because of timing or memory size constraints. Therefore, you are left with fixed-point math where binary integer values are scaled.

Why Use the Fixed-Point Designer Software?

The Fixed-Point Designer software allows you to efficiently design control systems and digital filters that you will implement using fixed-point arithmetic. With the Fixed-Point Designer software, you can construct Simulink and Stateflow models that contain detailed fixed-point information about your systems. You can then perform bit-true simulations with the models to observe the effects of limited range and precision on your designs.

You can configure the Fixed-Point Tool to automatically log the overflows, saturations, and signal extremes of your simulations. You can also use it to automate data typing and scaling decisions and to compare your fixed-point implementations against idealized, floating-point benchmarks.

You can use the Fixed-Point Designer software with the Simulink Coder product to automatically generate efficient, integer-only C code representations of your designs. You can use this C code in a production target or for rapid prototyping. In addition, you can use the Fixed-Point Designer software with the Embedded Coder product to generate real-time C code for use on an integer production, embedded target. You can also use Fixed-Point Designer with HDL Coder to generate portable, synthesizable VHDL and Verilog code from Simulink models and Stateflow charts.

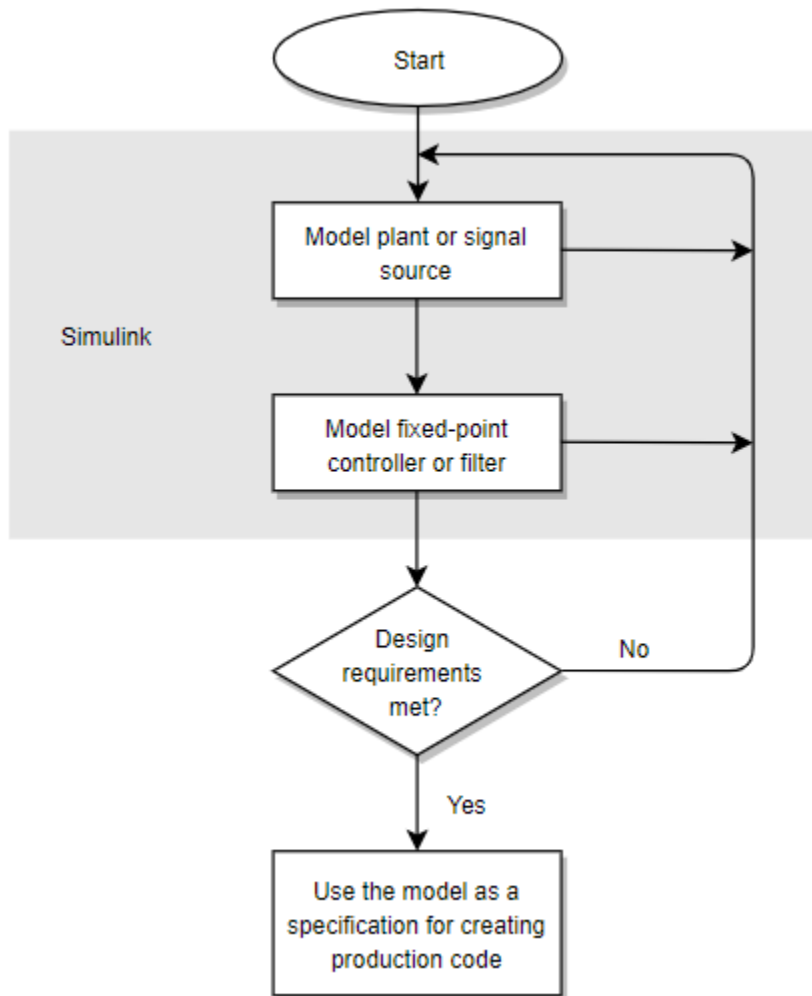
See Also

More About

- “Why Use Fixed-Point Hardware?” on page 34-9

Developing and Testing Fixed-Point Systems

The Fixed-Point Designer software provides tools that aid in the development and testing of fixed-point dynamic systems. You directly design dynamic system models in the Simulink software that are ready for implementation on fixed-point hardware. The development cycle is illustrated below.



Using the MATLAB, Simulink, and Fixed-Point Designer software, you follow these steps of the development cycle:

- 1 Model the system (plant or signal source) within the Simulink software using double-precision numbers. Typically, the model will contain nonlinear elements.
- 2 Design and simulate a fixed-point dynamic system (for example, a control system or digital filter) with fixed-point Simulink blocks that meets the design, performance, and other constraints.
- 3 Analyze the results and go back to step 1 if needed.

When you have met the design requirements, you can use the model as a specification for creating production code using the Simulink Coder product or generating HDL code using the HDL Coder product.

The above steps interact strongly. In steps 1 and 2, there is a significant amount of freedom to select different solutions. Generally, you fine-tune the model based on feedback from the results of the current implementation (step 3). There is no specific modeling approach. For example, you may obtain models from first principles such as equations of motion, or from a frequency response such as a sine sweep. There are many controllers that meet the same frequency-domain or time-domain specifications. Additionally, for each controller there are an infinite number of realizations.

The Fixed-Point Designer software helps expedite the design cycle by allowing you to simulate the effects of various fixed-point controller and digital filter structures.

See Also

More About

- “Why Use Fixed-Point Hardware?” on page 34-9

Supported Data Types

The Fixed-Point Designer software supports the following integer and fixed-point data types for simulation and code generation:

- Unsigned data types from 1 bit to 128 bits
- Signed data types from 2 bits to 128 bits
- Boolean, double, single, and half
- Scaled doubles

The `fi` object in MATLAB has a word length limit of 65535 bits.

The software supports all scaling choices including pure integer, binary point, and slope bias. For slope bias scaling, it does not support complex fixed-point types that have non-zero bias or non-trivial slope.

The save data type support extends to signals, parameters, and states.

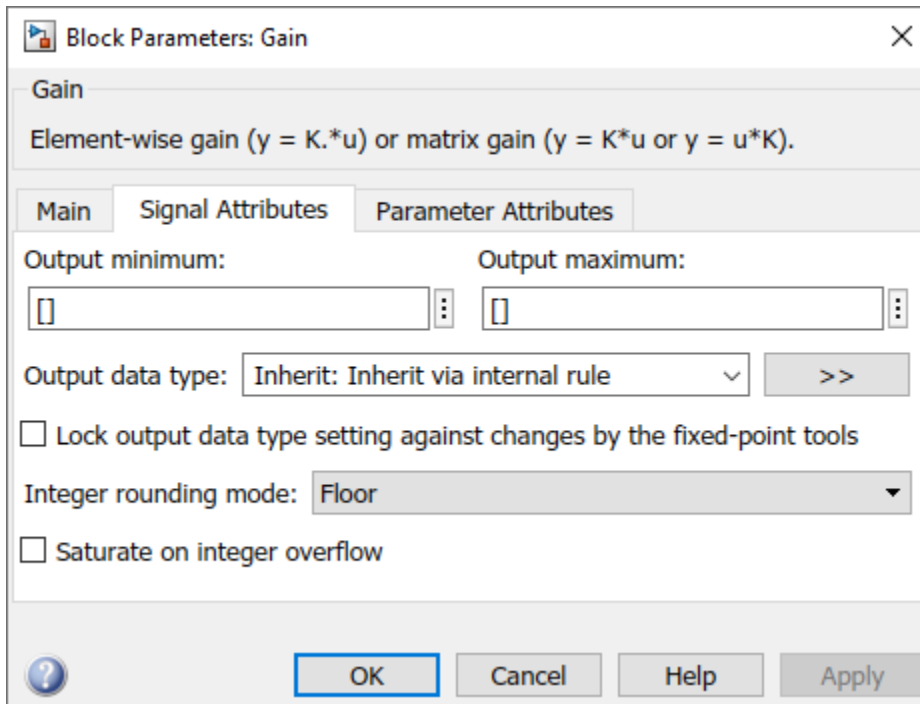
See Also

More About

- “Specify Data Types Using Data Type Assistant”

Configure Blocks with Fixed-Point Output

To create a fixed-point model, configure Simulink blocks to output fixed-point signals. Simulink blocks that support fixed-point output provide parameters that allow you to specify whether a block should output fixed-point signals and, if so, the size, scaling, and other attributes of the fixed-point output. These parameters typically appear on the **Signal Attributes** pane of the block's parameter dialog box.



The following sections explain how to use these parameters to configure a block for fixed-point output.

Specify the Output Data Type and Scaling

Many Simulink blocks allow you to specify an output data type and scaling using a parameter that appears on the block dialog box. This parameter (typically named **Output data type**) provides a pull-down menu that lists the data types a particular block supports. In general, you can specify the output data type as a rule that inherits a data type, a built-in data type, an expression that evaluates to a data type, or a Simulink data type object. For more information, see “Control Data Types of Signals”.

The Fixed-Point Designer software enables you to configure Simulink blocks with:

- **Fixed-point data types**

Fixed-point data types are characterized by their word size in bits and by their binary point — the means by which fixed-point values are scaled.

- **Floating-point data types**

Floating-point data types are characterized by their sign bit, fraction (mantissa) field, and exponent field.

To configure blocks with Fixed-Point Designer data types, specify the data type parameter on a block dialog box as an expression that evaluates to a data type. Alternatively, you can use an assistant that simplifies the task of entering data type expressions (see “Specify Fixed-Point Data Types with the Data Type Assistant” on page 34-15). The sections that follow describe varieties of fixed-point and floating-point data types and example of how to specify these types by using the `fixdt` function. The `fixdt` function also allows you to specify scaling for fixed-point data types.

Integers

To configure a 16-bit unsigned integer via the block dialog box, specify the **Output data type** parameter as `fixdt(0,16,0)`. To configure a 16-bit signed integer, specify the **Output data type** parameter as `fixdt(1,16,0)`.

For integer data types, the default binary point is assumed to lie to the right of all bits.

Fractional Numbers

To configure the output as a 16-bit unsigned fractional number via the block dialog box, specify the **Output data type** parameter to be `fixdt(0,16,16)`. To configure a 16-bit signed fractional number, specify **Output data type** to be `fixdt(1,16,15)`.

Fractional numbers are distinguished from integers by their default scaling. Whereas signed and unsigned integer data types have a default binary point to the right of all bits, unsigned fractional data types have a default binary point to the left of all bits, while signed fractional data types have a default binary point to the right of the sign bit.

Both unsigned and signed fractional data types support *guard bits*, which act to guard against overflow. For example, `fixdt(1,16,11)` specifies a 16-bit signed fractional number with 4 guard bits. The guard bits lie to the left of the default binary point.

Generalized Fixed-Point Numbers

To configure the output as a 16-bit unsigned generalized fixed-point number via the block dialog box, specify the **Output data type** parameter to be `fixdt(0,16)`. To configure a 16-bit signed generalized fixed-point number, specify **Output data type** to be `fixdt(1,16)`.

Generalized fixed-point numbers are distinguished from integers and fractionals by the absence of a default scaling. For these data types, a block typically inherits its scaling from another block.

Floating-Point Numbers

The Fixed-Point Designer software supports single-precision and double-precision floating-point numbers as defined by the IEEE Standard 754-1985 for Binary Floating-Point Arithmetic.

To configure the output as a single-precision floating-point number via the block dialog box, specify the **Output data type** parameter as `fixdt('single')`. To configure a double-precision floating-point number, specify **Output data type** as `fixdt('double')`.

Specify Fixed-Point Data Types with the Data Type Assistant

The **Data Type Assistant** is an interactive graphical tool that simplifies the task of specifying data types for Simulink blocks and data objects. The assistant appears on block and object dialog boxes,

adjacent to parameters that provide data type control, such as the **Output data type** parameter. For more information about accessing and interacting with the assistant, see “Specify Data Types Using Data Type Assistant”.

You can use the **Data Type Assistant** to specify a fixed-point data type. When you select **Fixed point** in the **Mode** field, the assistant displays fields for describing additional attributes of a fixed-point data type, as shown in this example:

The screenshot shows the 'Block Parameters: Constant' dialog box with the 'Signal Attributes' tab selected. The 'Output data type' is set to 'fixdt(1,16,2^0,0)'. The 'Data Type Assistant' section is expanded, showing the following settings:

- Mode: Fixed point
- Signedness: Signed
- Word length: 16
- Scaling: Slope and bias
- Slope: 2^0
- Bias: 0
- Data type override: Inherit
- Buttons: Calculate Best-Precision Scaling

There is also a checkbox for 'Lock output data type setting against changes by the fixed-point tools' which is currently unchecked. At the bottom of the dialog are buttons for '?', 'OK', 'Cancel', 'Help', and 'Apply'.

You can set the following fixed-point attributes:

Signedness

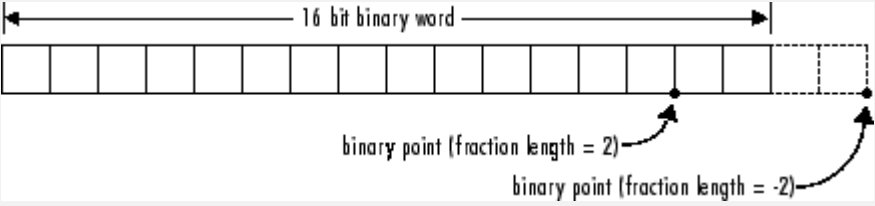
Select whether you want the fixed-point data to be **Signed** or **Unsigned**. Signed data can represent positive and negative quantities. Unsigned data represents positive values only.

Word Length

Specify the size (in bits) of the word that will hold the quantized integer. Large word sizes represent large quantities with greater precision than small word sizes. Fixed-point word sizes up to 128 bits are supported for simulation.

Scaling

Specify the method for scaling your fixed-point data to avoid overflow conditions and minimize quantization errors. You can select the following scaling modes:

| Scaling Mode | Description |
|----------------|---|
| Binary point | <p>If you select this mode, the assistant displays the Fraction length field, specifying the binary point location.</p> <p>Binary points can be positive or negative integers. A positive integer moves the binary point left of the rightmost bit by that amount. For example, an entry of 2 sets the binary point in front of the second bit from the right. A negative integer moves the binary point further right of the rightmost bit by that amount.</p>  <p>See “Binary-Point-Only Scaling” on page 35-5 for more information.</p> |
| Slope and bias | <p>If you select this mode, the assistant displays fields for entering the Slope and Bias.</p> <ul style="list-style-type: none"> • Slope can be any <i>positive</i> real number. • Bias can be any real number. <p>See “Slope and Bias Scaling” on page 35-6 for more information.</p> |
| Best precision | <p>If you select this mode, the block scales a constant vector or matrix such that the precision of its elements is maximized. This mode is available only for particular blocks.</p> <p>See “Constant Scaling for Best Precision” on page 36-19 for more information.</p> |

Calculate Best-Precision Scaling

The Fixed-Point Designer software can automatically calculate “best-precision” values for both Binary point and Slope and bias scaling, based on the values that you specify for other parameters on the dialog box. To calculate best-precision-scaling values automatically, enter values for the block’s **Output minimum** and **Output maximum** parameters. Then click the **Calculate Best-Precision Scaling** button in the assistant.

Rounding

You specify how fixed-point numbers are rounded with the **Integer rounding mode** parameter. The following rounding modes are supported:

- **Ceiling** — This mode rounds toward positive infinity and is equivalent to the MATLAB `ceil` function.
- **Convergent** — This mode rounds toward the nearest representable number, with ties rounding to the nearest even integer. Convergent rounding is equivalent to the Fixed-Point Designer `convergent` function.

- **Floor** — This mode rounds toward negative infinity and is equivalent to the MATLAB `floor` function.
- **Nearest** — This mode rounds toward the nearest representable number, with the exact midpoint rounded toward positive infinity. Rounding toward nearest is equivalent to the Fixed-Point Designer `nearest` function.
- **Round** — This mode rounds to the nearest representable number, with ties for positive numbers rounding in the direction of positive infinity and ties for negative numbers rounding in the direction of negative infinity. This mode is equivalent to the Fixed-Point Designer `round` function.
- **Simplest** — This mode automatically chooses between round toward floor and round toward zero to produce generated code that is as efficient as possible.
- **Zero** — This mode rounds toward zero and is equivalent to the MATLAB `fix` function.

For more information about each of these rounding modes, see “Rounding” on page 36-2.

Overflow Handling

To control how overflow conditions are handled for fixed-point operations, use the **Saturate on integer overflow** check box.

If this box is selected, overflows saturate to either the maximum or minimum value represented by the data type. For example, an overflow associated with a signed 8-bit integer can saturate to -128 or 127.

If this box is not selected, overflows wrap to the appropriate value that is representable by the data type. For example, the number 130 does not fit in a signed 8-bit integer, and would wrap to -126.

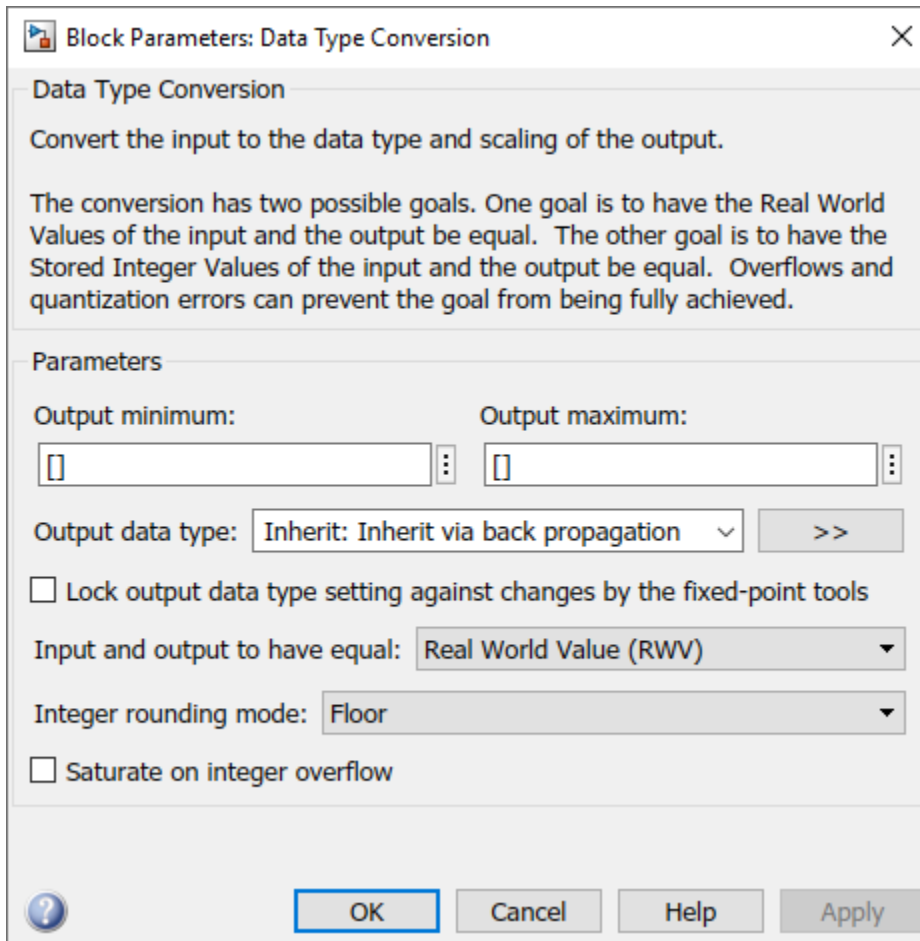
Lock the Output Data Type Setting

If the output data type is a generalized fixed-point number, you have the option of locking its output data type setting by selecting the **Lock output data type setting against changes by the fixed-point tools** check box.

When locked, the Fixed-Point Tool and automatic scaling script `autofixexp` do not change the output data type setting. Otherwise, the Fixed-Point Tool and `autofixexp` script are free to adjust the output data type setting.

Real-World Values Versus Stored Integer Values

You can configure Data Type Conversion blocks to treat signals as real-world values or as stored integers with the **Input and output to have equal** parameter.



The possible values are Real World Value (RWV) and Stored Integer (SI).

In terms of the variables defined in “Scaling” on page 35-5, the real-world value is given by V and the stored integer value is given by Q . You may want to treat numbers as stored integer values if you are modeling hardware that produces integers as output.

See Also

More About

- “Configure Blocks with Fixed-Point Parameters” on page 34-20

Configure Blocks with Fixed-Point Parameters

Certain Simulink blocks allow you to specify fixed-point numbers as the values of parameters used to compute the block's output, for example, the **Gain** parameter of a Gain block.

Note S-functions and the Stateflow Chart block do not support fixed-point parameters.

You can specify a fixed-point parameter value either directly by setting the value of the parameter to an expression that evaluates to a `fi` object, or indirectly by setting the value of the parameter to an expression that refers to a fixed-point `Simulink.Parameter` object.

Note Simulating or performing data type override on a model with `fi` objects requires a Fixed-Point Designer software license. See “Sharing Fixed-Point Models” on page 34-2 for more information.

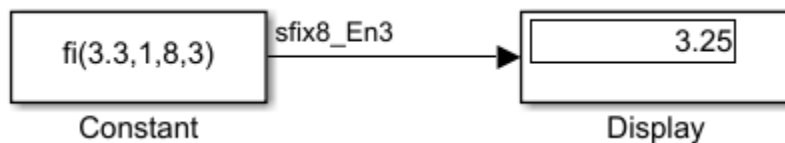
Specify Fixed-Point Values Directly

You can specify fixed-point values for block parameters using `fi` objects. In the block dialog's parameter field, simply enter the name of a `fi` object or an expression that includes the `fi` constructor function.

For example, entering the expression

```
fi(3.3,1,8,3)
```

as the **Constant value** parameter for the Constant block specifies a signed fixed-point value of 3.3, with a word length of 8 bits and a fraction length of 3 bits.



Specify Fixed-Point Values Via Parameter Objects

You can specify fixed-point parameter objects for block parameters using instances of the `Simulink.Parameter` class. To create a fixed-point parameter object, either specify a `fi` object as the parameter object's `Value` property, or specify the relevant fixed-point data type for the parameter object's `DataType` property.

For example, suppose that you want to create a fixed-point constant in your model. You could do this using a fixed-point parameter object and a Constant block as follows:

- 1 Enter the following command at the MATLAB prompt to create an instance of the `Simulink.Parameter` class:

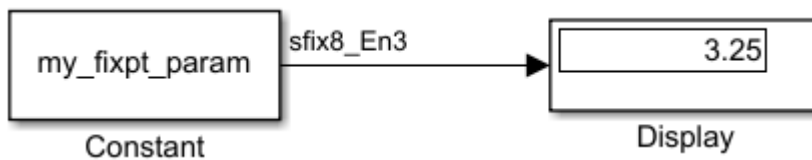

```
my_fixpt_param = Simulink.Parameter
```
- 2 Specify either the name of a `fi` object or an expression that includes the `fi` constructor function as the parameter object's `Value` property:

```
my_fixpt_param.Value = fi(3.3,1,8,3)
```

Alternatively, you can set the parameter object's `Value` and `DataType` properties separately. In this case, specify the relevant fixed-point data type using a `Simulink.AliasType` object, a `Simulink.NumericType` object, or a `fixdt` expression. For example, the following commands independently set the parameter object's value and data type, using a `fixdt` expression as the `DataType`:

```
my_fixpt_param.Value = 3.3;
my_fixpt_param.DataType = 'fixdt(1,8,2^-3,0)'
```

- 3 Specify the parameter object as the value of a block's parameter. For example, `my_fixpt_param` specifies the **Constant value** parameter for the Constant block in the following model:



The Constant block outputs a signed fixed-point value of 3.3, with a word length of 8 bits and a fraction length of 3 bits.

See Also

More About

- “Configure Blocks with Fixed-Point Output” on page 34-14

Pass Fixed-Point Data Between Simulink Models and MATLAB

You can read fixed-point data from the MATLAB software into your Simulink models, and there are several ways in which you can log fixed-point information from your models and simulations to the workspace.

Read Fixed-Point Data from the Workspace

Use the From Workspace block to read fixed-point data from the MATLAB workspace into a Simulink model. To do this, the data must be in structure format with a Fixed-Point Designer `fi` object in the `values` field. In array format, the From Workspace block only accepts real, double-precision data.

To read in `fi` data, the **Interpolate data** parameter of the From Workspace block must not be selected, and the **Form output after final data value by** parameter must be set to anything other than Extrapolation.

Write Fixed-Point Data to the Workspace

You can write fixed-point output from a model to the MATLAB workspace via the To Workspace block in either array or structure format. Fixed-point data written by a To Workspace block to the workspace in structure format can be read back into a Simulink model in structure format by a From Workspace block.

Note To write fixed-point data to the workspace as a `fi` object, select the **Log fixed-point data as a `fi` object** check box on the To Workspace block dialog. Otherwise, fixed-point data is converted to `double` and written to the workspace as `double`.

For example, you can use the following code to create a structure in the MATLAB workspace with a `fi` object in the `values` field. You can then use the From Workspace block to bring the data into a Simulink model.

```
a = fi([sin(0:10)' sin(10:-1:0)'])
```

```
a =
```

```

      0   -0.5440
  0.8415   0.4121
  0.9093   0.9893
  0.1411   0.6570
 -0.7568  -0.2794
 -0.9589  -0.9589
 -0.2794  -0.7568
  0.6570   0.1411
  0.9893   0.9093
  0.4121   0.8415
 -0.5440   0

```

```

      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 16
      FractionLength: 15

```

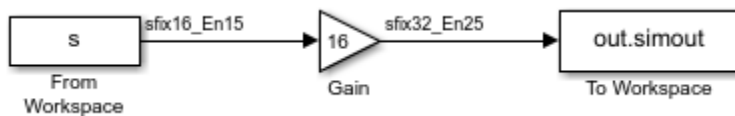
```

s.signals.values = a
s =
  struct with fields:
    signals: [1x1 struct]
s.signals.dimensions = 2
s =
  struct with fields:
    signals: [1x1 struct]
s.time = [0:10]'
s =
  struct with fields:
    signals: [1x1 struct]
    time: [11x1 double]

```

The From Workspace block in the following model has the `fi` structure `s` in the **Data** parameter. In the model, the following parameters in the **Solver** pane of the Configuration Parameters dialog box have the indicated settings:

- **Start time** — 0.0
- **Stop time** — 10.0
- **Type** — Fixed-step
- **Solver** — discrete (no continuous states)
- **Fixed-step size (fundamental sample time)** — 1.0



The To Workspace block writes the result of the simulation to the MATLAB workspace as a `fi` structure.

```
out.simout.data
```

```
ans =
```

```

      0   -8.7041
 13.4634   6.5938
 14.5488  15.8296
   2.2578  10.5117
-12.1089  -4.4707
-15.3428 -15.3428
  -4.4707 -12.1089
 10.5117   2.2578

```

```
15.8296    14.5488
 6.5938    13.4634
-8.7041         0
```

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 32
FractionLength: 25
```

Log Fixed-Point Signals

When fixed-point signals are logged to the MATLAB workspace via signal logging, they are always logged as Fixed-Point Designer `fi` objects.

To enable signal logging, first select the signal. Then, in the **Simulation** tab, click **Log Signals**.

For more information, refer to “Save Signal Data Using Signal Logging”.

When you log signals from a referenced model or Stateflow chart in your model, the word lengths of `fi` objects may be larger than you expect. The word lengths of fixed-point signals in referenced models and Stateflow charts are logged as the next larger data storage container size.

Access Fixed-Point Block Data During Simulation

Simulink provides an application programming interface (API) that enables programmatic access to block data, such as block inputs and outputs, parameters, states, and work vectors, while a simulation is running. You can use this interface to develop MATLAB programs capable of accessing block data while a simulation is running or to access the data from the MATLAB command line. Fixed-point signal information is returned to you via this API as `fi` objects. For more information about the API, refer to “Access Block Data During Simulation”.

See Also

To Workspace | From Workspace

More About

- “Save Signal Data Using Signal Logging”

Cast from Doubles to Fixed Point

In this section...

“Simulate Using Binary-Point-Only Scaling” on page 34-25

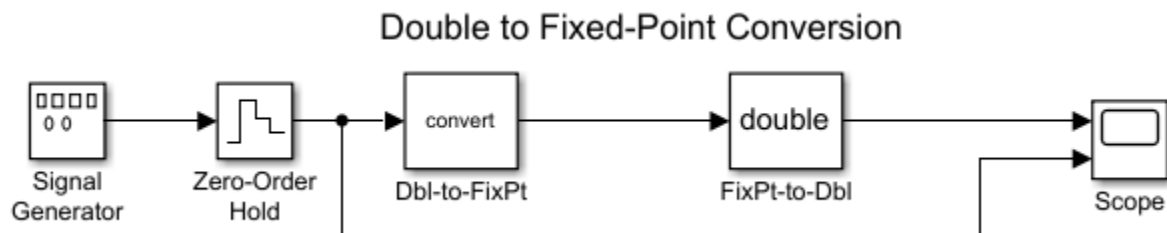
“Simulate Using [Slope Bias] Scaling” on page 34-27

This example shows you how to simulate a continuous real-world doubles signal using a generalized fixed-point data type. Using the `fxpdemo_dbl2fix` model, you can explore many of the important features of the Fixed-Point Designer software, including

- Data types
- Scaling
- Rounding
- Logging minimum and maximum simulation values to the workspace
- Overflow handling

To open the model, at the MATLAB command line, enter

```
openExample('fixedpoint/DoubleToFixedPointConversionExample')
```



In this example, you configure the Signal Generator block to output a sine wave signal with an amplitude defined on the interval $[-5 \ 5]$. The Signal Generator block always outputs double-precision numbers.

The `Dbl-to-FixPt` Data Type Conversion block converts the double-precision numbers from the Signal Generator block into one of the Fixed-Point Designer data types. For simplicity, the size of the output signal is 5 bits in this example.

The `FixPt-to-Dbl` Data Type Conversion block converts one of the Fixed-Point Designer data types into a Simulink data type. In this example, it outputs double-precision numbers.

Simulate Using Binary-Point-Only Scaling

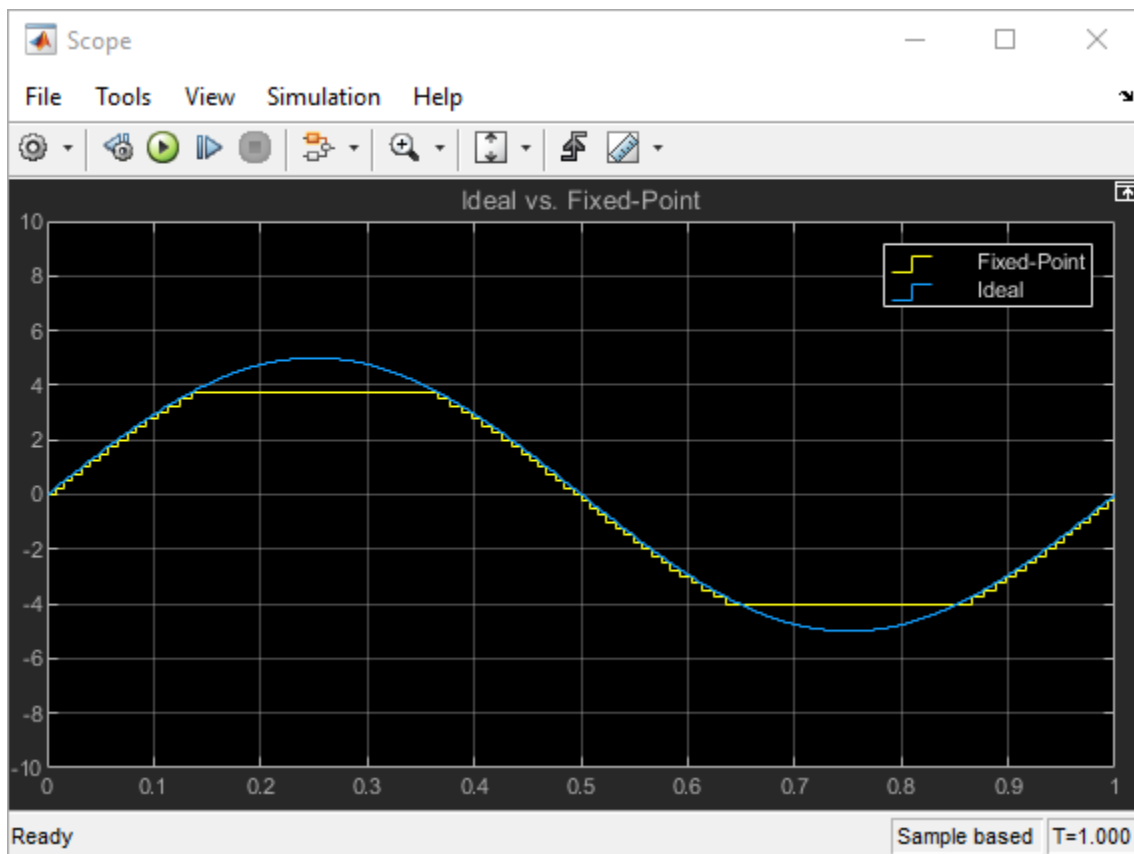
When using binary-point-only scaling, your goal is to find the optimal power-of-two exponent E , as defined in “Scaling” on page 35-5. For this scaling mode, the fractional slope F is 1 and there is no bias.

To set up the model to use binary-point-only scaling:

- 1 Configure the Signal Generator block to output a sine wave signal with an amplitude defined on the interval $[-5 \ 5]$.

- a Double-click the Signal Generator block to open the **Block Parameters** dialog.
 - b Set the **Wave form** parameter to `sine`.
 - c Set the **Amplitude** parameter to 5.
 - d Click **OK**.
- 2 Configure the Data Type Conversion (Dbl-to-FixPt) block.
 - a Double-click the **Dbl-to-FixPt** block to open the **Block Parameters** dialog.
 - b Verify that the **Output data type** parameter is `fixdt(1,5,2)`. This specifies a 5-bit, signed, fixed-point number with scaling 2^{-2} , which puts the binary point two places to the left of the rightmost bit. Hence the maximum value is $011.11 = 3.75$, the minimum value is $100.00 = -4.00$, and the precision is $(1/2)^2 = 0.25$.
 - c Verify that the **Integer rounding mode** parameter is set to `Floor`. This rounds the fixed-point result toward negative infinity.
 - d Select the **Saturate on integer overflow** check box to prevent the block from wrapping on overflow.
 - e Click **OK**.
 - 3 To simulate the model, in the **Simulation** tab, click **Run**.

The Scope displays the ideal and the fixed-point simulation results.



The simulation shows the quantization effects of fixed-point arithmetic. Using a 5-bit word with a precision of $(1/2)^2 = 0.25$ produces a discretized output that does not span the full range of the input signal.

To span the complete range of the input signal with 5 bits using binary-point-only scaling, set the output scaling to 2^{-1} . This puts the binary point one place to the left of the rightmost bit, giving a maximum value of $0111.1 = 7.5$ and a minimum value of $1000.0 = -8.0$. However, the precision is reduced to $(1/2)^1 = 0.5$. If you want to span the complete range of the input signal with 5 bits using binary-point-only scaling, then your only option is to sacrifice precision. Hence, the output scaling is 2^{-1} , which puts the binary point one place to the left of the rightmost bit. This scaling gives a maximum value of $0111.1 = 7.5$, a minimum value of $1000.0 = -8.0$, and a precision of $(1/2)^1 = 0.5$.

To see the effect of reducing the precision to 0.5, set the **Output data type** parameter of the **Dbl-to-FixPt Data Type Conversion** block to `fixdt(1,5,1)` and rerun the simulation.

Simulate Using [Slope Bias] Scaling

When using [Slope Bias] scaling, your goal is to find the optimal fractional slope F and fixed power-of-two exponent E , as defined in “Scaling” on page 35-5. There is no bias for this example because the sine wave is defined on the interval $[-5 \ 5]$.

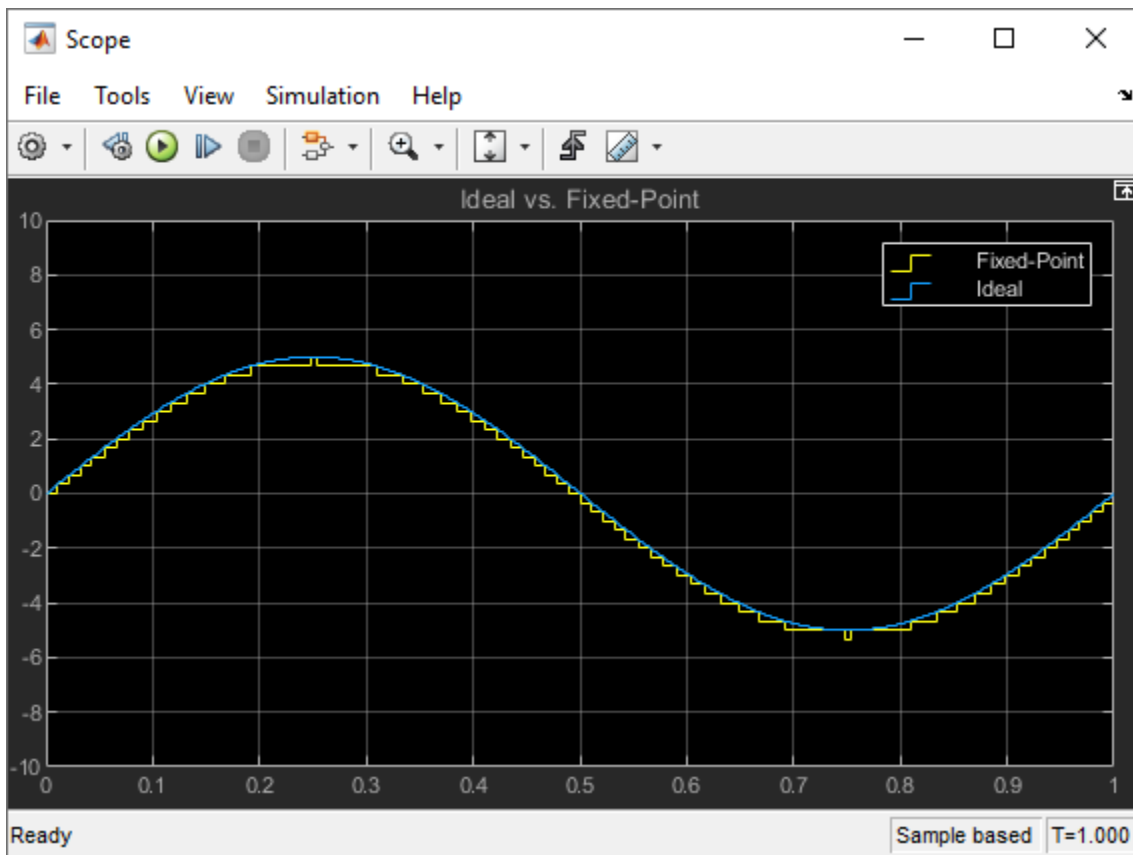
To find the slope, you begin by assuming a fixed power-of-two exponent of -2. To find the fractional slope, divide the maximum value of the sine wave by the maximum value of the scaled 5-bit number. The result is $5.00/3.75 = 1.3333$. The slope (and precision) is $1.3333(0.25) = 0.3333$. You specify the [Slope Bias] scaling as `[0.3333 0]` by entering the expression `fixdt(1,5,0.3333,0)` as the value of the **Output data type** parameter.

You could also specify a fixed power-of-two exponent of -1 and a corresponding fractional slope of 0.6667. The resulting slope is the same since E is reduced by 1 bit but F is increased by 1 bit. The Fixed-Point Designer software would automatically store F as 1.3332 and E as -2 because of the normalization condition of $1 \leq F < 2$.

To set up the model to use [Slope Bias] scaling:

- 1 Configure the Signal Generator block to output a sine wave signal with an amplitude defined on the interval $[-5 \ 5]$.
 - a Double-click the Signal Generator block to open the **Block Parameters** dialog.
 - b Set the **Wave form** parameter to `sine`.
 - c Set the **Amplitude** parameter to 5.
 - d Click **OK**.
- 2 Configure the **Dbl-to-FixPt Data Type Conversion** block.
 - a Double-click the **Dbl-to-FixPt** block to open the **Block Parameters** dialog.
 - b To specify a [Slope Bias] scaling of `[0.3333 0]`, set the **Output data type** parameter to `fixdt(1,5,0.3333,0)`.
 - c Verify that the **Integer rounding mode** parameter is `Floor`. This rounds the fixed-point result toward negative infinity.
 - d Select the **Saturate on integer overflow** check box to prevent the block from wrapping on overflow.
 - e Click **OK**.
- 3 To simulate the model, in the **Simulation** tab, click **Run**.

The Scope displays the ideal and the fixed-point simulation results.



If you do not know the slope, you can achieve reasonable simulation results by selecting your scaling based on the formula

$$\frac{(max_value - min_value)}{2^{ws} - 1},$$

where

- *max_value* is the maximum value to be simulated
- *min_value* is the minimum value to be simulated
- *ws* is the word size in bits
- $2^{ws} - 1$ is the largest value of a word with size *ws*

For this example, the formula produces a slope of 0.32258.

See Also

More About

- “Physical Quantities and Measurement Scales” on page 34-3
- “Scaling” on page 1-3

Data Types and Scaling

- “Data Types and Scaling in Digital Hardware” on page 35-2
- “Scaling” on page 35-5
- “Quantization” on page 35-7
- “Range and Precision” on page 35-9
- “Fixed-Point Numbers in Simulink” on page 35-13
- “Display Port Data Types” on page 35-15
- “Scaled Doubles” on page 35-16
- “Use Scaled Doubles to Avoid Precision Loss” on page 35-18
- “Floating-Point Numbers” on page 35-20

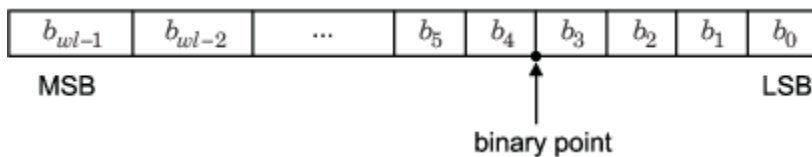
Data Types and Scaling in Digital Hardware

Fixed-Point Data Types

In digital hardware, numbers are stored in binary words. A binary word is a fixed-length sequence of bits (1's and 0's). How hardware components or software functions interpret this sequence of 1's and 0's is defined by the data type. Binary numbers are represented as either fixed-point or floating-point data types.

A fixed-point data type is characterized by the word length in bits, the position of the binary point, and whether it is signed or unsigned. The position of the binary point is the means by which fixed-point values are scaled and interpreted.

For example, a binary representation of a generalized fixed-point number (either signed or unsigned) is shown below:



where

- b_i is the i^{th} binary digit.
- wl is the word length in bits.
- b_{wl-1} is the location of the most significant, or highest, bit (MSB).
- b_0 is the location of the least significant, or lowest, bit (LSB).
- The binary point is shown four places to the left of the LSB. In this example, the number is said to have four fractional bits, or a fraction length of four.

Fixed-point data types can be either signed or unsigned. Whether a fixed-point value is signed or unsigned is usually not encoded explicitly within the binary word; that is, there is no sign bit. Instead, the sign information is implicitly defined within the computer architecture.

Signed binary fixed-point numbers are typically represented in computer hardware in one of three ways:

- Sign/magnitude - One bit of a binary word is always the dedicated sign bit, while the remaining bits of the word encode the magnitude of the number. Negation using sign/magnitude representation consists of flipping the sign bit from 0 (positive) to 1 (negative), or from 1 to 0.
- One's complement - Negating a binary number in one's complement requires a bitwise complement. That is, all 0's are flipped to 1's and all 1's are flipped to 0's. In one's complement notation there are two ways to represent zero. A binary word of all 0's represents "positive" zero, while a binary word of all 1's represents "negative" zero.
- Two's complement - Negation using signed two's complement representation consists of a bit inversion (translation into one's complement) followed by the binary addition of a one. For example, the two's complement of 000101 is 111011.

Two's complement is the most common representation of signed fixed-point numbers and is the only representation used by Fixed-Point Designer documentation.

Binary Point Interpretation

The binary point is the means by which fixed-point numbers are scaled. It is usually the software that determines the binary point. When performing basic math functions such as addition or subtraction, the hardware uses the same logic circuits regardless of the value of the scale factor. In essence, the logic circuits have no knowledge of a scale factor. They are performing signed or unsigned fixed-point binary algebra as if the binary point is to the right of b_0 .

Fixed-Point Designer supports general binary point scaling on page 35-5 $V = Q \times 2^E$, where V is the real-world value, Q is the stored integer value, and the fixed exponent E is equal to the negative of the fraction length. In other words, $RealWorldValue = StoredInteger \times 2^{-FractionLength}$.

The fraction length defines the scaling of the stored integer value. The word length limits the values that the stored integer can take, but it does not limit the values that the fraction length can take. The software does not restrict the value of the exponent E based on the word length of the stored integer Q . Because E is equal to $-FractionLength$, restricting the binary point to being contiguous with the fraction is unnecessary; the fraction length can be negative or greater than the word length.

For example, a word consisting of three unsigned bits is usually represented in scientific notation on page 35-20 in one of the following ways:

$$\begin{aligned}bbb. &= bbb. \times 2^0 \\bb.b &= bbb. \times 2^{-1} \\b.bb &= bbb. \times 2^{-2} \\\.bbb &= bbb. \times 2^{-3}\end{aligned}$$

If the exponent were greater than 0 or less than -3, then the representation would involve additional zeros:

$$\begin{aligned}bbb00000. &= bbb. \times 2^5 \\bbb00. &= bbb. \times 2^2 \\\.00bbb &= bbb. \times 2^{-5} \\\.00000bbb &= bbb. \times 2^{-8}\end{aligned}$$

These extra zeros never change to ones, so they do not show up in the hardware. Unlike floating-point exponents, a fixed-point exponent never shows up in the hardware, so fixed-point exponents are not limited by a finite number of bits.

Consider a signed value with a word length of 8, a fraction length of 10, and a stored integer value of 5 (binary value `00000101`). The real-world value is calculated using the formula $RealWorldValue = StoredInteger \times 2^{-FractionLength}$. In this case, $RealWorldValue = 5 \times 2^{-10} = 0.0048828125$. Because the fraction length is 2 bits longer than the word length, the binary value of the stored integer is `x.xx00000101`, where `x` is a placeholder for implicit zeros. `0.0000000101` (binary) is equivalent to `0.0048828125` (decimal). For an example using a `fi` object, see “Fraction Length Greater Than Word Length” on page 49-5.

Floating-Point Data Types

Floating-point data types are characterized by a sign bit, a fraction (or mantissa) field, and an exponent field. Fixed-Point Designer adheres to the IEEE Standard 754-1985 for Binary Floating-

Point Arithmetic (referred to simply as the IEEE Standard 754 throughout this guide) and supports half-, single- and double-precision data types.

When choosing a data type, you must consider these factors:

- The numerical range of the result
- The precision required of the result
- The associated quantization error (i.e., the rounding mode)
- The method for dealing with exceptional arithmetic conditions

These choices depend on your specific application, the computer architecture used, and the cost of development, among others.

With Fixed-Point Designer, you can explore the relationship between data types, range, precision, and quantization error in the modeling of dynamic digital systems. With Simulink Coder, you can generate production code based on that model. With HDL Coder, you can generate portable, synthesizable VHDL and Verilog code from Simulink models and Stateflow charts.

See Also

More About

- “Floating-Point Numbers” on page 35-20
- “Benefits of Fixed-Point Hardware”
- “Fixed-Point Numbers in Simulink” on page 35-13

Scaling

The dynamic range of fixed-point numbers is much less than floating-point numbers with equivalent word sizes. To avoid overflow conditions and minimize quantization errors, fixed-point numbers must be scaled.

With the Fixed-Point Designer software, you can select a fixed-point data type whose scaling is defined by its binary point, or you can select an arbitrary linear scaling that suits your needs. This section presents the scaling choices available for fixed-point data types.

You can represent a fixed-point number by a general slope and bias encoding scheme.

$$\text{real-world value} = (\text{slope} \times \text{integer}) + \text{bias}$$

where the slope can be expressed as

$$\text{slope} = \text{slope adjustment factor} \times 2^{\text{fixed exponent}}$$

The integer is sometimes called the *stored integer*. This is the raw binary number, in which the binary point is assumed to be at the far right of the word. In Fixed-Point Designer documentation, the negative of the fixed exponent is often referred to as the *fraction length*.

The slope and bias together represent the scaling of the fixed-point number. In a number with zero bias, only the slope affects the scaling. A fixed-point number that is only scaled by binary point position is equivalent to a number in slope bias representation that has a bias equal to zero and a slope adjustment factor equal to one. This is referred to as binary point-only scaling or power-of-two scaling:

$$\text{real-world value} = 2^{\text{fixed exponent}} \times \text{integer}$$

or

$$\text{real-world value} = 2^{-\text{fraction length}} \times \text{integer}$$

Binary-Point-Only Scaling

Binary-point-only or power-of-two scaling involves moving the binary point within the fixed-point word. The advantage of this scaling mode is to minimize the number of processor arithmetic operations.

With binary-point-only scaling, the components of the general slope and bias formula have the following values:

- bias = 0
- slope adjustment factor = 1
- slope = slope adjustment factor $\times 2^{\text{fixed exponent}} = 2^{\text{fixed exponent}}$

The scaling of a quantized real-world number is defined by the slope S , which is restricted to a power of two. The negative of the power-of-two exponent is called the fraction length. The fraction length is the number of bits to the right of the binary point. For Binary-Point-Only scaling, specify fixed-point data types as

- signed types — `fixdt(1,WordLength,FractionLength)`
- unsigned types — `fixdt(0,WordLength,FractionLength)`

Integers are a special case of fixed-point data types. Integers have a trivial scaling with slope 1 and bias 0, or equivalently with fraction length 0. Specify integers as

- signed integer — `fixdt(1,WordLength,0)`
- unsigned integer — `fixdt(0,WordLength,0)`

Slope and Bias Scaling

When you scale by slope and bias, the slope S and bias B of the quantized real-world number can take on any value. The slope must be a positive number. Using slope and bias, specify fixed-point data types as

- `fixdt(Signed,WordLength,Slope,Bias)`

Unspecified Scaling

Specify fixed-point data types with an unspecified scaling as

- `fixdt(Signed,WordLength)`

Simulink signals, parameters, and states must never have unspecified scaling. When scaling is unspecified, you must use some other mechanism such as automatic best precision scaling to determine the scaling that the Simulink software uses.

See Also

More About

- “Recommendations for Arithmetic and Scaling” on page 36-31

Quantization

The quantization Q of a real-world value V is represented by a weighted sum of bits. Within the context of the general slope and bias encoding scheme, the value of an unsigned fixed-point quantity is given by

$$\tilde{V} = S \cdot \left[\sum_{i=0}^{ws-1} b_i 2^i \right] + B,$$

while the value of a signed fixed-point quantity is given by

$$\tilde{V} = S \cdot \left[-b_{ws-1} 2^{ws-1} + \sum_{i=0}^{ws-2} b_i 2^i \right] + B,$$

where

- b_i are binary digits, with $b_i = 1, 0$, for $i = 0, 1, \dots, ws - 1$
- The word size in bits is given by ws , with $ws = 1, 2, 3, \dots, 128$.
- S is given by $F = 2^E$, where the scaling is unrestricted because the binary point does not have to be contiguous with the word.

b_i are called *bit multipliers* and 2^i are called the *weights*.

Fixed-Point Format

Formats for 8-bit signed and unsigned fixed-point values are shown in the following figure.

| | | | | | | | | |
|---|---|---|---|---|---|---|---|--------------------|
| 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | Unsigned data type |
| 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | Signed data type |

Note that you cannot discern whether these numbers are signed or unsigned data types merely by inspection since this information is not explicitly encoded within the word.

The binary number 0011.0101 yields the same value for the unsigned and two's complement representation because the MSB = 0. Setting $B = 0$ and using the appropriate weights, bit multipliers, and scaling, the value is

$$\begin{aligned} \tilde{V} &= (F2^E)Q = 2^E \left[\sum_{i=0}^{ws-1} b_i 2^i \right] \\ &= 2^{-4} (0 \times 2^7 + 0 \times 2^6 + 1 \times 2^5 + 1 \times 2^4 + 0 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0) \\ &= 3.3125. \end{aligned}$$

Conversely, the binary number 1011.0101 yields different values for the unsigned and two's complement representation since the MSB = 1.

Setting $B = 0$ and using the appropriate weights, bit multipliers, and scaling, the unsigned value is

$$\begin{aligned}\tilde{V} &= (F2^E)Q = 2^E \left[\sum_{i=0}^{ws-1} b_i 2^i \right] \\ &= 2^{-4} (1 \times 2^7 + 0 \times 2^6 + 1 \times 2^5 + 1 \times 2^4 + 0 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0) \\ &= 11.3125,\end{aligned}$$

while the two's complement value is

$$\begin{aligned}\tilde{V} &= (F2^E)Q = 2^E \left[-b_{ws-1} 2^{ws-1} + \sum_{i=0}^{ws-2} b_i 2^i \right] \\ &= 2^{-4} (-1 \times 2^7 + 0 \times 2^6 + 1 \times 2^5 + 1 \times 2^4 + 0 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0) \\ &= -4.6875.\end{aligned}$$

Range and Precision

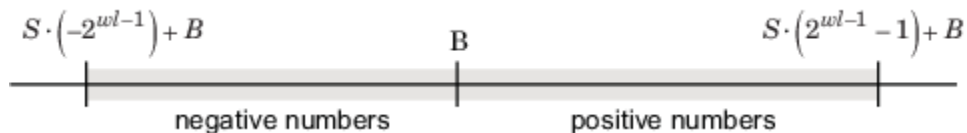
The *range* of a number gives the limits of the representation, while the *precision* gives the distance between successive numbers in the representation. The range and precision of a fixed-point number depend on the length of the word and the scaling.

Note You must pay attention to the precision and range of the fixed-point data types and scalings you choose in order to know whether rounding methods will be invoked or if overflows or underflows will occur.

Range

The range is the span of numbers that a fixed-point data type and scaling can represent. Range is limited because fixed-point words have limited size.

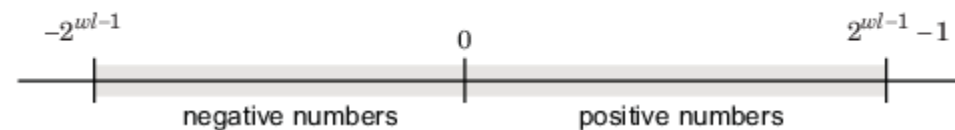
The range of representable numbers for a two's complement fixed-point number of word length wl , scaling S and bias B is illustrated below, where the values of wl , S , and B allow for both negative and positive numbers.



For both signed and unsigned fixed-point numbers of any data type, the number of different bit patterns is 2^{wl} .

For example, in two's complement, negative numbers must be represented as well as zero, so the maximum value is $2^{wl-1} - 1$. Because there is only one representation for zero, there are an unequal number of positive and negative numbers. This means there is a representation for -2^{wl-1} but not for 2^{wl-1} :

For slope = 1 and bias = 0:



Limitations on Range

Because a fixed-point data type represents numbers within a finite range, overflows and underflows can occur if the result of an operation is larger or smaller than the numbers in that range.

In binary arithmetic, a processor might need to take an n -bit fixed-point number and store it in m bits, where $m \neq n$. If $m < n$, the range of the number has been reduced and an operation can produce an overflow condition. Some processors identify this condition as Inf or NaN. For other processors, especially digital signal processors (DSPs), the value *saturates* or *wraps*.

Fixed-Point Designer software allows you to either *saturate* or *wrap* overflows. Saturation represents positive overflows as the largest positive number in the range being used, and negative overflows as

the largest negative number in the range being used. Wrapping uses modulo arithmetic to cast an overflow back into the representable range of the data type.

When you create a `fi` object, any overflows are saturated. The `OverflowAction` property of the default `fimath` is `saturate`. You can log overflows and underflows by setting the `LoggingMode` property of the `fipref` object to `on`.

If $m > n$, the range of the number has been extended. Extending the range of a word requires the inclusion of *guard bits*, which act to guard against potential overflow.

The Simulink software supports saturation and wrapping for all fixed-point data types, while guard bits are supported only for fractional data types.

Precision

The precision of a fixed-point number is the difference between successive values representable by its data type and scaling. The value of the least significant bit, and therefore the precision of the number, is determined by the number of fractional bits. A fixed-point value can be represented to within half of the precision of its data type and scaling.

For example, a fixed-point representation with four bits to the right of the binary point has a precision of 2^{-4} or 0.0625, which is the value of its least significant bit. Any number within the range of this data type and scaling can be represented to within $(2^{-4})/2$ or 0.03125, which is half the precision. This is an example of representing a number with finite precision.

Limitations on Precision

The precision of a fixed-point word depends on the word size and binary point location. For example, suppose you must represent the real-world number 35.375 with a fixed-point number. Using a slope bias encoding scheme, the representation is

$$V \approx \tilde{V} = SQ + B = 2^{-2}Q + 32,$$

where $V = 35.375$.

The two closest approximations to the real-world value are $Q = 13$ and $Q = 14$:

$$\tilde{V} = 2^{-2}(13) + 32 = 35.25,$$

$$\tilde{V} = 2^{-2}(14) + 32 = 35.50.$$

In either case, the absolute error is the same:

$$|\tilde{V} - V| = 0.125 = \frac{S}{2} = \frac{F2^E}{2}.$$

For fixed-point values within the limited range, this represents the worst-case error if round-to-nearest is used. If other rounding modes are used, the worst-case error can be twice as large:

$$|\tilde{V} - V| < F2^E.$$

Extending the precision of a word can be accomplished with more bits, but you face practical limitations with this approach. Instead, you must carefully select the data type, word size, and scaling

such that numbers are accurately represented. Rounding and padding with trailing zeros are typical methods implemented on processors to deal with the precision of binary words.

Fixed-Point Data Type Parameters

The low limit, high limit, and default binary-point-only scaling for the supported fixed-point data types discussed in “Binary-Point-Only Scaling” on page 35-5 are given in the following table.

Fixed-Point Data Type Range and Default Scaling

| Name | Data Type | Low Limit | High Limit | Default Scaling (~Precision) |
|-----------------------|-----------------------------|--------------------|-------------------------|------------------------------|
| Unsigned Integer | $\text{fixdt}(0, ws, 0)$ | 0 | $2^{ws} - 1$ | 1 |
| Signed Integer | $\text{fixdt}(1, ws, 0)$ | -2^{ws-1} | $2^{ws-1} - 1$ | 1 |
| Unsigned Binary Point | $\text{fixdt}(0, ws, fl)$ | 0 | $(2^{ws} - 1)2^{-fl}$ | 2^{-fl} |
| Signed Binary Point | $\text{fixdt}(1, ws, fl)$ | $-2^{ws-1-fl}$ | $(2^{ws-1} - 1)2^{-fl}$ | 2^{-fl} |
| Unsigned Slope Bias | $\text{fixdt}(0, ws, s, b)$ | b | $s(2^{ws} - 1) + b$ | s |
| Signed Slope Bias | $\text{fixdt}(1, ws, s, b)$ | $-s(2^{ws-1}) + b$ | $s(2^{ws-1} - 1) + b$ | s |

s = Slope, b = Bias, ws = WordLength, fl = FractionLength

Range and Precision of an 8-Bit Fixed-Point Data Type – Binary-Point-Only Scaling

The precisions, range of signed values, and range of unsigned values for an 8-bit generalized fixed-point data type with binary-point-only scaling are listed in the follow table. Note that the first scaling value (2^1) represents a binary point that is not contiguous with the word.

| Scaling | Precision | Range of Signed Values (Low, High) | Range of Unsigned Values (Low, High) |
|----------|-----------|------------------------------------|--------------------------------------|
| 2^1 | 2.0 | -256, 254 | 0, 510 |
| 2^0 | 1.0 | -128, 127 | 0, 255 |
| 2^{-1} | 0.5 | -64, 63.5 | 0, 127.5 |
| 2^{-2} | 0.25 | -32, 31.75 | 0, 63.75 |
| 2^{-3} | 0.125 | -16, 15.875 | 0, 31.875 |
| 2^{-4} | 0.0625 | -8, 7.9375 | 0, 15.9375 |
| 2^{-5} | 0.03125 | -4, 3.96875 | 0, 7.96875 |
| 2^{-6} | 0.015625 | -2, 1.984375 | 0, 3.984375 |
| 2^{-7} | 0.0078125 | -1, 0.9921875 | 0, 1.9921875 |

| Scaling | Precision | Range of Signed Values (Low, High) | Range of Unsigned Values (Low, High) |
|----------|------------|------------------------------------|--------------------------------------|
| 2^{-8} | 0.00390625 | -0.5, 0.49609375 | 0, 0.99609375 |

Range and Precision of an 8-Bit Fixed-Point Data Type — Slope and Bias Scaling

The precision and ranges of signed and unsigned values for an 8-bit fixed-point data type using slope and bias scaling are listed in the following table. The slope starts at a value of 1.25 with a bias of 1.0 for all slopes. Note that the slope is the same as the precision.

| Bias | Slope/Precision | Range of Signed Values (low, high) | Range of Unsigned Values (low, high) |
|------|-----------------|------------------------------------|--------------------------------------|
| 1 | 1.25 | -159, 159.75 | 1, 319.75 |
| 1 | 0.625 | -79, 80.375 | 1, 160.375 |
| 1 | 0.3125 | -39, 40.6875 | 1, 80.6875 |
| 1 | 0.15625 | -19, 20.84375 | 1, 40.84375 |
| 1 | 0.078125 | -9, 10.921875 | 1, 20.921875 |
| 1 | 0.0390625 | -4, 5.9609375 | 1, 10.9609375 |
| 1 | 0.01953125 | -1.5, 3.48046875 | 1, 5.98046875 |
| 1 | 0.009765625 | -0.25, 2.240234375 | 1, 3.490234375 |
| 1 | 0.0048828125 | 0.375, 1.6201171875 | 1, 2.2451171875 |

See Also

More About

- “Saturation and Wrapping” on page 36-25
- “Rounding” on page 36-2
- “Guard Bits” on page 36-28

Fixed-Point Numbers in Simulink

Simulink data type names must be valid MATLAB identifiers with less than 128 characters. The data type name provides information about container type, number encoding, and scaling.

You can represent a fixed-point number using the fixed-point scaling equation

$$V \approx \tilde{V} = SQ + B,$$

where

- V is the real-world value.
- \tilde{V} is the approximate real-world value.
- $S = F2^E$ is the slope.
- F is the slope adjustment factor.
- E is the fixed power-of-two exponent.
- Q is the stored integer.
- B is the bias.

Fixed-Point Data Type and Scaling Notation

The following table provides a key for various symbols that appear in Simulink products to indicate the data type and scaling of a fixed-point value.

| Symbol | Description | Example |
|-------------------------|--|---|
| Container Type | | |
| <code>ufix</code> | Unsigned fixed-point data type | <code>ufix8</code> is an 8-bit unsigned fixed-point data type |
| <code>sfix</code> | Signed fixed-point data type | <code>sfix128</code> is a 128-bit signed fixed-point data type |
| <code>fltu</code> | Scaled Doubles on page 35-16 override of an unsigned fixed-point data type (<code>ufix</code>) | <code>fltu32</code> is a scaled doubles override of <code>ufix32</code> |
| <code>flts</code> | Scaled Doubles on page 35-16 override of a signed fixed-point data type (<code>sfix</code>) | <code>flts64</code> is a scaled doubles override of <code>sfix64</code> |
| Number Encoding | | |
| <code>e</code> | 10^{\wedge} | <code>125e8</code> equals $125 * (10^{\wedge}(8))$ |
| <code>n</code> | Negative | <code>n31</code> equals -31 |
| <code>p</code> | Decimal point | <code>1p5</code> equals 1.5 <code>p2</code> equals 0.2 |
| Scaling Encoding | | |

| Symbol | Description | Example |
|---------------|---|--|
| S | Slope | <code>ufix16_S5_B7</code> is a 16-bit unsigned fixed-point data type with Slope of 5 and Bias of 7 |
| B | Bias | <code>ufix16_S5_B7</code> is a 16-bit unsigned fixed-point data type with Slope of 5 and Bias of 7 |
| E | Fixed exponent (2^E) A negative fixed exponent describes the fraction length | <code>sfix32_En31</code> is a 32-bit signed fixed-point data type with a fraction length of 31 |
| F | Slope adjustment factor | <code>ufix16_F1p5_En50</code> is a 16-bit unsigned fixed-point data type with a SlopeAdjustmentFactor of 1.5 and a FixedExponent of -50 |
| C, c, D, or d | Compressed encoding for Bias Note If you pass this symbol to the <code>sldataTypeAndScale</code> function, it returns a valid <code>fixdt</code> data type. | No example available. For backwards compatibility only. To identify and replace calls to <code>sldataTypeAndScale</code> , use the “Check for calls to <code>sldataTypeAndScale</code> ” Model Advisor check. |
| T or t | Compressed encoding for Slope Note If you pass this symbol to the <code>sldataTypeAndScale</code> , it returns a valid <code>fixdt</code> data type. | No example available. For backwards compatibility only. To identify and replace calls to <code>sldataTypeAndScale</code> , use the “Check for calls to <code>sldataTypeAndScale</code> ” Model Advisor check. |

See Also

More About

- “Scaling” on page 35-5
- “Data Types and Scaling in Digital Hardware” on page 35-2

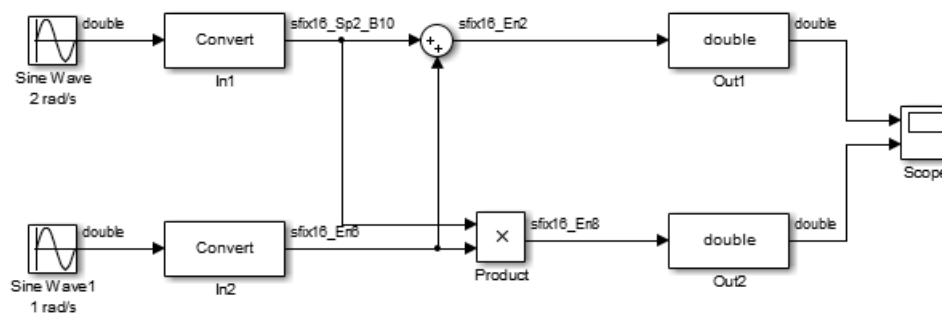
Display Port Data Types

To display the data types for the ports in your model.

- 1 On the Simulink **Debug** tab, select **Information Overlays > Base Data Types**.

The port display for fixed-point signals consists of three parts: the data type, the number of bits, and the scaling. These three parts reflect the block **Output data type** parameter value or the data type and scaling that is inherited from the driving block or through back propagation.

The following model displays its port data types.



In the model, the data type displayed with the In1 block indicates that the output data type name is `sfix16_Sp2_B10`. This corresponds to `fixdt(1, 16, 0.2, 10)` which is a signed 16 bit fixed-point number with slope 0.2 and bias 10.0. The data type displayed with the In2 block indicates that the output data type name is `sfix16_En6`. This corresponds to `fixdt(1, 16, 6)` which is a signed 16 bit fixed-point number with fraction length of 6.

See Also

More About

- “Fixed-Point Data Type and Scaling Notation” on page 35-13

Scaled Doubles

What Are Scaled Doubles?

Scaled doubles are a hybrid between floating-point and fixed-point numbers. The Fixed-Point Designer software stores them as doubles with the scaling, sign, and word length information retained. For example, the storage container for a fixed-point data type `sfix16_En14` is `int16`. The storage container of the equivalent scaled doubles data type, `flts16_En14` is floating-point double. The Fixed-Point Designer software applies the scaling information to the stored floating-point double to obtain the real-world value. Storing the value in a double almost always eliminates overflow and precision issues.

What is the Difference between Scaled Double and Double Data Types?

The storage container for both the scaled double and double data types is floating-point double. Therefore both data type override settings, `Double` and `Scaled double`, provide the range and precision advantages of floating-point doubles. Scaled doubles retain the information about the specified data type and scaling, but doubles do not retain this information. Because scaled doubles retain the information about the specified scaling, they can also be used for overflow detection.

Consider an example where you want to store a value of `0.75001` degrees Celsius in a data type `sfix16_En13`. For this data type:

- The slope is $S = 2^{-13}$.
- The bias is $B = 0$.

Using the scaling equation $V \approx \tilde{V} = SQ + B$, where V is the real-world value and Q is the stored integer value:

- $B = 0$.
- $\tilde{V} = SQ = 2^{-13}Q = 0.75001$.

Because the storage container of the data type `sfix16_En13` is 16 bits, the stored integer Q can only be represented as an integer within these 16 bits. Therefore, the ideal value of Q is quantized to `6144`, causing precision loss.

If you override the data type `sfix16_En13` with `Double`, the data type changes to `Double` and you lose the information about the scaling. The stored-value equals the real-world value `0.75001`.

If you override the data type `sfix16_En13` with `Scaled Double`, the data type changes to `flts16_En13`. The scaling is still given by `_En13` and is identical to that of the original data type. The only difference is the storage container used to hold the stored value which is now double so the stored-value is `6144.08192`. This example shows one advantage of using scaled doubles: the virtual elimination of quantization errors.

When to Use Scaled Doubles

The **Fixed-Point Tool** enables you to perform various data type overrides on fixed-point signals in your simulations. Use scaled doubles to override the fixed-point data types and scaling using double-

precision numbers to avoid quantization effects. Overriding the fixed-point data types provides a floating-point benchmark that represents the ideal output.

Scaled doubles are useful for:

- Testing and debugging
- Detecting overflows
- Applying data type overrides to individual subsystems

If you apply a data type override to subsystems in your model rather than to the whole model, Scaled doubles provide the information that the fixed-point portions of the model need for consistent data type propagation.

See Also

More About

- “Fixed-Point Data Type and Scaling Notation” on page 35-13
- “Use Scaled Doubles to Avoid Precision Loss” on page 35-18

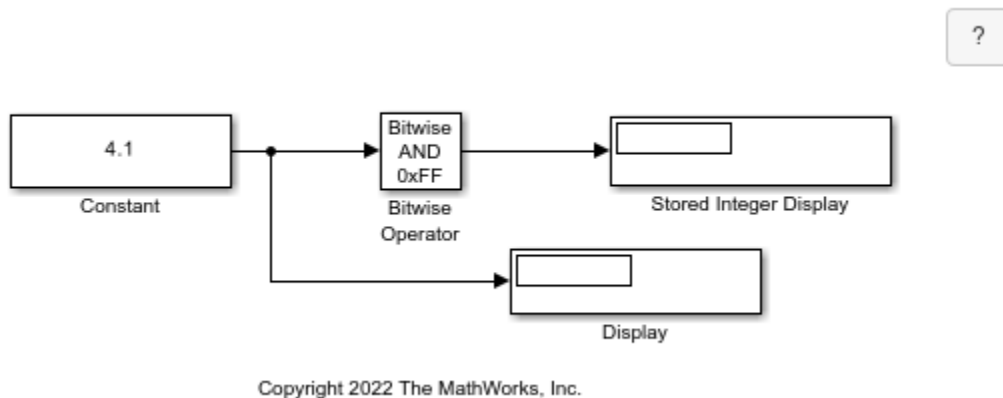
Use Scaled Doubles to Avoid Precision Loss

This example shows how you can avoid precision loss by overriding the data types in your model with scaled doubles.

Open the Model

Open the `ex_scaled_double` model.

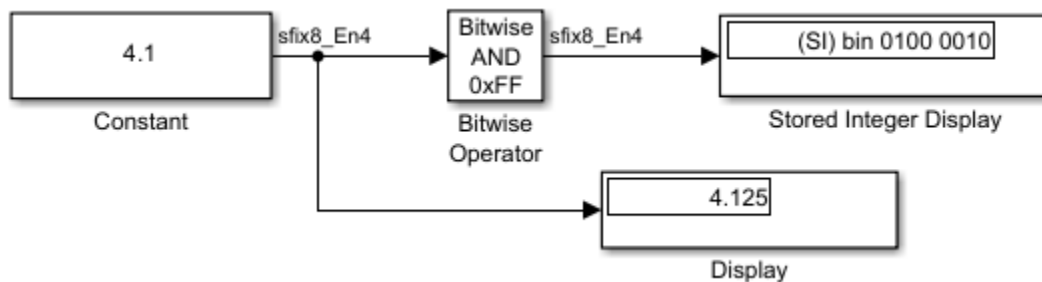
`ex_scaled_double`



In the `ex_scaled_double` model, the Constant block **Output data type** is `fixdt(1,8,4)`. The Bitwise Operator block uses the AND operator and the bit mask `0xFF` to pass the input value to the output. Because the **Treat mask as** parameter is set to **Stored Integer**, the block outputs the stored integer value Q of its input. The encoding scheme is $V = SQ + B$, where V is the real-world value and Q is the stored integer value.

Collect Ranges in the Fixed-Point Tool

- 1 From the Simulink® **Apps** tab, select **Fixed-Point Tool**.
- 2 Select the **Range Collection** workflow.
- 3 In the Fixed-Point Tool, select **Collect Ranges > Use current settings**. Click **Collect Ranges**.



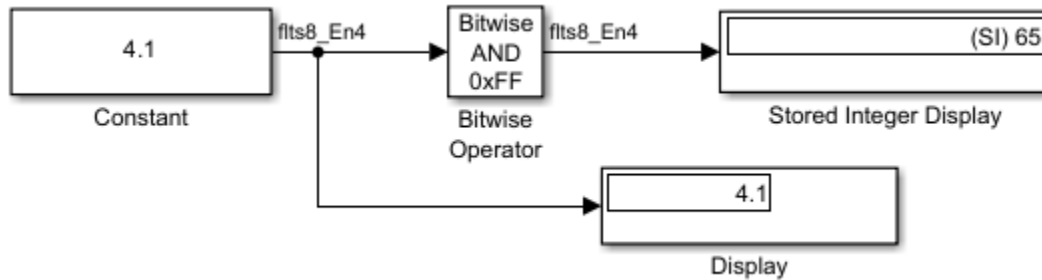
The Display block displays `4.125` as the output value of the Constant block. The Stored Integer Display block displays `(SI) bin 0100 0010`, which is the binary equivalent of the stored integer

value. Precision loss occurs because the output data type, `fixdt(1,8,4)`, cannot represent the output value 4.1 exactly.

Override Data Types with Scaled Doubles

Override data types in the model with scaled doubles.

- 1 In the Fixed-Point Tool, select **New > Range Collection** to start a new workflow.
- 2 Select **Collect Ranges > Scaled double precision**. Click **Collect Ranges**.



The Display block correctly displays 4.1 as the output value of the Constant block. The Stored Integer Display block displays (SI) 65, which is the binary equivalent of the stored integer value. Because the model uses scaled doubles to override the data type `fixdt(1,8,4)`, the compiled output data type changes to `flts8_En4`, which is the scaled double equivalent of `fixdt(1,8,4)`. No precision loss occurs because the scaled doubles use a double to hold the stored value and retain information about the specified data type and scaling. Note that you cannot use a data type override setting of **Double precision** or **Single precision** on this model because the Bitwise Operator block does not support floating-point data types.

See Also

More About

- “Scaled Doubles” on page 35-16
- “Fixed-Point Instrumentation and Data Type Override” on page 42-61
- “Fixed-Point Numbers in Simulink” on page 35-13

Floating-Point Numbers

In this section...

“Floating-Point Numbers” on page 35-20

“Scientific Notation” on page 35-20

“IEEE 754 Standard for Floating-Point Numbers” on page 35-21

“Range and Precision” on page 35-22

“Exceptional Arithmetic” on page 35-24

Floating-Point Numbers

Fixed-point numbers are limited in that they cannot simultaneously represent very large or very small numbers using a reasonable word size. This limitation can be overcome by using scientific notation. With scientific notation, you can dynamically place the binary point at a convenient location and use powers of the binary to keep track of that location. Thus, you can represent a range of very large and very small numbers with only a few digits.

You can represent any binary floating-point number in scientific notation form as $f2^e$, where f is the fraction (or mantissa), 2 is the radix or base (binary in this case), and e is the exponent of the radix. The radix is always a positive number, while f and e can be positive or negative.

When performing arithmetic operations, floating-point hardware must take into account that the sign, exponent, and fraction are all encoded within the same binary word. This results in complex logic circuits when compared with the circuits for binary fixed-point operations.

The Fixed-Point Designer software supports half-precision, single-precision, and double-precision floating-point numbers as defined by the IEEE Standard 754.

Scientific Notation

A direct analogy exists between scientific notation and radix point notation. For example, scientific notation using five decimal digits for the fraction would take the form

$$\pm d . dddd \times 10^p = \pm dddd.0 \times 10^{p-4} = \pm 0.ddd dd \times 10^{p+1},$$

where $d = 0, \dots, 9$ and p is an integer of unrestricted range.

Radix point notation using five bits for the fraction is the same except for the number base

$$\pm b . bbbb \times 2^q = \pm bbbbb.0 \times 2^{q-4} = \pm 0.bbbb b \times 2^{q+1},$$

where $b = 0, 1$ and q is an integer of unrestricted range.

For fixed-point numbers, the exponent is fixed but there is no reason why the binary point must be contiguous with the fraction. For more information, see “Binary Point Interpretation” on page 35-3.

IEEE 754 Standard for Floating-Point Numbers

The IEEE Standard 754 has been widely adopted, and is used with virtually all floating-point processors and arithmetic coprocessors, with the notable exception of many DSP floating-point processors.

This standard specifies several floating-point number formats, of which singles and doubles are the most widely used. Each format contains three components: a sign bit, a fraction field, and an exponent field.

The Sign Bit

IEEE floating-point numbers use sign/magnitude representation, where the sign bit is explicitly included in the word. Using sign/magnitude representation, a sign bit of 0 represents a positive number and a sign bit of 1 represents a negative number. This is in contrast to the two's complement representation preferred for signed fixed-point numbers.

The Fraction Field

Floating-point numbers can be represented in many different ways by shifting the number to the left or right of the binary point and decreasing or increasing the exponent of the binary by a corresponding amount.

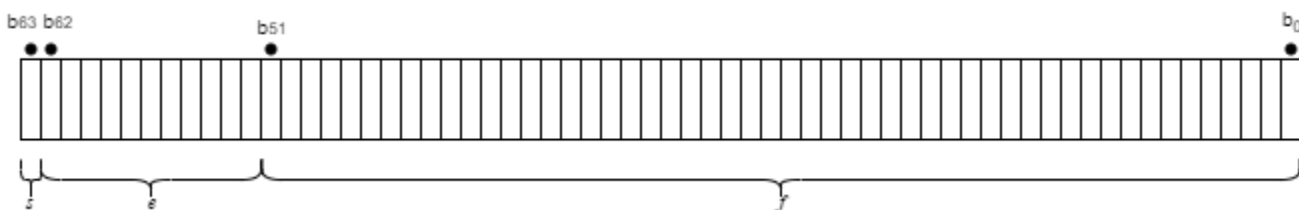
To simplify operations on floating-point numbers, they are normalized in the IEEE format. A normalized binary number has a fraction of the form $1.f$, where f has a fixed size for a given data type. Since the leftmost fraction bit is always a 1, it is unnecessary to store this bit and it is therefore implicit (or hidden). Thus, an n -bit fraction stores an $n+1$ -bit number. The IEEE format also supports "Denormalized Numbers" on page 35-24, which have a fraction of the form $0.f$.

The Exponent Field

In the IEEE format, exponent representations are biased. This means a fixed value, the bias, is subtracted from the exponent field to get the true exponent value. For example, if the exponent field is 8 bits, then the numbers 0 through 255 are represented, and there is a bias of 127. Note that some values of the exponent are reserved for flagging *Inf* (infinity), *NaN* (not-a-number), and denormalized numbers, so the true exponent values range from -126 to 127. See "Inf" on page 35-24 and "NaN" on page 35-24 for more information.

Double-Precision Format

The IEEE double-precision floating-point format is a 64-bit word divided into a 1-bit sign indicator s , an 11-bit biased exponent e , and a 52-bit fraction f .



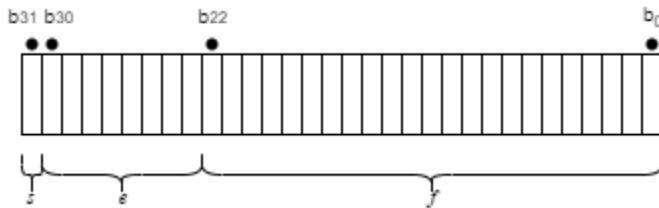
The relationship between double-precision format and the representation of real numbers is given by

$$\text{value} = \begin{cases} (-1)^s(2^{e-1023})(1.f) & \text{normalized, } 0 < e < 2047, \\ (-1)^s(2^{e-1022})(0.f) & \text{denormalized, } e = 0, f > 0, \\ \text{exceptional value} & \text{otherwise.} \end{cases}$$

See “Exceptional Arithmetic” on page 35-24 for more information.

Single-Precision Format

The IEEE single-precision floating-point format is a 32-bit word divided into a 1-bit sign indicator s , an 8-bit biased exponent e , and a 23-bit fraction f .



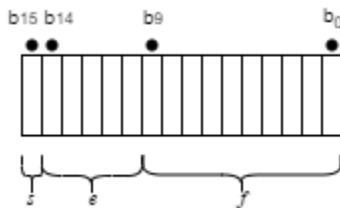
The relationship between single-precision format and the representation of real numbers is given by

$$\text{value} = \begin{cases} (-1)^s(2^{e-127})(1.f) & \text{normalized, } 0 < e < 255, \\ (-1)^s(2^{e-126})(0.f) & \text{denormalized, } e = 0, f > 0, \\ \text{exceptional value} & \text{otherwise.} \end{cases}$$

See “Exceptional Arithmetic” on page 35-24 for more information.

Half-Precision Format

The IEEE half-precision floating-point format is a 16-bit word divided into a 1-bit sign indicator s , a 5-bit biased exponent e , and a 10-bit fraction f .



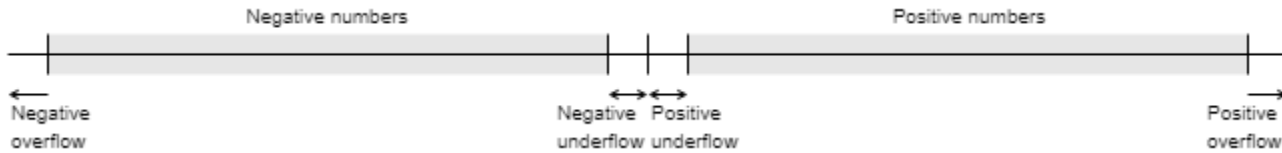
Half-precision numbers are supported in MATLAB and Simulink. For more information, see `half` and “The Half-Precision Data Type in Simulink” on page 51-2.

Range and Precision

The range of a number gives the limits of the representation. The precision gives the distance between successive numbers in the representation. The range and precision of an IEEE floating-point number depends on the specific format.

Range

The range of representable numbers for an IEEE floating-point number with f bits allocated for the fraction, e bits allocated for the exponent, and the bias of e given by $bias = 2^{(e-1)} - 1$ is given below.



where

- Normalized positive numbers are defined within the range $2^{(1-bias)}$ to $(2-2^{-f})2^{bias}$.
- Normalized negative numbers are defined within the range $-2^{(1-bias)}$ to $-(2-2^{-f})2^{bias}$.
- Positive numbers greater than $(2-2^{-f})2^{bias}$ and negative numbers less than $-(2-2^{-f})2^{bias}$ are overflows.
- Positive numbers less than $2^{(1-bias)}$ and negative numbers greater than $-2^{(1-bias)}$ are either underflows or denormalized numbers.
- Zero is given by a special bit pattern, where $e = 0$ and $f = 0$.

Overflows and underflows result from exceptional arithmetic conditions. Floating-point numbers outside the defined range are always mapped to $\pm Inf$.

Note You can use the MATLAB commands `realmin` and `realmax` to determine the dynamic range of double-precision floating-point values for your computer.

Precision

A floating-point number is only an approximation of the “true” value because of a finite word size. Therefore, it is important to have an understanding of the precision (or accuracy) of a floating-point result. A value v with an accuracy q is specified by $v \pm q$. For IEEE floating-point numbers,

$$v = (-1)^s (2^{e-bias}) (1.f)$$

and

$$q = 2^{-f} \times 2^{e-bias}$$

Thus, the precision is associated with the number of bits in the fraction field.

Note In the MATLAB software, floating-point relative accuracy is given by the command `eps`, which returns the distance from 1.0 to the next larger floating-point number. For a computer that supports the IEEE Standard 754, $eps = 2^{-52}$ or $2.22045 \cdot 10^{-16}$.

Floating-Point Data Type Parameters

The range, bias, and precision for supported floating-point data types are given in the table below.

| Data Type | Low Limit | High Limit | Exponent Bias | Precision |
|-----------|---------------------------------------|---|---------------|----------------------------|
| Half | $2^{-14} \approx 6.1 \cdot 10^{-5}$ | $(2-2^{-10}) \cdot 2^{15} \approx 6.5 \cdot 10^4$ | 15 | $2^{-10} \approx 10^{-3}$ |
| Single | $2^{-126} \approx 10^{-38}$ | $2^{128} \approx 3 \cdot 10^{38}$ | 127 | $2^{-23} \approx 10^{-7}$ |
| Double | $2^{-1022} \approx 2 \cdot 10^{-308}$ | $2^{1024} \approx 2 \cdot 10^{308}$ | 1023 | $2^{-52} \approx 10^{-16}$ |

Because floating-point numbers are represented using sign/magnitude, there are two representations of zero, one positive and one negative. For both representations $e = 0$ and $f.0 = 0.0$.

Exceptional Arithmetic

The IEEE Standard 754 specifies practices and procedures so that predictable results are produced independently of the hardware platform. Denormalized numbers, `Inf`, and `NaN` are defined to deal with exceptional arithmetic (underflow and overflow).

If an underflow or overflow is handled as `Inf` or `NaN`, then significant processor overhead is required to deal with this exception. Although the IEEE Standard 754 specifies practices and procedures to deal with exceptional arithmetic conditions in a consistent manner, microprocessor manufacturers might handle these conditions in ways that depart from the standard.

Denormalized Numbers

Denormalized numbers are used to handle cases of exponent underflow. When the exponent of the result is too small (i.e., a negative exponent with too large a magnitude), the result is denormalized by right-shifting the fraction and leaving the exponent at its minimum value. The use of denormalized numbers is also referred to as gradual underflow. Without denormalized numbers, the gap between the smallest representable nonzero number and zero is much wider than the gap between the smallest representable nonzero number and the next larger number. Gradual underflow fills that gap and reduces the impact of exponent underflow to a level comparable with round-off among the normalized numbers. Denormalized numbers provide extended range for small numbers at the expense of precision.

Inf

Arithmetic involving `Inf` (infinity) is treated as the limiting case of real arithmetic, with infinite values defined as those outside the range of representable numbers, or $-\infty \leq$ (representable numbers) $< \infty$. With the exception of the special cases discussed below (`NaN`), any arithmetic operation involving `Inf` yields `Inf`. `Inf` is represented by the largest biased exponent allowed by the format and a fraction of zero.

NaN

A `NaN` (not-a-number) is a symbolic entity encoded in floating-point format. There are two types of `NaN`: signaling and quiet. A signaling `NaN` signals an invalid operation exception. A quiet `NaN` propagates through almost every arithmetic operation without signaling an exception. The following operations result in a `NaN`: $\infty - \infty$, $-\infty + \infty$, $0 \times \infty$, $0/0$, and ∞/∞ .

Both signaling `NaN` and quiet `NaN` are represented by the largest biased exponent allowed by the format and a nonzero fraction. The bit pattern for a quiet `NaN` is given by $0.f$, where the most significant bit in f must be a one. The bit pattern for a signaling `NaN` is given by $0.f$, where the most significant bit in f must be zero and at least one of the remaining bits must be nonzero.

See Also**More About**

- “Data Types and Scaling in Digital Hardware” on page 35-2

Arithmetic Operations

- “Rounding” on page 36-2
- “Rounding Modes for Fixed-Point Simulink Blocks” on page 36-5
- “Rounding Mode: Ceiling” on page 36-8
- “Rounding Mode: Convergent” on page 36-9
- “Rounding Mode: Floor” on page 36-10
- “Rounding Mode: Nearest” on page 36-11
- “Rounding Mode: Round” on page 36-12
- “Rounding Mode: Simplest” on page 36-14
- “Rounding Mode: Zero” on page 36-17
- “Maximize Precision” on page 36-18
- “Net Slope and Net Bias Precision” on page 36-21
- “Detect Fixed-Point Constant Precision Loss” on page 36-24
- “Saturation and Wrapping” on page 36-25
- “Guard Bits” on page 36-28
- “Determine the Range of Fixed-Point Numbers” on page 36-29
- “Handle Overflows in Simulink Models” on page 36-30
- “Recommendations for Arithmetic and Scaling” on page 36-31
- “Parameter and Signal Conversions” on page 36-39
- “Rules for Arithmetic Operations” on page 36-42
- “The Summation Process” on page 36-49
- “The Multiplication Process” on page 36-51
- “The Division Process” on page 36-53
- “Shifts” on page 36-54
- “Conversions and Arithmetic Operations” on page 36-55

Rounding

When you represent numbers with finite precision, not every number in the available range can be represented exactly. The result of any operation on a fixed-point number is typically stored in a register that is longer than the number's original format. When the result is put back into the original format, a rounding method is used to cast the value to a representable number. Precision is always lost in the rounding operation, and produces quantization errors and computational noise. The cost of the rounding operation and the amount of bias that is introduced depends on the rounding method itself.

Choosing a Rounding Method

Each rounding method has a set of inherent properties. Depending on the requirements of your design, these properties could make the rounding method more or less desirable to you. By knowing the requirements of your design and understanding the properties of each rounding method, you can determine which is the best fit for your needs. The most important properties to consider are:

- Cost — Independent of the hardware being used, how much processing expense does the rounding method require?
 - Low — The method requires few processing cycles.
 - Moderate — The method requires a moderate number of processing cycles.
 - High — The method requires more processing cycles.

Note The cost estimates provided here are hardware independent. Some processors have rounding modes built-in, so consider carefully the hardware you are using before calculating the true cost of each rounding mode.

- Bias — What is the expected value of the rounded values minus the original values: $E(\hat{\theta} - \theta)$?
 - $E(\hat{\theta} - \theta) < 0$ — The rounding method introduces a negative bias.
 - $E(\hat{\theta} - \theta) = 0$ — The rounding method is unbiased.
 - $E(\hat{\theta} - \theta) > 0$ — The rounding method introduces a positive bias.

Fixed-Point Designer Rounding Modes

To provide you with greater flexibility in the trade-off between cost and bias, the Fixed-Point Designer product currently supports the following rounding methods:

| Fixed-Point Designer Rounding Mode | Description | Tie Handling | Cost | Bias |
|------------------------------------|---|--------------|------|----------------|
| Ceiling on page 36-8 | Rounds to the nearest representable number in the direction of positive infinity. | N/A | Low | Large positive |

| Fixed-Point Designer Rounding Mode | Description | Tie Handling | Cost | Bias |
|--|--|--|----------|---|
| Convergent on page 36-9 | Rounds to the nearest representable number. | Ties are rounded to nearest even number. | High | Unbiased |
| Floor on page 36-10 | Rounds to the nearest representable number in the direction of negative infinity. Equivalent to two's complement truncation. | N/A | Low | Large negative |
| Nearest on page 36-11 | Rounds to the nearest representable number. | Ties are rounded to the closest representable number in the direction of positive infinity. | Moderate | Small positive |
| Round on page 36-12 | Rounds to the nearest representable number. | <ul style="list-style-type: none"> For positive numbers, ties are rounded to the nearest representable number in the direction of positive infinity. For negative numbers, ties are rounded to the nearest representable number in the direction of negative infinity. | High | <ul style="list-style-type: none"> Small negative for negative samples Unbiased for samples with evenly distributed positive and negative values Small positive for positive samples |
| Simplest on page 36-14 (Simulink only) | Automatically chooses between Floor and Zero to produce generated code that is as efficient as possible. | N/A | Low | Depends on the operation |
| Zero on page 36-17 | Rounds to the nearest representable number in the direction of zero. | N/A | Low | <ul style="list-style-type: none"> Large positive for negative samples Unbiased for samples with evenly distributed positive and negative values Large negative for positive samples |

Choosing a Rounding Mode for Diagnostic Purposes

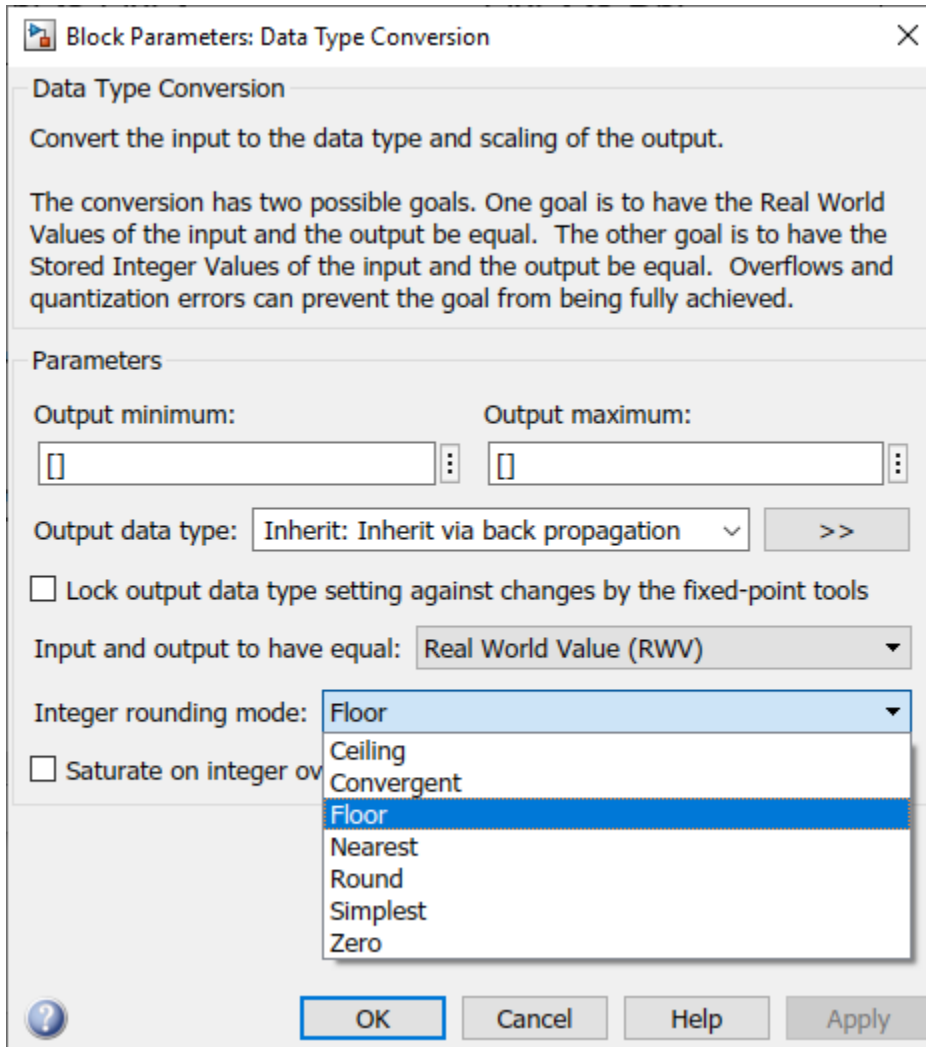
Rounding toward ceiling and rounding toward floor are sometimes useful for diagnostic purposes. For example, after a series of arithmetic operations, you may not know the exact answer because of word-size limitations, which introduce rounding. If every operation in the series is performed twice, once rounding to positive infinity and once rounding to negative infinity, you obtain an upper limit and a lower limit on the correct answer. You can then decide if the result is sufficiently accurate or if additional analysis is necessary.

See Also**More About**

- “Range and Precision” on page 35-9

Rounding Modes for Fixed-Point Simulink Blocks

Fixed-point Simulink blocks support the rounding modes shown in the expanded drop-down menu of the following dialog box.



Fixed-Point Designer Rounding Modes

To provide you with greater flexibility in the trade-off between cost and bias, the Fixed-Point Designer product currently supports the following rounding methods:

| Fixed-Point Designer Rounding Mode | Description | Tie Handling | Cost | Bias |
|---|--|--|-------------|---|
| Ceiling on page 36-8 | Rounds to the nearest representable number in the direction of positive infinity. | N/A | Low | Large positive |
| Convergent on page 36-9 | Rounds to the nearest representable number. | Ties are rounded to nearest even number. | High | Unbiased |
| Floor on page 36-10 | Rounds to the nearest representable number in the direction of negative infinity. Equivalent to two's complement truncation. | N/A | Low | Large negative |
| Nearest on page 36-11 | Rounds to the nearest representable number. | Ties are rounded to the closest representable number in the direction of positive infinity. | Moderate | Small positive |
| Round on page 36-12 | Rounds to the nearest representable number. | <ul style="list-style-type: none"> For positive numbers, ties are rounded to the nearest representable number in the direction of positive infinity. For negative numbers, ties are rounded to the nearest representable number in the direction of negative infinity. | High | <ul style="list-style-type: none"> Small negative for negative samples Unbiased for samples with evenly distributed positive and negative values Small positive for positive samples |
| Simplest on page 36-14 (Simulink only) | Automatically chooses between Floor and Zero to produce generated code that is as efficient as possible. | N/A | Low | Depends on the operation |

| Fixed-Point Designer Rounding Mode | Description | Tie Handling | Cost | Bias |
|------------------------------------|--|--------------|------|---|
| Zero on page 36-17 | Rounds to the nearest representable number in the direction of zero. | N/A | Low | <ul style="list-style-type: none">• Large positive for negative samples• Unbiased for samples with evenly distributed positive and negative values• Large negative for positive samples |

See Also

More About

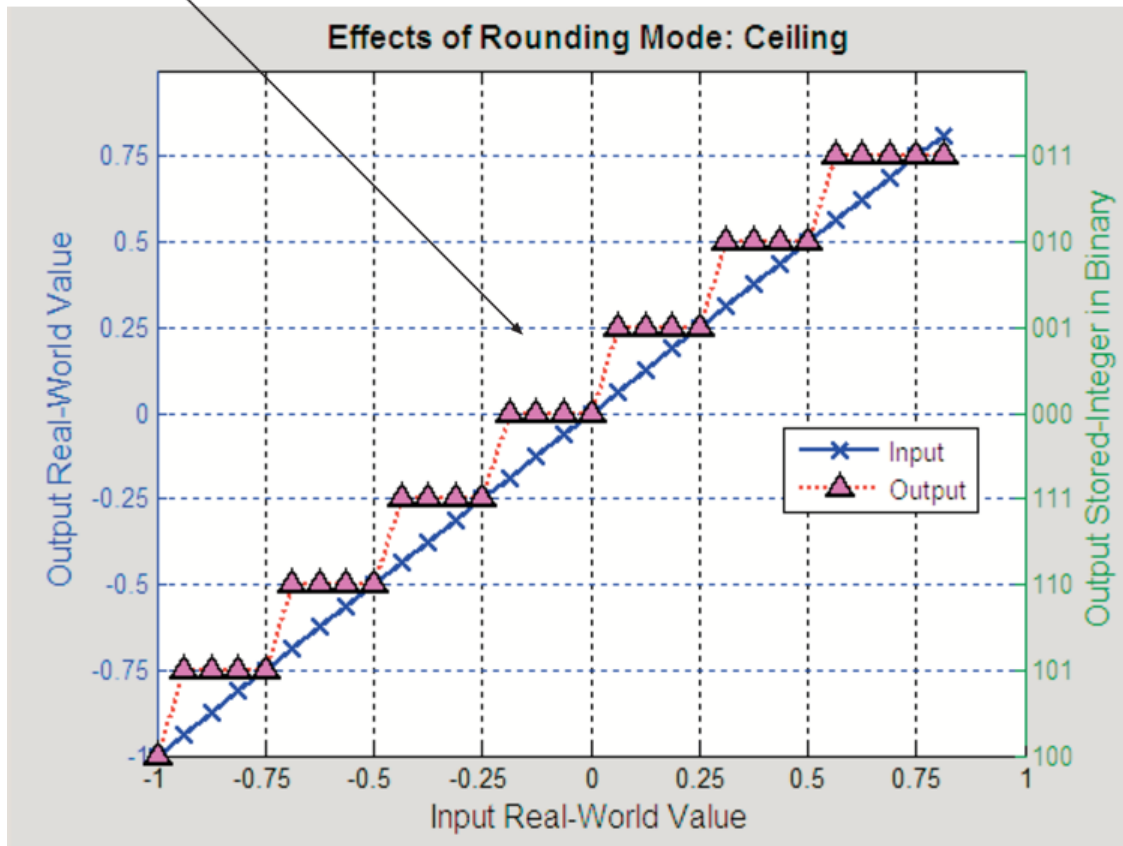
- “Rounding” on page 36-2
- “Range and Precision” on page 35-9

Rounding Mode: Ceiling

When you round toward ceiling, both positive and negative numbers are rounded toward positive infinity. As a result, a positive cumulative bias is introduced in the number.

In the MATLAB software, you can round to ceiling using the `ceil` function. Rounding toward ceiling is shown in the following figure.

All numbers are rounded toward positive infinity

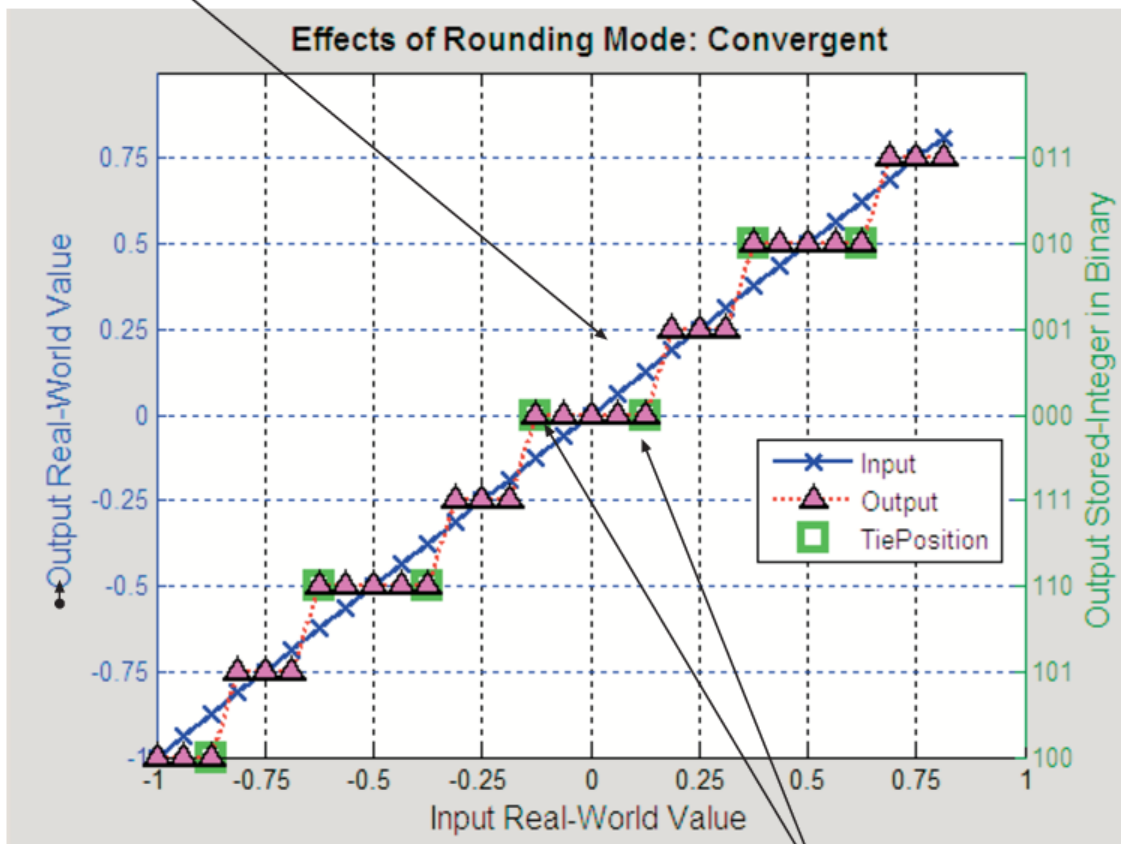


Rounding Mode: Convergent

Convergent rounds toward the nearest representable value with ties rounding toward the nearest even integer. It eliminates bias due to rounding. However, it introduces the possibility of overflow.

In the MATLAB software, you can perform convergent rounding using the `convergent` function. Convergent rounding is shown in the following figure.

All numbers are rounded to the nearest representable number

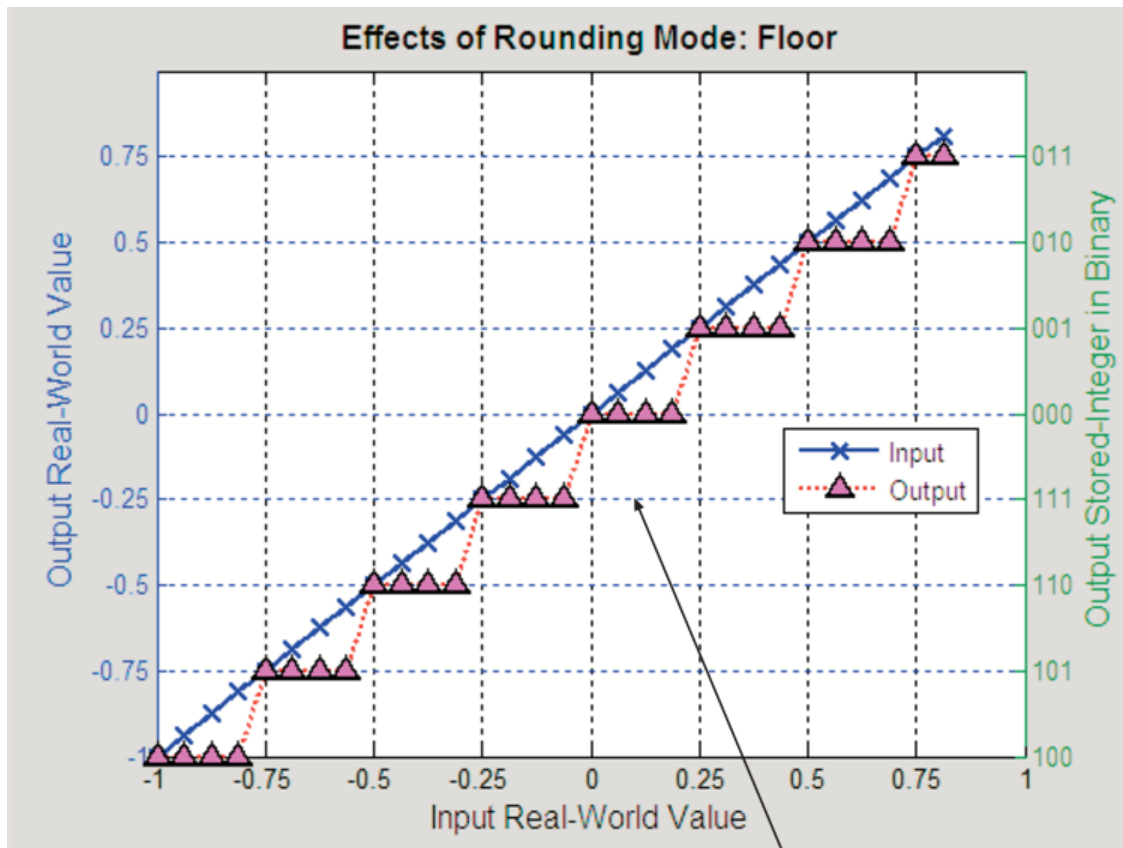


Ties are rounded to the nearest even number

Rounding Mode: Floor

When you round toward floor, both positive and negative numbers are rounded to negative infinity. As a result, a negative cumulative bias is introduced in the number.

In the MATLAB software, you can round to floor using the `floor` function. Rounding toward floor is shown in the following figure.

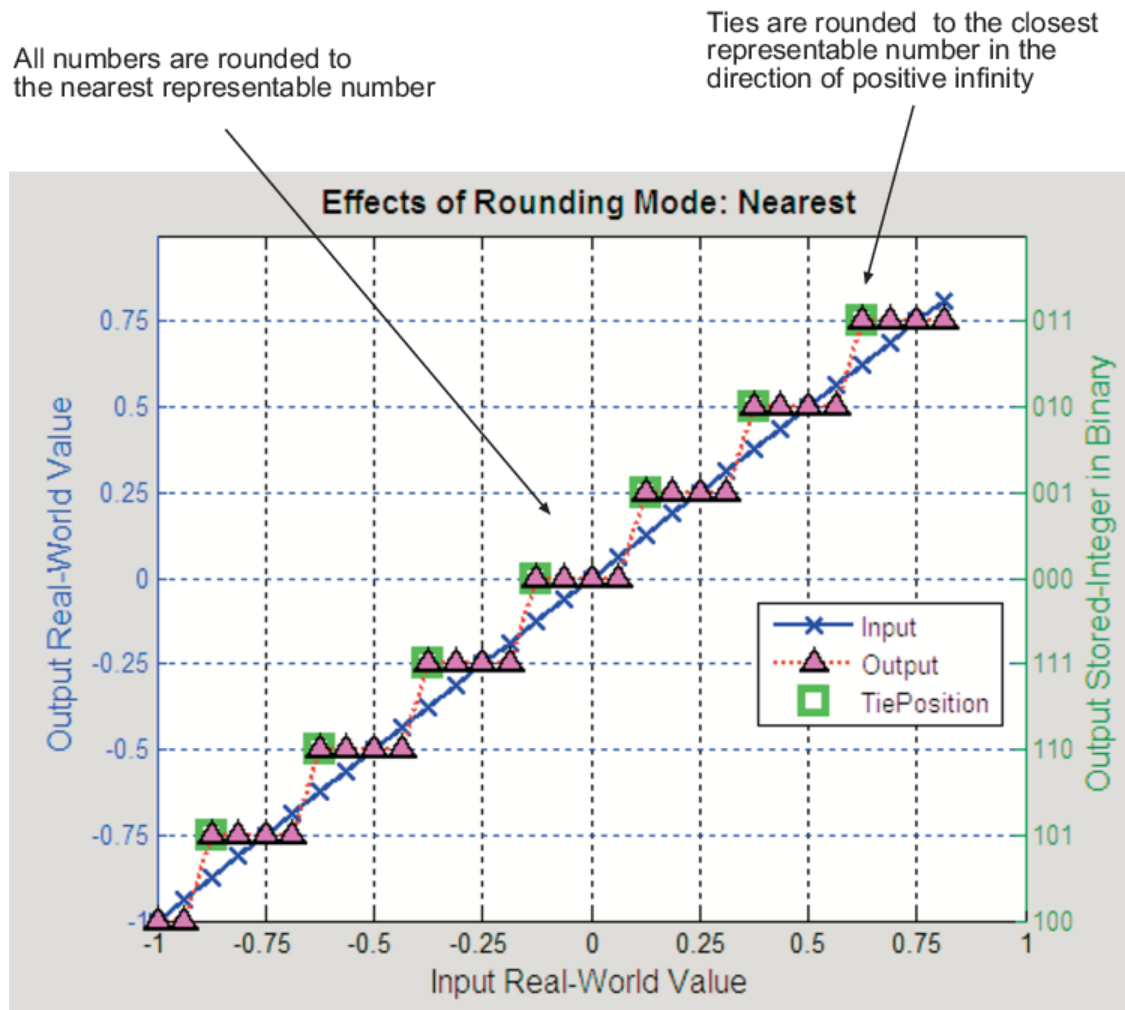


All numbers are rounded toward negative infinity

Rounding Mode: Nearest

When you round toward nearest, the number is rounded to the nearest representable value. In the case of a tie, nearest rounds to the closest representable number in the direction of positive infinity.

In the Fixed-Point Designer software, you can round to nearest using the `nearest` function. Rounding toward nearest is shown in the following figure.



Rounding Mode: Round

Round rounds to the closest representable number. In the case of a tie, it rounds:

- Positive numbers to the closest representable number in the direction of positive infinity.
- Negative numbers to the closest representable number in the direction of negative infinity.

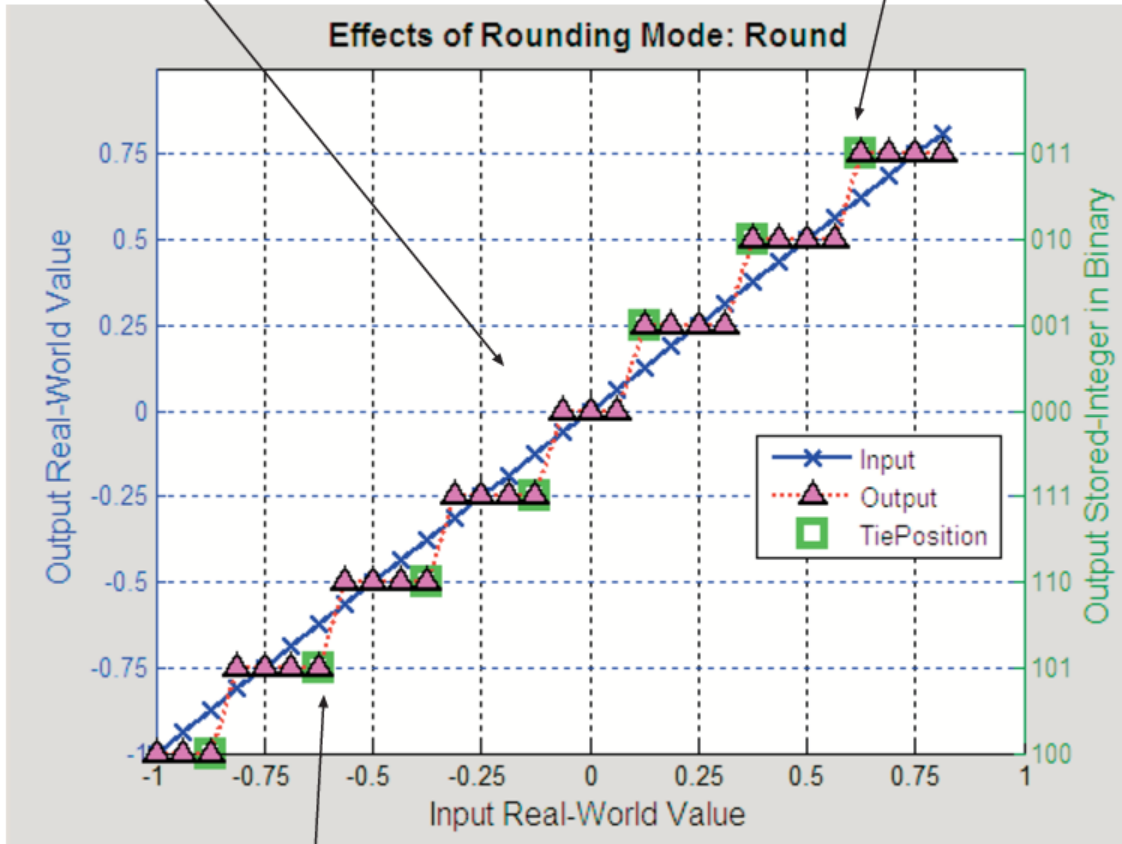
As a result:

- A small negative bias is introduced for negative samples.
- No bias is introduced for samples with evenly distributed positive and negative values.
- A small positive bias is introduced for positive samples.

In the MATLAB software, you can perform this type of rounding using the `round` function. The rounding mode Round is shown in the following figure.

All numbers are rounded to the nearest representable number

For positive numbers, ties are rounded to the closest representable number in the direction of positive infinity



For negative numbers, ties are rounded to the closest representable number in the direction of negative infinity

Rounding Mode: Simplest

The simplest rounding mode attempts to reduce or eliminate the need for extra rounding code in your generated code using a combination of techniques. In nearly all cases, the simplest rounding mode produces the most efficient generated code.

For a very specialized case of division that meets three specific criteria, round to floor might be more efficient. These three criteria are:

- Fixed-point/integer signed division
- Denominator is an invariant constant
- Denominator is an exact power of two

For this case, set the rounding mode to floor and the **Model Configuration Parameters > Hardware Implementation > Production Hardware > Signed integer division rounds to** parameter to describe the rounding behavior of your production target.

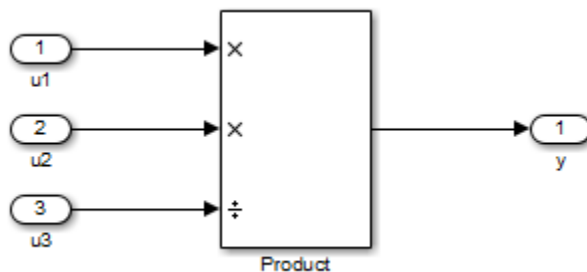
Optimize Rounding for Casts

The Data Type Conversion block casts a signal with one data type to another data type. When the block casts the signal to a data type with a shorter word length than the original data type, precision is lost and rounding occurs. The simplest rounding mode automatically chooses the best rounding for these cases based on the following rules:

- When casting from one integer or fixed-point data type to another, the simplest mode rounds toward floor.
- When casting from a floating-point data type to an integer or fixed-point data type, the simplest mode rounds toward zero.

Optimize Rounding for High-Level Arithmetic Operations

The simplest rounding mode chooses the best rounding for each high-level arithmetic operation. For example, consider the operation $y = u_1 \times u_2 / u_3$ implemented using a Product block:



As stated in the C standard, the most efficient rounding mode for multiplication operations is always floor. However, the C standard does not specify the rounding mode for division in cases where at least one of the operands is negative. Therefore, the most efficient rounding mode for a divide operation with signed data types can be floor or zero, depending on your production target.

The simplest rounding mode:

- Rounds to floor for all non-division operations.
- Rounds to zero or floor for division, depending on the setting of the **Model Configuration Parameters > Hardware Implementation > Production Hardware > Signed integer division rounds to** parameter.

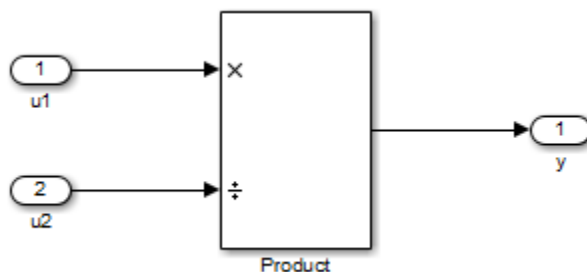
To get the most efficient code, you must set the **Signed integer division rounds to** parameter to specify whether your production target rounds to zero or to floor for integer division. Most production targets round to zero for integer division operations. Note that **Simplest** rounding enables “mixed-mode” rounding for such cases, as it rounds to floor for multiplication and to zero for division.

If the **Signed integer division rounds to** parameter is set to **Undefined**, the simplest rounding mode might not be able to produce the most efficient code. The simplest mode rounds to zero for division for this case, but it cannot rely on your production target to perform the rounding, because the parameter is **Undefined**. Therefore, you need additional rounding code to ensure rounding to zero behavior.

Note For signed fixed-point division where the denominator is an invariant constant power of 2, the simplest rounding mode does not generate the most efficient code. In this case, set the rounding mode to floor.

Optimize Rounding for Intermediate Arithmetic Operations

For fixed-point arithmetic with nonzero slope and bias, the simplest rounding mode also chooses the best rounding for each intermediate arithmetic operation. For example, consider the operation $y = u_1 / u_2$ implemented using a Product block, where u_1 and u_2 are fixed-point quantities:



As discussed in “Data Types and Scaling in Digital Hardware” on page 35-2, each fixed-point quantity is calculated using its slope, bias, and stored integer. In this example, the high-level divide operation specified by the block results in intermediate addition and multiplication operations:

$$y = \frac{u_1}{u_2} = \frac{S_1 Q_1 + B_1}{S_2 Q_2 + B_2}$$

The simplest rounding mode performs the best rounding for each of these operations, high-level and intermediate, to produce the most efficient code. The rules used to select the appropriate rounding for intermediate arithmetic operations are the same as those described in “Optimize Rounding for High-Level Arithmetic Operations” on page 36-14. Again, this enables mixed-mode rounding, with the most common case being round toward floor used for additions, subtractions, and multiplies, and round toward zero used for divides.

Remember that generating the most efficient code using the simplest rounding mode requires you to set the **Model Configuration Parameters > Hardware Implementation > Production Hardware > Signed integer division rounds to** parameter to describe the rounding behavior of your production target.

Note For signed fixed-point division where the denominator is an invariant constant power of 2, the simplest rounding mode does not generate the most efficient code. In this case, set the rounding mode to floor.

See Also

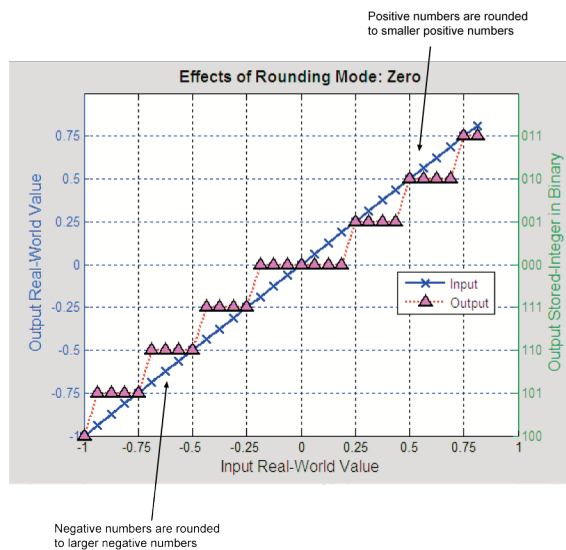
More About

- “Precision and Range” on page 1-6

Rounding Mode: Zero

Rounding towards zero is the simplest rounding mode computationally. All digits beyond the number required are dropped. Rounding towards zero results in a number whose magnitude is always less than or equal to the more precise original value. In the MATLAB software, you can round to zero using the `fix` function.

Rounding toward zero introduces a cumulative downward bias in the result for positive numbers and a cumulative upward bias in the result for negative numbers. That is, all positive numbers are rounded to smaller positive numbers, while all negative numbers are rounded to larger negative numbers. Rounding toward zero is shown in the following figure.



Rounding to Zero Versus Truncation

Rounding to zero and *truncation* or *chopping* are sometimes thought to mean the same thing. However, the results produced by rounding to zero and truncation are different for unsigned and two's complement numbers. For this reason, the ambiguous term "truncation" is not used in this guide, and explicit rounding modes are used instead.

To illustrate this point, consider rounding a 5-bit unsigned number to zero by dropping (truncating) the two least significant bits. For example, the unsigned number $100.01 = 4.25$ is truncated to $100 = 4$. Therefore, truncating an unsigned number is equivalent to rounding to zero *or* rounding to floor.

Now consider rounding a 5-bit two's complement number by dropping the two least significant bits. At first glance, you may think truncating a two's complement number is the same as rounding to zero. For example, dropping the last two digits of -3.75 yields -3.00 . However, digital hardware performing two's complement arithmetic yields a different result. Specifically, the number $100.01 = -3.75$ truncates to $100 = -4$, which is rounding to floor.

Maximize Precision

Precision is limited by slope. To achieve maximum precision, you should make the slope as small as possible while keeping the range adequately large. The bias is adjusted in coordination with the slope.

Assume the maximum and minimum real-world values are given by $\max(V)$ and $\min(V)$, respectively. These limits might be known based on physical principles or engineering considerations. To maximize the precision, you must decide upon a rounding scheme and whether overflows saturate or wrap. To simplify matters, this example assumes the minimum real-world value corresponds to the minimum encoded value, and the maximum real-world value corresponds to the maximum encoded value. Using the encoding scheme described in “Scaling” on page 35-5, these values are given by

$$\begin{aligned}\max(V) &= F2^E(\max(Q)) + B \\ \min(V) &= F2^E(\min(Q)) + B.\end{aligned}$$

Solving for the slope, you get

$$F2^E = \frac{\max(V) - \min(V)}{\max(Q) - \min(Q)} = \frac{\max(V) - \min(V)}{2^{ws} - 1}.$$

This formula is independent of rounding and overflow issues, and depends only on the word size, ws .

Pad with Trailing Zeros

Padding with trailing zeros involves extending the least significant bit (LSB) of a number with extra bits. This method involves going from low precision to higher precision.

For example, suppose two numbers are subtracted from each other. First, the exponents must be aligned, which typically involves a right shift of the number with the smaller value. In performing this shift, significant digits can “fall off” to the right. However, when the appropriate number of extra bits is appended, the precision of the result is maximized. Consider two 8-bit fixed-point numbers that are close in value and subtracted from each other:

$$1.0000000 \times 2^q - 1.1111111 \times 2^{q-1},$$

where q is an integer. To perform this operation, the exponents must be equal:

$$\begin{array}{r} 1.0000000 \times 2^q \\ -0.1111111 \times 2^q \\ \hline 0.0000001 \times 2^q \end{array}.$$

If the top number is padded by two zeros and the bottom number is padded with one zero, then the above equation becomes

$$\begin{array}{r} 1.00000000 \times 2^q \\ -0.11111110 \times 2^q \\ \hline 0.00000010 \times 2^q \end{array},$$

which produces a more precise result. An example of padding with trailing zeros in a Simulink model is illustrated in “Digital Controller Realization” on page 42-41.

Constant Scaling for Best Precision

The following fixed-point Simulink blocks provide a mode for scaling parameters whose values are constant vectors or matrices:

- Constant
- Discrete FIR Filter
- Gain
- Relay
- Repeating Sequence Stair

This scaling mode is based on binary-point-only scaling. Using this mode, you can scale a constant vector or matrix such that a common binary point is found based on the best precision for the largest value in the vector or matrix.

Constant scaling for best precision is available only for fixed-point data types with unspecified scaling. All other fixed-point data types use their specified scaling. You can use the **Data Type Assistant** (see “Specify Data Types Using Data Type Assistant”) on a block dialog box to enable the best precision scaling mode.

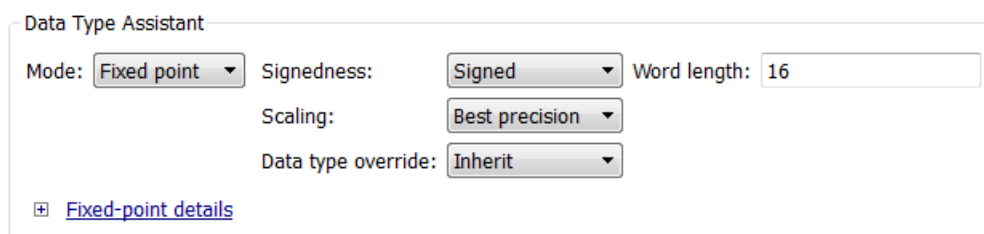
- 1 On a block dialog box, click the **Show data type assistant** button  .

The **Data Type Assistant** appears.

- 2 In the **Data Type Assistant**, and from the **Mode** list, select **Fixed point**.

The **Data Type Assistant** displays additional options associated with fixed-point data types.

- 3 From the **Scaling** list, select **Best precision**.



To understand how you might use this scaling mode, consider a 3-by-3 matrix of doubles, M , defined as

```
3.3333e-003  3.3333e-004  3.3333e-005
3.3333e-002  3.3333e-003  3.3333e-004
3.3333e-001  3.3333e-002  3.3333e-003
```

Now suppose you specify M as the value of the **Gain** parameter for a Gain block. The results of specifying your own scaling versus using the constant scaling mode are described here:

- **Specified Scaling**

Suppose the matrix elements are converted to a signed, 10-bit generalized fixed-point data type with binary-point-only scaling of 2^{-7} (that is, the binary point is located seven places to the left of the right most bit). With this data format, M becomes

| | | |
|-------------|-------------|---|
| 0 | 0 | 0 |
| 3.1250e-002 | 0 | 0 |
| 3.3594e-001 | 3.1250e-002 | 0 |

Note that many of the matrix elements are zero, and for the nonzero entries, the scaled values differ from the original values. This is because a double is converted to a binary word of fixed size and limited precision for each element. The larger and more precise the conversion data type, the more closely the scaled values match the original values.

- **Constant Scaling for Best Precision**

If M is scaled based on its largest matrix value, you obtain

| | | |
|-------------|-------------|-------------|
| 2.9297e-003 | 0 | 0 |
| 3.3203e-002 | 2.9297e-003 | 0 |
| 3.3301e-001 | 3.3203e-002 | 2.9297e-003 |

Best precision would automatically select the fraction length that minimizes the quantization error. Even though precision was maximized for the given word length, quantization errors can still occur. In this example, a few elements still quantize to zero.

See Also

More About

- “Range and Precision” on page 35-9
- “Detect Fixed-Point Constant Precision Loss” on page 36-24

Net Slope and Net Bias Precision

What are Net Slope and Net Bias?

You can represent a fixed-point number by a general slope and bias encoding scheme,

$$V \approx \tilde{V} = SQ + B,$$

where:

- V is an arbitrarily precise real-world value.
- \tilde{V} is the approximate real-world value.
- Q , the stored value, is an integer that encodes V .
- $S = F2^E$ is the slope.
- B is the bias.

For a cast operation,

$$S_a Q_a + B_a = S_b Q_b + B_b$$

or

$$Q_a = \frac{S_b Q_b}{S_a} + \left(\frac{B_b - B_a}{S_a} \right),$$

where:

- $\frac{S_b}{S_a}$ is the net slope.
- $\frac{B_b - B_a}{S_a}$ is the net bias.

Detect Net Slope and Net Bias Precision Issues

Precision issues might occur in the fixed-point constants, net slope and net bias, due to quantization errors when you convert from floating point to fixed point. These fixed-point constant precision issues can result in numerical inaccuracy in your model.

You can configure your model to alert you when fixed-point constant precision issues occur.

You can configure your model to alert you when fixed-point constant precision issues occur. To receive alerts when fixed-point constant precision issues occur, use these options available in the Simulink Configuration Parameters dialog box, on the **Diagnostics > Type Conversion** pane. Set the parameters to **warning** or **error** so that Simulink alerts you when precision issues occur.

| Configuration Parameter | Specifies | Default |
|-------------------------|---|---------------------------------------|
| "Detect underflow" | Diagnostic action when a fixed-point constant underflow occurs during simulation | Does not generate a warning or error. |
| "Detect overflow" | Diagnostic action when a fixed-point constant overflow occurs during simulation | Does not generate a warning or error. |
| "Detect precision loss" | Diagnostic action when a fixed-point constant precision loss occurs during simulation | Does not generate a warning or error. |

The Fixed-Point Designer software provides the following information:

- The type of precision issue: underflow, overflow, or precision loss.
- The original value of the fixed-point constant.
- The quantized value of the fixed-point constant.
- The error in the value of the fixed-point constant.
- The block that introduced the error.

This information warns you that the outputs from this block are not accurate. If possible, change the data types in your model to fix the issue.

Fixed-Point Constant Underflow

Fixed-point constant underflow occurs when the Fixed-Point Designer software encounters a fixed-point constant whose data type does not have enough precision to represent the ideal value of the constant, because the ideal value is too close to zero. Casting the ideal value to the fixed-point data type causes the value of the fixed-point constant to become zero. Therefore the value of the fixed-point constant differs from its ideal value.

Fixed-Point Constant Overflow

Fixed-point constant overflow occurs when the Fixed-Point Designer software converts a fixed-point constant to a data type whose range is not large enough to accommodate the ideal value of the constant with reasonable precision. The data type cannot accurately represent the ideal value because the ideal value is either too large or too small. Casting the ideal value to the fixed-point data type causes overflow. For example, suppose the ideal value is 200 and the converted data type is `int8`. Overflow occurs in this case because the maximum value that `int8` can represent is 127.

The Fixed-Point Designer software reports an overflow error if the quantized value differs from the ideal value by more than the precision for the data type. The precision for a data type is approximately equal to the default scaling (for more information, see "Fixed-Point Data Type Parameters" on page 35-11.) Therefore, for positive values, the Fixed-Point Designer software treats errors greater than the slope as overflows. For negative values, it treats errors greater than or equal to the slope as overflows.

For example, the maximum value that `int8` can represent is 127. The precision for `int8` is 1.0. An ideal value of 127.3 quantizes to 127 with an absolute error of 0.3. Although the ideal value 127.3 is greater than the maximum representable value for `int8`, the quantization error is small relative to the precision of `int8`. Therefore the Fixed-Point Designer software does not report an overflow.

However, an ideal value of 128.1 does cause an overflow because the quantization error is 1.1, which is larger than the precision for int8.

Note Fixed-point constant overflow differs from fixed-point constant precision loss. Precision loss occurs when the ideal fixed-point constant value is within the range of the current data type and scaling, but the software cannot represent this value exactly.

Fixed-Point Constant Precision Loss

Fixed-point constant precision loss occurs when the Fixed-Point Designer software converts a fixed-point constant to a data type without enough precision to represent the exact value of the constant. As a result, the quantized value differs from the ideal value. For an example of this behavior, see “Detect Fixed-Point Constant Precision Loss” on page 36-24.

Note Fixed-point constant precision loss differs from fixed-point constant overflow. Overflow occurs when the range of the parameter data type, that is, the maximum value that it can represent, is smaller than the ideal value of the parameter.

See Also

More About

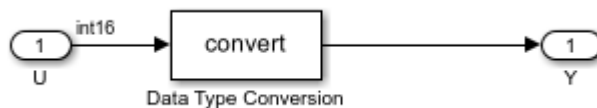
- “Detect Fixed-Point Constant Precision Loss” on page 36-24

Detect Fixed-Point Constant Precision Loss

This example shows how to detect fixed-point constant precision loss.

Open the Model

```
open_system('ex_fixed_point_constant_precision_loss');
```



Detect Fixed-Point Constant Precision Loss

For the Data Type Conversion block in this model:

- Input slope, $S_U = 1$
- Output slope, $S_Y = 1.000001$
- Net slope, $S_U/S_Y = 1/1.000001$

To set up the model and run the simulation:

- 1 For the Inport block, set the **Data type** to int16.
- 2 For the Data Type Conversion block, set the **Output data type** to fixdt(1,16,1.000001,0).
- 3 In the **Configuration Parameters** dialog box, set the **Diagnostics > Type Conversion > Detect precision loss** configuration parameter to error.
- 4 In your Simulink model window, in the **Simulation** tab, click **Run**.

When you simulate the model, a net slope quantization error occurs.

The Fixed-Point Designer software generates an error informing you that net scaling quantization caused precision loss. The message provides the following information:

- The block that introduced the error.
- The original value of the net slope.
- The quantized value of the net slope.
- The error in the value of the net slope.

See Also

More About

- “Net Slope and Net Bias Precision” on page 36-21
- “Range and Precision” on page 35-9
- “Fixed-Point Numbers in Simulink” on page 35-13

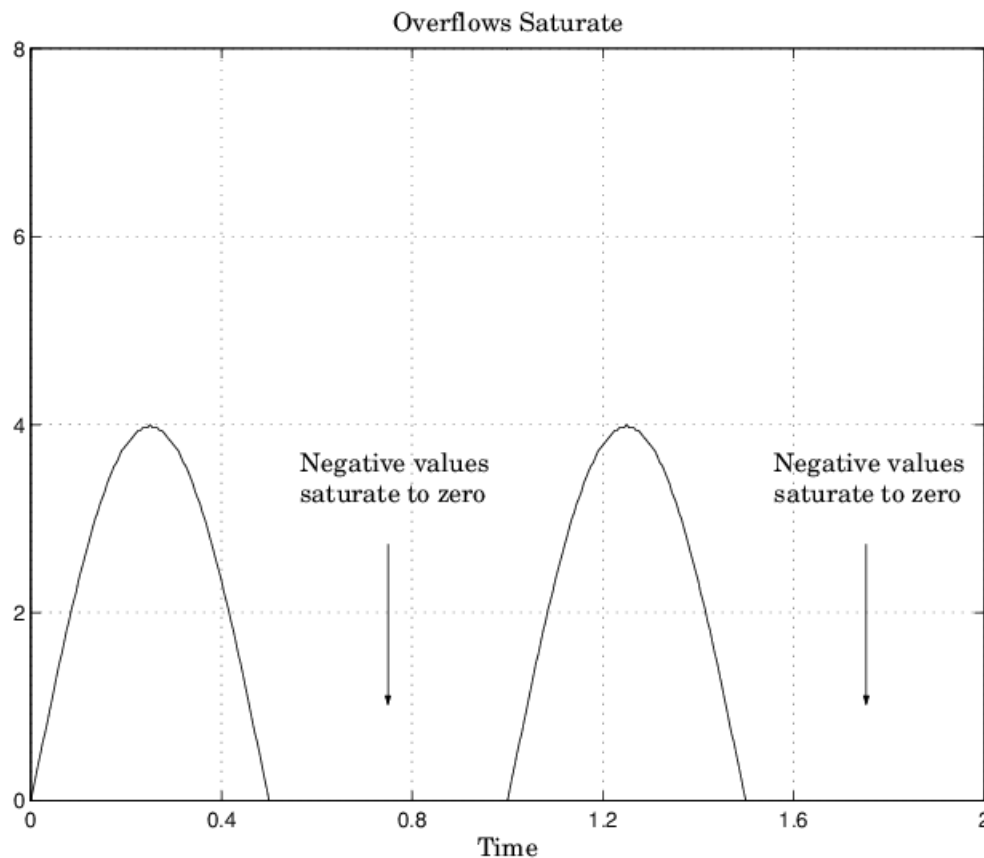
Saturation and Wrapping

What Are Saturation and Wrapping?

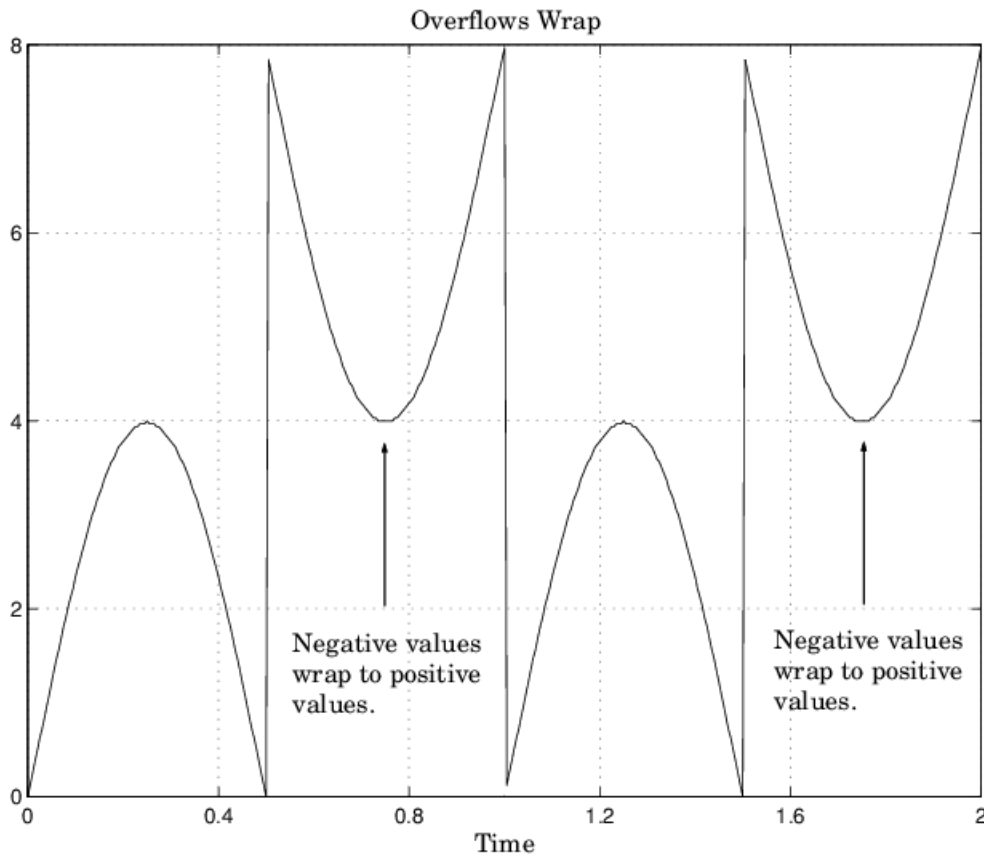
Saturation and wrapping describe a particular way that some processors deal with overflow conditions. For example, the ADSP-2100 family of processors from Analog Devices® supports either of these modes. If a register has a saturation mode of operation, then an overflow condition is set to the maximum positive or negative value allowed. Conversely, if a register has a wrapping mode of operation, an overflow condition is set to the appropriate value within the range of the representation.

Saturation and Wrapping

Consider an 8-bit unsigned word with binary-point-only scaling of 2^{-5} . Suppose this data type must represent a sine wave that ranges from -4 to 4. For values between 0 and 4, the word can represent these numbers without regard to overflow. This is not the case with negative numbers. If overflows saturate, all negative values are set to zero, which is the smallest number representable by the data type. The saturation of overflows is shown in the following figure.

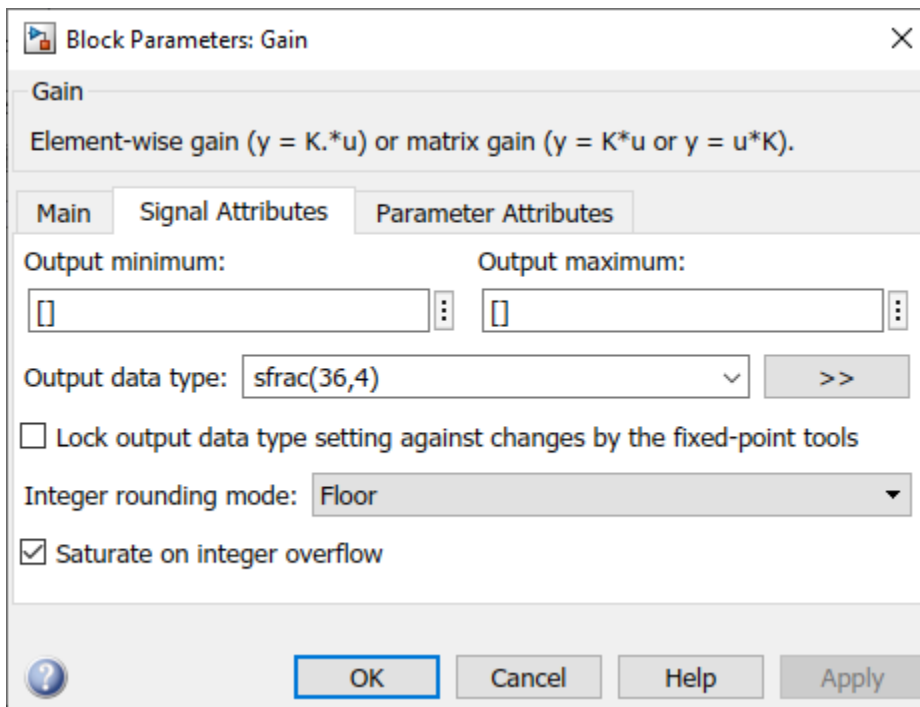


If overflows wrap, all negative values are set to the appropriate positive value. The wrapping of overflows is shown in the following figure.



Note For most control applications, saturation is the safer way of dealing with fixed-point overflow. However, some processor architectures allow automatic saturation by hardware. If hardware saturation is not available, then extra software is required, resulting in larger, slower programs. This cost is justified in some designs — perhaps for safety reasons. Other designs accept wrapping to obtain the smallest, fastest software.

The Simulink software supports saturation and wrapping for all fixed-point data types. You can select saturation or wrapping for fixed-point Simulink blocks with the **Saturate on integer overflow** check box.



See Also

More About

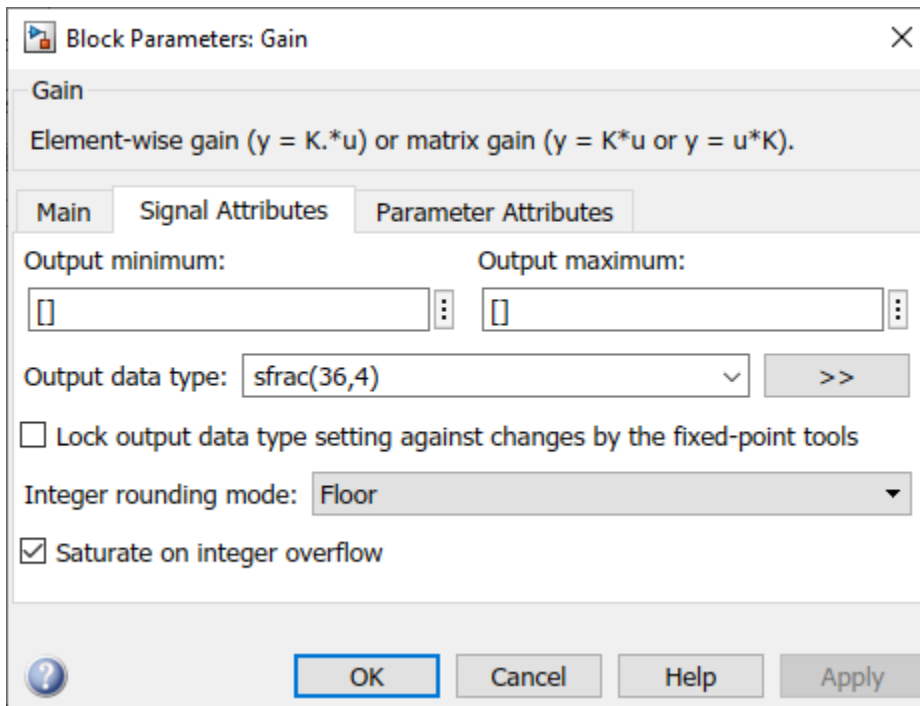
- “Range and Precision” on page 35-9

Guard Bits

You can eliminate the possibility of overflow by appending the appropriate number of guard bits to a binary word.

For a two's complement signed value, the guard bits are filled with either 0's or 1's depending on the value of the most significant bit (MSB). This is called *sign extension*. For example, consider a 4-bit two's complement number with value 1011. If this number is extended in range to 7 bits with sign extension, then the number becomes 1111101 and the value remains the same.

The Simulink software supports guard bits only for fractional data types. For both signed and unsigned fractionals, the guard bits lie to the left of the default binary point. For example, by setting **Output data type** to `sfrac(36,4)`, you specify a 36-bit signed fractional data type with 4 guard bits (total word size is 40 bits).



Determine the Range of Fixed-Point Numbers

Fixed-point variables have a limited range for the same reason they have limited precision — because digital systems represent numbers with a finite number of bits. As a general example, consider the case where an integer is represented as a fixed-point word of size ws . The range for signed and unsigned words is given by

$$\max(Q) - \min(Q),$$

where

$$\min(Q) = \begin{cases} 0 & \text{unsigned,} \\ -2^{ws-1} & \text{signed,} \end{cases}$$

$$\max(Q) = \begin{cases} 2^{ws} - 1 & \text{unsigned,} \\ 2^{ws-1} - 1 & \text{signed.} \end{cases}$$

Using the general [Slope Bias] encoding scheme described in “Scaling” on page 35-5, the approximate real-world value has the range

$$\max(\tilde{V}) - \min(\tilde{V}),$$

where

$$\min(\tilde{V}) = \begin{cases} B & \text{unsigned,} \\ -F2^E(2^{ws-1}) + B & \text{signed,} \end{cases}$$

$$\max(\tilde{V}) = \begin{cases} F2^E(2^{ws} - 1) + B & \text{unsigned,} \\ F2^E(2^{ws-1} - 1) + B & \text{signed.} \end{cases}$$

If the real-world value exceeds the limited range of the approximate value, then the accuracy of the representation can become significantly worse.

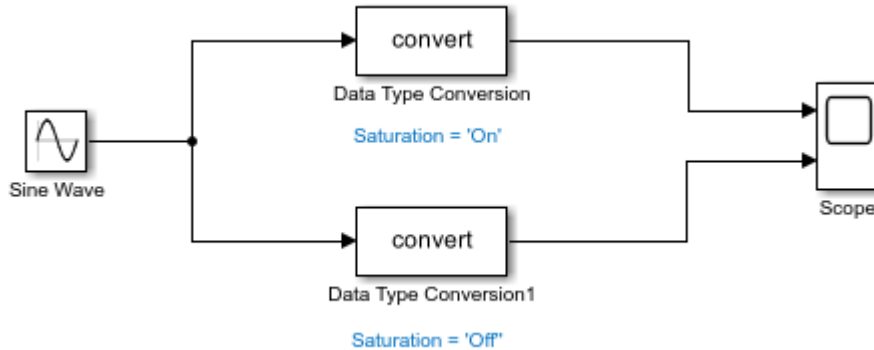
Handle Overflows in Simulink Models

This example shows how to control the warning messages you receive when a model contains an overflow. This diagnostic control can simplify debugging models in which only one type of overflow is of interest.

Open the Model

This model contains a sine wave with an amplitude of 1.5 passed through two Data Type Conversion blocks. In the Data Type Conversion block, the **Saturate on integer overflow** parameter is selected. The Data Type Conversion1 block wraps when the signal is too large to fit into the output data type.

```
open_system('ex_detect_overflows')
```



Simulate Model with Original Diagnostic Settings

Simulate the model.

The Diagnostic Viewer displays two overflow warnings. The first overflow saturated and the second overflow wrapped.

Adjust Diagnostic Settings

In the Configuration Parameters dialog box:

- Set **Diagnostics > Data Validity > Wrap on overflow** to Error.
- Set **Diagnostics > Data Validity > Saturate on overflow** to Warning.

Simulate the model again.

The Diagnostic Viewer displays an error message for the overflow that wrapped, and a warning message for the overflow that saturated.

See Also

“Wrap on overflow” | “Saturate on overflow”

Recommendations for Arithmetic and Scaling

In this section...

“Arithmetic Operations and Fixed-Point Scaling” on page 36-31

“Addition” on page 36-31

“Accumulation” on page 36-33

“Multiplication” on page 36-34

“Gain” on page 36-35

“Division” on page 36-36

“Summary” on page 36-37

Arithmetic Operations and Fixed-Point Scaling

The sections that follow describe the relationship between arithmetic operations and fixed-point scaling, and offer some basic recommendations that may be appropriate for your fixed-point design. For each arithmetic operation,

- The general [Slope Bias] encoding scheme described in “Scaling” on page 35-5 is used.
- The scaling of the result is automatically selected based on the scaling of the two inputs. In other words, the scaling is *inherited*.
- Scaling choices are based on
 - Minimizing the number of arithmetic operations of the result
 - Maximizing the precision of the result

Additionally, binary-point-only scaling is presented as a special case of the general encoding scheme.

In embedded systems, the scaling of variables at the hardware interface (the ADC or DAC) is fixed. However for most other variables, the scaling is something you can choose to give the best design. When scaling fixed-point variables, it is important to remember that

- Your scaling choices depend on the particular design you are simulating.
- There is no best scaling approach. All choices have associated advantages and disadvantages. It is the goal of this section to expose these advantages and disadvantages to you.

Addition

Consider the addition of two real-world values:

$$V_a = V_b + V_c.$$

These values are represented by the general [Slope Bias] encoding scheme described in “Scaling” on page 35-5:

$$V_i = F_i 2^{E_i} Q_i + B_i.$$

In a fixed-point system, the addition of values results in finding the variable Q_a :

$$Q_a = \frac{F_b}{F_a} 2^{E_b - E_a} Q_b + \frac{F_c}{F_a} 2^{E_c - E_a} Q_c + \frac{B_b + B_c - B_a}{F_a} 2^{-E_a}.$$

This formula shows

- In general, Q_a is not computed through a simple addition of Q_b and Q_c .
- In general, there are two multiplications of a constant and a variable, two additions, and some additional bit shifting.

Inherited Scaling for Speed

In the process of finding the scaling of the sum, one reasonable goal is to simplify the calculations. Simplifying the calculations should reduce the number of operations, thereby increasing execution speed. The following choices can help to minimize the number of arithmetic operations:

- Set $B_a = B_b + B_c$. This eliminates one addition.
- Set $F_a = F_b$ or $F_a = F_c$. Either choice eliminates one of the two constant times variable multiplications.

The resulting formula is

$$Q_a = 2^{E_b - E_a} Q_b + \frac{F_c}{F_a} 2^{E_c - E_a} Q_c$$

or

$$Q_a = \frac{F_b}{F_a} 2^{E_b - E_a} Q_b + 2^{E_c - E_a} Q_c.$$

These equations appear to be equivalent. However, your choice of rounding and precision may make one choice stand out over the other. To further simplify matters, you could choose $E_a = E_c$ or $E_a = E_b$. This will eliminate some bit shifting.

Inherited Scaling for Maximum Precision

In the process of finding the scaling of the sum, one reasonable goal is maximum precision. You can determine the maximum-precision scaling if the range of the variable is known. “Maximize Precision” on page 36-18 shows that you can determine the range of a fixed-point operation from $\max(V_a)$ and $\min(V_a)$. For a summation, you can determine the range from

$$\begin{aligned} \min(\tilde{V}_a) &= \min(\tilde{V}_b) + \min(\tilde{V}_c), \\ \max(\tilde{V}_a) &= \max(\tilde{V}_b) + \max(\tilde{V}_c). \end{aligned}$$

You can now derive the maximum-precision slope:

$$\begin{aligned} F_a 2^{E_a} &= \frac{\max(\tilde{V}_a) - \min(\tilde{V}_a)}{2^{ws_a} - 1} \\ &= \frac{F_b 2^{E_b} (2^{ws_b} - 1) + F_c 2^{E_c} (2^{ws_c} - 1)}{2^{ws_a} - 1}. \end{aligned}$$

In most cases the input and output word sizes are much greater than one, and the slope becomes

$$F_a 2^{E_a} \approx F_b 2^{E_b + w_{sb} - w_{sa}} + F_c 2^{E_c + w_{sc} - w_{sa}},$$

which depends only on the size of the input and output words. The corresponding bias is

$$B_a = \min(\tilde{V}_a) - F_a 2^{E_a} \min(Q_a).$$

The value of the bias depends on whether the inputs and output are signed or unsigned numbers.

If the inputs and output are all unsigned, then the minimum values for these variables are all zero and the bias reduces to a particularly simple form:

$$B_a = B_b + B_c.$$

If the inputs and the output are all signed, then the bias becomes

$$B_a \approx B_b + B_c + F_b 2^{E_b} (-2^{w_{sb}-1} + 2^{w_{sb}-1}) + F_c 2^{E_c} (-2^{w_{sc}-1} + 2^{w_{sc}-1}),$$

$$B_a \approx B_b + B_c.$$

Binary-Point-Only Scaling

For binary-point-only scaling, finding Q_a results in this simple expression:

$$Q_a = 2^{E_b - E_a} Q_b + 2^{E_c - E_a} Q_c.$$

This scaling choice results in only one addition and some bit shifting. The avoidance of any multiplications is a big advantage of binary-point-only scaling.

Note The subtraction of values produces results that are analogous to those produced by the addition of values.

Accumulation

The accumulation of values is closely associated with addition:

$$V_{a_new} = V_{a_old} + V_b.$$

Finding Q_{a_new} involves one multiplication of a constant and a variable, two additions, and some bit shifting:

$$Q_{a_new} = Q_{a_old} + \frac{F_b}{F_a} 2^{E_b - E_a} Q_b + \frac{B_b}{F_a} 2^{-E_a}.$$

The important difference for fixed-point implementations is that the scaling of the output is identical to the scaling of the first input.

Binary-Point-Only Scaling

For binary-point-only scaling, finding Q_{a_new} results in this simple expression:

$$Q_{a_new} = Q_{a_old} + 2^{E_b - E_a} Q_b.$$

This scaling option only involves one addition and some bit shifting.

Note The negative accumulation of values produces results that are analogous to those produced by the accumulation of values.

Multiplication

Consider the multiplication of two real-world values:

$$V_a = V_b V_c.$$

These values are represented by the general [Slope Bias] encoding scheme described in “Scaling” on page 35-5:

$$V_i = F_i 2^{E_i} Q_i + B_i.$$

In a fixed-point system, the multiplication of values results in finding the variable Q_a :

$$Q_a = \frac{F_b F_c}{F_a} 2^{E_b + E_c - E_a} Q_b Q_c + \frac{F_b B_c}{F_a} 2^{E_b - E_a} Q_b + \frac{F_c B_b}{F_a} 2^{E_c - E_a} Q_c + \frac{B_b B_c - B_a}{F_a} 2^{-E_a}.$$

This formula shows

- In general, Q_a is not computed through a simple multiplication of Q_b and Q_c .
- In general, there is one multiplication of a constant and two variables, two multiplications of a constant and a variable, three additions, and some additional bit shifting.

Inherited Scaling for Speed

The number of arithmetic operations can be reduced with these choices:

- Set $B_a = B_b B_c$. This eliminates one addition operation.
- Set $F_a = F_b F_c$. This simplifies the triple multiplication—certainly the most difficult part of the equation to implement.
- Set $E_a = E_b + E_c$. This eliminates some of the bit shifting.

The resulting formula is

$$Q_a = Q_b Q_c + \frac{B_c}{F_c} 2^{-E_c} Q_b + \frac{B_b}{F_b} 2^{-E_b} Q_c.$$

Inherited Scaling for Maximum Precision

You can determine the maximum-precision scaling if the range of the variable is known. “Maximize Precision” on page 36-18 shows that you can determine the range of a fixed-point operation from

$$\max(\tilde{V}_a)$$

and

$$\min(\tilde{V}_a).$$

For multiplication, you can determine the range from

$$\begin{aligned}\min(\tilde{V}_a) &= \min(V_{LL}, V_{LH}, V_{HL}, V_{HH}), \\ \max(\tilde{V}_a) &= \max(V_{LL}, V_{LH}, V_{HL}, V_{HH}),\end{aligned}$$

where

$$\begin{aligned}V_{LL} &= \min(\tilde{V}_b) \cdot \min(\tilde{V}_c), \\ V_{LH} &= \min(\tilde{V}_b) \cdot \max(\tilde{V}_c), \\ V_{HL} &= \max(\tilde{V}_b) \cdot \min(\tilde{V}_c), \\ V_{HH} &= \max(\tilde{V}_b) \cdot \max(\tilde{V}_c).\end{aligned}$$

Binary-Point-Only Scaling

For binary-point-only scaling, finding Q_a results in this simple expression:

$$Q_a = 2^{E_b + E_c - E_a} Q_b Q_c.$$

Gain

Consider the multiplication of a constant and a variable

$$V_a = K V_b,$$

where K is a constant called the gain. Since V_a results from the multiplication of a constant and a variable, finding Q_a is a simplified version of the general fixed-point multiplication formula:

$$Q_a = \left(\frac{K F_b 2^{E_b}}{F_a 2^{E_a}} \right) Q_b + \left(\frac{K B_b - B_a}{F_a 2^{E_a}} \right).$$

Note that the terms in the parentheses can be calculated offline. Therefore, there is only one multiplication of a constant and a variable and one addition.

To implement the above equation without changing it to a more complicated form, the constants need to be encoded using a binary-point-only format. For each of these constants, the range is the trivial case of only one value. Despite the trivial range, the binary point formulas for maximum precision are still valid. The maximum-precision representations are the most useful choices unless there is an overriding need to avoid any shifting. The encoding of the constants is

$$\begin{aligned}\left(\frac{K F_b 2^{E_b}}{F_a 2^{E_a}} \right) &= 2^{E_X} Q_X \\ \left(\frac{K B_b - B_a}{F_a 2^{E_a}} \right) &= 2^{E_Y} Q_Y\end{aligned}$$

resulting in the formula

$$Q_a = 2^{E_X} Q_X Q_B + 2^{E_Y} Q_Y.$$

Inherited Scaling for Speed

The number of arithmetic operations can be reduced with these choices:

- Set $B_a = KB_b$. This eliminates one constant term.
- Set $F_a = KF_b$ and $E_a = E_b$. This sets the other constant term to unity.

The resulting formula is simply

$$Q_a = Q_b.$$

If the number of bits is different, then either handling potential overflows or performing sign extensions is the only possible operation involved.

Inherited Scaling for Maximum Precision

The scaling for maximum precision does not need to be different from the scaling for speed unless the output has fewer bits than the input. If this is the case, then saturation should be avoided by dividing the slope by 2 for each lost bit. This prevents saturation but causes rounding to occur.

Division

Division of values is an operation that should be avoided in fixed-point embedded systems, but it can occur in places. Therefore, consider the division of two real-world values:

$$V_a = V_b/V_c.$$

These values are represented by the general [Slope Bias] encoding scheme described in "Scaling" on page 35-5:

$$V_i = F_i 2^{E_i} Q_i + B_i.$$

In a fixed-point system, the division of values results in finding the variable Q_a :

$$Q_a = \frac{F_b 2^{E_b} Q_b + B_b}{F_c F_a 2^{E_c + E_a} Q_c + B_c F_a 2^{E_a}} - \frac{B_a}{F_a} 2^{-E_a}.$$

This formula shows

- In general, Q_a is not computed through a simple division of Q_b by Q_c .
- In general, there are two multiplications of a constant and a variable, two additions, one division of a variable by a variable, one division of a constant by a variable, and some additional bit shifting.

Inherited Scaling for Speed

The number of arithmetic operations can be reduced with these choices:

- Set $B_a = 0$. This eliminates one addition operation.
- If $B_c = 0$, then set the fractional slope $F_a = F_b/F_c$. This eliminates one constant times variable multiplication.

The resulting formula is

$$Q_a = \frac{Q_b}{Q_c} 2^{E_b - E_c - E_a} + \frac{(B_b/F_b)}{Q_c} 2^{-E_c - E_a}.$$

If $B_c \neq 0$, then no clear recommendation can be made.

Inherited Scaling for Maximum Precision

You can determine the maximum-precision scaling if the range of the variable is known. “Maximize Precision” on page 36-18 shows that you can determine the range of a fixed-point operation from

$$\max(\tilde{V}_a)$$

and

$$\min(\tilde{V}_a).$$

For division, you can determine the range from

$$\min(\tilde{V}_a) = \min(V_{LL}, V_{LH}, V_{HL}, V_{HH}),$$

$$\max(\tilde{V}_a) = \max(V_{LL}, V_{LH}, V_{HL}, V_{HH}),$$

where for nonzero denominators

$$V_{LL} = \min(\tilde{V}_b) / \min(\tilde{V}_c),$$

$$V_{LH} = \min(\tilde{V}_b) / \max(\tilde{V}_c),$$

$$V_{HL} = \max(\tilde{V}_b) / \min(\tilde{V}_c),$$

$$V_{HH} = \max(\tilde{V}_b) / \max(\tilde{V}_c).$$

Binary-Point-Only Scaling

For binary-point-only scaling, finding Q_a results in this simple expression:

$$Q_a = \frac{Q_b}{Q_c} 2^{E_b - E_c - E_a}.$$

Note For the last two formulas involving Q_a , a divide by zero and zero divided by zero are possible. In these cases, the hardware will give some default behavior but you must make sure that these default responses give meaningful results for the embedded system.

Summary

From the previous analysis of fixed-point variables scaled within the general [Slope Bias] encoding scheme, you can conclude

- Addition, subtraction, multiplication, and division can be very involved unless certain choices are made for the biases and slopes.

- Binary-point-only scaling guarantees simpler math, but generally sacrifices some precision.

Note that the previous formulas don't show the following:

- Constants and variables are represented with a finite number of bits.
- Variables are either signed or unsigned.
- Rounding and overflow handling schemes. You must make these decisions before an actual fixed-point realization is achieved.

See Also

“Scaling” on page 35-5

Parameter and Signal Conversions

In this section...

“Introduction” on page 36-39

“Parameter Conversions” on page 36-39

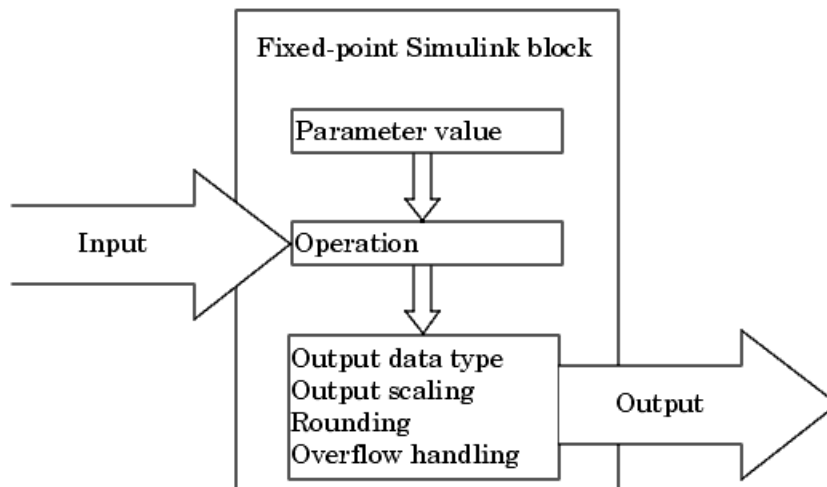
“Signal Conversions” on page 36-40

Introduction

To completely understand the results generated by fixed-point Simulink blocks, you must be aware of these issues:

- When numerical block parameters are converted from doubles to fixed-point data types
- When input signals are converted from one fixed-point data type to another (if at all)
- When arithmetic operations on input signals and parameters are performed

For example, suppose a fixed-point Simulink block performs an arithmetic operation on its input signal and a parameter, and then generates output having characteristics that are specified by the block. The following diagram illustrates how these issues are related.



The sections that follow describe parameter and signal conversions. “Rules for Arithmetic Operations” on page 36-42 discusses arithmetic operations.

Parameter Conversions

Parameters of fixed-point blocks that accept numerical values are always converted from `double` to a fixed-point data type. Parameters can be converted to the input data type, the output data type, or to a data type explicitly specified by the block. For example, the Discrete FIR Filter block converts its **Initial states** parameter to the input data type, and converts its **Numerator coefficient** parameter to a data type you explicitly specify via the block dialog box.

Parameters are always converted before any arithmetic operations are performed. Additionally, parameters are always converted *offline* using round-to-nearest and saturation. Offline conversions are discussed below.

Note Because parameters of fixed-point blocks begin as `double`, they are never precise to more than 53 bits. Therefore, if the output of your fixed-point block is longer than 53 bits, your result might be less precise than you anticipated.

Offline Conversions

An offline conversion is a conversion performed by your development platform (for example, the processor on your PC), and not by the fixed-point processor you are targeting. For example, suppose you are using a PC to develop a program to run on a fixed-point processor, and you need the fixed-point processor to compute

$$y = \left(\frac{ab}{c}\right)u = Cu$$

over and over again. If a , b , and c are constant parameters, it is inefficient for the fixed-point processor to compute ab/c every time. Instead, the PC's processor should compute ab/c offline one time, and the fixed-point processor computes only $C \cdot u$. This eliminates two costly fixed-point arithmetic operations.

Signal Conversions

Consider the conversion of a real-world value from one fixed-point data type to another. Ideally, the values before and after the conversion are equal.

$$V_a = V_b,$$

where V_b is the input value and V_a is the output value. To see how the conversion is implemented, the two ideal values are replaced by the general [Slope Bias] encoding scheme described in "Scaling" on page 35-5:

$$V_i = F_i 2^{E_i} Q_i + B_i.$$

Solving for the output data type's stored integer value, Q_a is obtained:

$$\begin{aligned} Q_a &= \frac{F_b}{F_a} 2^{E_b - E_a} Q_b + \frac{B_b - B_a}{F_a} 2^{-E_a} \\ &= F_s 2^{E_b - E_a} Q_b + B_{net}, \end{aligned}$$

where F_s is the adjusted fractional slope and B_{net} is the net bias. The offline conversions and online conversions and operations are discussed below.

Offline Conversions

Both F_s and B_{net} are computed offline using round-to-nearest and saturation. B_{net} is then stored using the output data type and F_s is stored using an automatically selected data type.

Online Conversions and Operations

The remaining conversions and operations are performed *online* by the fixed-point processor, and depend on the slopes and biases for the input and output data types. The conversions and operations are given by these steps:

- 1 The initial value for Q_a is given by the net bias, B_{net} :

$$Q_a = B_{net}.$$

- 2 The input integer value, Q_b , is multiplied by the adjusted slope, F_s :

$$Q_{RawProduct} = F_s Q_b.$$

- 3 The result of step 2 is converted to the modified output data type where the slope is one and bias is zero:

$$Q_{Temp} = \text{convert}(Q_{RawProduct}).$$

This conversion includes any necessary bit shifting, rounding, or overflow handling.

- 4 The summation operation is performed:

$$Q_a = Q_{Temp} + Q_a.$$

This summation includes any necessary overflow handling.

Streamlining Simulations and Generated Code

Note that the maximum number of conversions and operations is performed when the slopes and biases of the input signal and output signal differ (are mismatched). If the scaling of these signals is identical (matched), the number of operations is reduced from the worst (most inefficient) case. For example, when an input has the same fractional slope and bias as the output, only step 3 is required:

$$Q_a = \text{convert}(Q_b).$$

Exclusive use of binary-point-only scaling for both input signals and output signals is a common way to eliminate mismatched slopes and biases, and results in the most efficient simulations and generated code.

Rules for Arithmetic Operations

Fixed-point arithmetic refers to how signed or unsigned binary words are operated on. The simplicity of fixed-point arithmetic functions such as addition and subtraction allows for cost-effective hardware implementations.

The sections that follow describe the rules that the Simulink software follows when arithmetic operations are performed on inputs and parameters. These rules are organized into four groups based on the operations involved: addition and subtraction, multiplication, division, and shifts. For each of these four groups, the rules for performing the specified operation are presented with an example using the rules.

Computational Units

The core architecture of many processors contains several computational units including arithmetic logic units (ALUs), multiply and accumulate units (MACs), and shifters. These computational units process the binary data directly and provide support for arithmetic computations of varying precision. The ALU performs a standard set of arithmetic and logic operations as well as division. The MAC performs multiply, multiply/add, and multiply/subtract operations. The shifter performs logical and arithmetic shifts, normalization, denormalization, and other operations.

Addition and Subtraction

Addition is the most common arithmetic operation a processor performs. When two n-bit numbers are added together, it is always possible to produce a result with n + 1 nonzero digits due to a carry from the leftmost digit. For two's complement addition of two numbers, there are three cases to consider:

- If both numbers are positive and the result of their addition has a sign bit of 1, then overflow has occurred; otherwise, the result is correct.
- If both numbers are negative and the sign of the result is 0, then overflow has occurred; otherwise, the result is correct.
- If the numbers are of unlike sign, overflow cannot occur and the result is always correct.

Fixed-Point Simulink Blocks Summation Process

Consider the summation of two numbers. Ideally, the real-world values obey the equation

$$V_a = \pm V_b \pm V_c,$$

where V_b and V_c are the input values and V_a is the output value. To see how the summation is actually implemented, the three ideal values should be replaced by the general [Slope Bias] encoding scheme described in "Scaling" on page 35-5:

$$V_i = F_i 2^{E_i} Q_i + B_i.$$

The equation in "Addition" on page 36-31 gives the solution of the resulting equation for the stored integer, Q_a . Using shorthand notation, that equation becomes

$$Q_a = \pm F_{sb} 2^{E_b - E_a} Q_b \pm F_{sc} 2^{E_c - E_a} Q_c + B_{net},$$

where F_{sb} and F_{sc} are the adjusted fractional slopes and B_{net} is the net bias. The offline conversions and online conversions and operations are discussed below.

Offline Conversions

F_{sb} , F_{sc} , and B_{net} are computed offline using round-to-nearest and saturation. Furthermore, B_{net} is stored using the output data type.

Online Conversions and Operations

The remaining operations are performed online by the fixed-point processor, and depend on the slopes and biases for the input and output data types. The worst (most inefficient) case occurs when the slopes and biases are mismatched. The worst-case conversions and operations are given by these steps:

- 1 The initial value for Q_a is given by the net bias, B_{net} :

$$Q_a = B_{net}.$$

- 2 The first input integer value, Q_b , is multiplied by the adjusted slope, F_{sb} :

$$Q_{RawProduct} = F_{sb}Q_b.$$

- 3 The previous product is converted to the modified output data type where the slope is one and the bias is zero:

$$Q_{Temp} = \text{convert}(Q_{RawProduct}).$$

This conversion includes any necessary bit shifting, rounding, or overflow handling.

- 4 The summation operation is performed:

$$Q_a = Q_a \pm Q_{Temp}.$$

This summation includes any necessary overflow handling.

- 5 Steps 2 to 4 are repeated for every number to be summed.

It is important to note that bit shifting, rounding, and overflow handling are applied to the intermediate steps (3 and 4) and not to the overall sum.

For more information, see “The Summation Process” on page 36-49.

Streamlining Simulations and Generated Code

If the scaling of the input and output signals is matched, the number of summation operations is reduced from the worst (most inefficient) case. For example, when an input has the same fractional slope as the output, step 2 reduces to multiplication by one and can be eliminated. Trivial steps in the summation process are eliminated for both simulation and code generation. Exclusive use of binary-point-only scaling for both input signals and output signals is a common way to eliminate mismatched slopes and biases, and results in the most efficient simulations and generated code.

Multiplication

The multiplication of an n-bit binary number with an m-bit binary number results in a product that is up to $m + n$ bits in length for both signed and unsigned words. Most processors perform n-bit by n-bit multiplication and produce a 2n-bit result (double bits) assuming there is no overflow condition.

Fixed-Point Simulink Blocks Multiplication Process

Consider the multiplication of two numbers. Ideally, the real-world values obey the equation

$$V_a = V_b V_c.$$

where V_b and V_c are the input values and V_a is the output value. To see how the multiplication is actually implemented, the three ideal values should be replaced by the general [Slope Bias] encoding scheme described in “Scaling” on page 35-5:

$$V_i = F_i 2^{E_i} Q_i + B_i.$$

The solution of the resulting equation for the output stored integer, Q_a , is given below:

$$Q_a = \frac{F_b F_c}{F_a} 2^{E_b + E_c - E_a} Q_b Q_c + \frac{F_b B_c}{F_a} 2^{E_b - E_a} Q_b + \frac{F_c B_b}{F_a} 2^{E_c - E_a} Q_c + \frac{B_b B_c - B_a}{F_a} 2^{-E_a}.$$

Multiplication with Nonzero Biases and Mismatched Fractional Slopes

The worst-case implementation of the above equation occurs when the slopes and biases of the input and output signals are mismatched. In such cases, several low-level integer operations are required to carry out the high-level multiplication (or division). Implementation choices made about these low-level computations can affect the computational efficiency, rounding errors, and overflow.

In Simulink blocks, the actual multiplication or division operation is always performed on fixed-point variables that have zero biases. If an input has nonzero bias, it is converted to a representation that has binary-point-only scaling before the operation. If the result is to have nonzero bias, the operation is first performed with temporary variables that have binary-point-only scaling. The result is then converted to the data type and scaling of the final output.

If both the inputs and the output have nonzero biases, then the operation is broken down as follows:

$$\begin{aligned} V_{1Temp} &= V_1, \\ V_{2Temp} &= V_2, \\ V_{3Temp} &= V_{1Temp} V_{2Temp}, \\ V_3 &= V_{3Temp}, \end{aligned}$$

where

$$\begin{aligned} V_{1Temp} &= 2^{E_{1Temp}} Q_{1Temp}, \\ V_{2Temp} &= 2^{E_{2Temp}} Q_{2Temp}, \\ V_{3Temp} &= 2^{E_{3Temp}} Q_{3Temp}. \end{aligned}$$

These equations show that the temporary variables have binary-point-only scaling. However, the equations do not indicate the signedness, word lengths, or values of the fixed exponent of these variables. The Simulink software assigns these properties to the temporary variables based on the following goals:

- Represent the original value without overflow.

The data type and scaling of the original value define a maximum and minimum real-world value:

$$V_{Max} = F 2^E Q_{MaxInteger} + B,$$

$$V_{Min} = F2^E Q_{MinInteger} + B.$$

The data type and scaling of the temporary value must be able to represent this range without overflow. Precision loss is possible, but overflow is never allowed.

- Use a data type that leads to efficient operations.

This goal is relative to the target that you will use for production deployment of your design. For example, suppose that you will implement the design on a 16-bit fixed-point processor that provides a 32-bit long, 16-bit int, and 8-bit short or char. For such a target, preserving efficiency means that no more than 32 bits are used, and the smaller sizes of 8 or 16 bits are used if they are sufficient to maintain precision.

- Maintain precision.

Ideally, every possible value defined by the original data type and scaling is represented perfectly by the temporary variable. However, this can require more bits than is efficient. Bits are discarded, resulting in a loss of precision, to the extent required to preserve efficiency.

For example, consider the following, assuming a 16-bit microprocessor target:

$$V_{Original} = Q_{Original} + -43.25,$$

where $Q_{Original}$ is an 8-bit, unsigned data type. For this data type,

$$Q_{MaxInteger} = 225,$$

$$Q_{MinInteger} = 0,$$

so

$$V_{Max} = 211.75,$$

$$V_{Min} = -43.25.$$

The minimum possible value is negative, so the temporary variable must be a signed integer data type. The original variable has a slope of 1, but the bias is expressed with greater precision with two digits after the binary point. To get full precision, the fixed exponent of the temporary variable has to be -2 or less. The Simulink software selects the least possible precision, which is generally the most efficient, unless overflow issues arise. For a scaling of 2^{-2} , selecting signed 16-bit or signed 32-bit avoids overflow. For efficiency, the Simulink software selects the smaller choice of 16 bits. If the original variable is an input, then the equations to convert to the temporary variable are

$$\text{uint8_T } Q_{Original},$$

$$\text{uint16_T } Q_{Temp},$$

$$Q_{Temp} = ((\text{int16_T})Q_{Original} \ll 2) - 173.$$

Multiplication with Zero Biases and Mismatched Fractional Slopes

When the biases are zero and the fractional slopes are mismatched, the implementation reduces to

$$Q_a = \frac{F_b F_c}{F_a} 2^{E_b + E_c - E_a} Q_b Q_c.$$

Offline Conversions

The quantity

$$F_{Net} = \frac{F_b F_c}{F_a}$$

is calculated offline using round-to-nearest and saturation. F_{Net} is stored using a fixed-point data type of the form

$$2^{E_{Net}} Q_{Net},$$

where E_{Net} and Q_{Net} are selected automatically to best represent F_{Net} .

Online Conversions and Operations

- 1 The integer values Q_b and Q_c are multiplied:

$$Q_{RawProduct} = Q_b Q_c.$$

To maintain the full precision of the product, the binary point of $Q_{RawProduct}$ is given by the sum of the binary points of Q_b and Q_c .

- 2 The previous product is converted to the output data type:

$$Q_{Temp} = \text{convert}(Q_{RawProduct}).$$

This conversion includes any necessary bit shifting, rounding, or overflow handling. "Signal Conversions" on page 36-40 discusses conversions.

- 3 The multiplication

$$Q_{2RawProduct} = Q_{Temp} Q_{Net}$$

is performed.

- 4 The previous product is converted to the output data type:

$$Q_a = \text{convert}(Q_{2RawProduct}).$$

This conversion includes any necessary bit shifting, rounding, or overflow handling. "Signal Conversions" on page 36-40 discusses conversions.

- 5 Steps 1 through 4 are repeated for each additional number to be multiplied.

Multiplication with Zero Biases and Matching Fractional Slopes

When the biases are zero and the fractional slopes match, the implementation reduces to

$$Q_a = 2^{E_b + E_c - E_a} Q_b Q_c.$$

Offline Conversions

No offline conversions are performed.

Online Conversions and Operations

- 1 The integer values Q_b and Q_c are multiplied:

$$Q_{RawProduct} = Q_b Q_c.$$

To maintain the full precision of the product, the binary point of $Q_{RawProduct}$ is given by the sum of the binary points of Q_b and Q_c .

- 2 The previous product is converted to the output data type:

$$Q_a = \text{convert}(Q_{RawProduct}).$$

This conversion includes any necessary bit shifting, rounding, or overflow handling. “Signal Conversions” on page 36-40 discusses conversions.

- 3 Steps 1 and 2 are repeated for each additional number to be multiplied.

For more information, see “The Multiplication Process” on page 36-51.

Division

This section discusses the division of quantities with zero bias.

Note When any input to a division calculation has nonzero bias, the operations performed exactly match those for multiplication described in “Multiplication with Nonzero Biases and Mismatched Fractional Slopes” on page 36-44.

Fixed-Point Simulink Blocks Division Process

Consider the division of two numbers. Ideally, the real-world values obey the equation

$$V_a = V_b/V_c,$$

where V_b and V_c are the input values and V_a is the output value. To see how the division is actually implemented, the three ideal values should be replaced by the general [Slope Bias] encoding scheme described in “Scaling” on page 35-5:

$$V_i = F_i 2^{E_i} Q_i + B_i.$$

For the case where the slope adjustment factors are one and the biases are zero for all signals, the solution of the resulting equation for the output stored integer, Q_a , is given by the following equation:

$$Q_a = 2^{E_b - E_c - E_a} (Q_b/Q_c).$$

This equation involves an integer division and some bit shifts. If $E_a > E_b - E_c$, then any bit shifts are to the right and the implementation is simple. However, if $E_a < E_b - E_c$, then the bit shifts are to the left and the implementation can be more complicated. The essential issue is that the output has more precision than the integer division provides. To get full precision, a *fractional* division is needed. The C programming language provides access to integer division only for fixed-point data types. Depending on the size of the numerator, you can obtain some of the fractional bits by performing a shift prior to the integer division. In the worst case, it might be necessary to resort to repeated subtractions in software.

In general, division of values is an operation that should be avoided in fixed-point embedded systems. Division where the output has more precision than the integer division (i.e., $E_a < E_b - E_c$) should be used with even greater reluctance.

For more information, see “The Division Process” on page 36-53.

Shifts

Nearly all microprocessors and digital signal processors support well-defined *bit-shift* (or simply *shift*) operations for integers. For example, consider the 8-bit unsigned integer 00110101. The results of a 2-bit shift to the left and a 2-bit shift to the right are shown in the following table.

| Shift Operation | Binary Value | Decimal Value |
|----------------------------|--------------|---------------|
| No shift (original number) | 00110101 | 53 |
| Shift left by 2 bits | 11010100 | 212 |
| Shift right by 2 bits | 00001101 | 13 |

You can perform a shift using the Simulink Shift Arithmetic block. Use this block to perform a bit shift, a binary point shift, or both

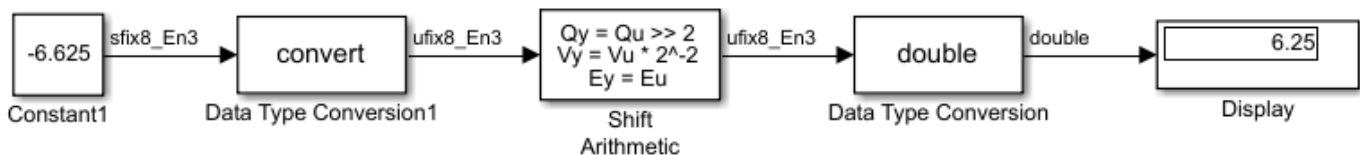
Shifting Bits to the Right

The special case of shifting bits to the right requires consideration of the treatment of the leftmost bit, which can contain sign information. A shift to the right can be classified either as a *logical* shift right or an *arithmetic* shift right. For a logical shift right, a 0 is incorporated into the most significant bit for each bit shift. For an arithmetic shift right, the most significant bit is recycled for each bit shift.

The Shift Arithmetic block performs an arithmetic shift right and, therefore, recycles the most significant bit for each bit shift right. For example, given the fixed-point number 11001.011 (-6.625), a bit shift two places to the right with the binary point unmoved yields the number 11110.010 (-1.75), as shown in the model below:



To perform a logical shift right on a signed number using the Shift Arithmetic block, use the Data Type Conversion block to cast the number as an unsigned number of equivalent length and scaling. This model shows that the fixed-point signed number 11001.001 (-6.625) becomes 00110.010 (6.25).

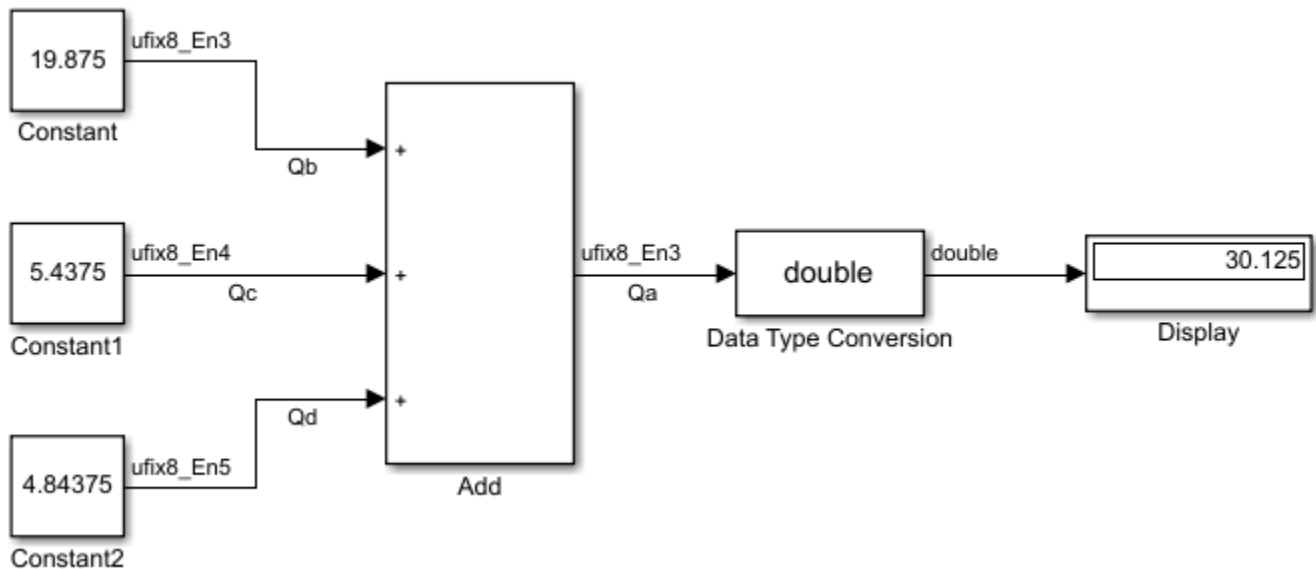


The Summation Process

Addition is the most common arithmetic operation a processor performs. When two n -bit numbers are added together, it is always possible to produce a result with $n + 1$ nonzero digits due to a carry from the leftmost digit.

Suppose you want to sum three numbers. Each of these numbers is represented by an 8-bit word, and each has a different binary-point-only scaling. Additionally, the output is restricted to an 8-bit word with binary-point-only scaling of 2^{-3} .

The summation is shown in the following model for the input values 19.875, 5.4375, and 4.84375.



The sum follows these steps:

- 1 Because the biases are matched, the initial value of Q_a is trivial:

$$Q_a = 00000.000.$$
- 2 The first number to be summed (19.875) has a fractional slope that matches the output fractional slope. Furthermore, the binary points and storage types are identical, so the conversion is trivial:

$$Q_b = 10011.111,$$

$$Q_{Temp} = Q_b.$$
- 3 The summation operation is performed:

$$Q_a = Q_a + Q_{Temp} = 10011.111.$$
- 4 The second number to be summed (5.4375) has a fractional slope that matches the output fractional slope, so a slope adjustment is not needed. The storage data types also match, but the difference in binary points requires that both the bits and the binary point be shifted one place to the right:

$$Q_c = 0101.0111,$$

$$Q_{Temp} = \text{convert}(Q_c)$$

$$Q_{Temp} = 00101.011.$$

Note that a loss in precision of one bit occurs, with the resulting value of Q_{Temp} determined by the rounding mode. For this example, round-to-floor is used. Overflow cannot occur in this case because the bits and binary point are both shifted to the right.

- 5 The summation operation is performed:

$$\begin{aligned} Q_a &= Q_a + Q_{Temp} \\ &10011.111 \\ &= \frac{+00101.011}{11001.010} = 25.250. \end{aligned}$$

Note that overflow did not occur, but it is possible for this operation.

- 6 The third number to be summed (4.84375) has a fractional slope that matches the output fractional slope, so a slope adjustment is not needed. The storage data types also match, but the difference in binary points requires that both the bits and the binary point be shifted two places to the right:

$$\begin{aligned} Q_d &= 100.11011, \\ Q_{Temp} &= \text{convert}(Q_d) \\ Q_{Temp} &= 00100.110. \end{aligned}$$

Note that a loss in precision of two bit occurs, with the resulting value of Q_{Temp} determined by the rounding mode. For this example, round-to-floor is used. Overflow cannot occur in this case because the bits and binary point are both shifted to the right.

- 7 The summation operation is performed:

$$\begin{aligned} Q_a &= Q_a + Q_{Temp} \\ &11001.010 \\ &= \frac{+00100.110}{11110.000} = 30.000. \end{aligned}$$

Note that overflow did not occur, but it is possible for this operation.

As shown here, the result of step 7 differs from the ideal sum:

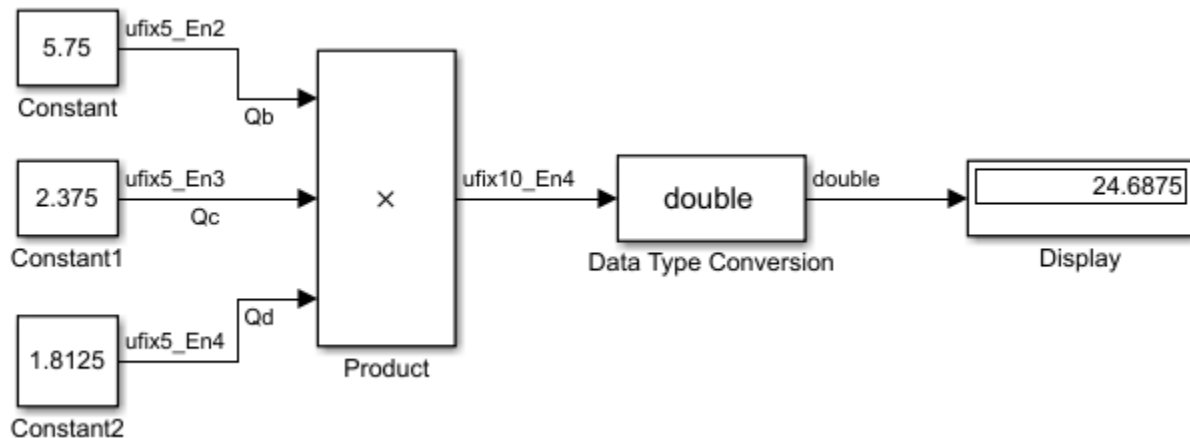
$$\begin{aligned} &10011.111 \\ &0101.0111 \\ &= \frac{+100.11011}{11110.001} = 30.125. \end{aligned}$$

Blocks that perform addition and subtraction include the Add, Gain, and Discrete FIR Filter blocks.

The Multiplication Process

The multiplication of an n-bit binary number with an m-bit binary number results in a product that is up to m + n bits in length for both signed and unsigned words.

Suppose you want to multiply three numbers. Each of these numbers is represented by a 5-bit word, and each has a different binary-point-only scaling. Additionally, the output is restricted to a 10-bit word with binary-point-only scaling of 2^{-4} . The multiplication is shown in the following model for the input values 5.75, 2.375, and 1.8125.



Applying the rules from the previous section, the multiplication follows these steps:

- 1 The first two numbers (5.75 and 2.375) are multiplied:

$$\begin{aligned}
 Q_{RawProduct} &= 101.11 \\
 &\quad \times 10.011 \\
 &\hline
 &101.11 \cdot 2^{-3} \\
 &101.11 \cdot 2^{-2} \\
 &+ 101.11 \cdot 2^1 \\
 \hline
 01101.10101 &= 13.65625.
 \end{aligned}$$

Note that the binary point of the product is given by the sum of the binary points of the multiplied numbers.

- 2 The result of step 1 is converted to the output data type:

$$\begin{aligned}
 Q_{Temp} &= \text{convert}(Q_{RawProduct}) \\
 &= 001101.1010 = 13.6250.
 \end{aligned}$$

“Signal Conversions” on page 36-40 discusses conversions. Note that a loss in precision of one bit occurs, with the resulting value of Q_{Temp} determined by the rounding mode. For this example, round-to-floor is used. Furthermore, overflow did not occur but is possible for this operation.

- 3 The result of step 2 and the third number (1.8125) are multiplied:

$$\begin{aligned} Q_{RawProduct} &= 01101.1010 \\ &\quad \times 1.1101 \\ &\quad \hline &\quad 1101.1010 \cdot 2^{-4} \\ &1101.1010 \cdot 2^{-2} \\ &1101.1010 \cdot 2^{-1} \\ &\hline &+ 1101.1010 \cdot 2^0 \\ &\hline 0011000.10110010 &= 24.6953125. \end{aligned}$$

Note that the binary point of the product is given by the sum of the binary points of the multiplied numbers.

- 4 The product is converted to the output data type:

$$\begin{aligned} Q_a &= \text{convert}(Q_{RawProduct}) \\ &= 011000.1011 = 24.6875. \end{aligned}$$

“Signal Conversions” on page 36-40 discusses conversions. Note that a loss in precision of 4 bits occurred, with the resulting value of Q_{Temp} determined by the rounding mode. For this example, round-to-floor is used. Furthermore, overflow did not occur but is possible for this operation.

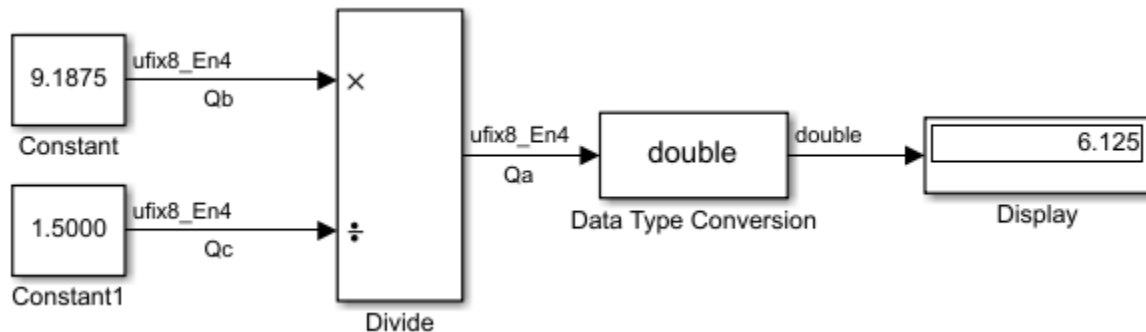
Blocks that perform multiplication include the Product, Discrete FIR Filter, and Gain blocks.

The Division Process

The C programming language provides access to integer division only for fixed-point data types. Depending on the size of the numerator, you can obtain some of the fractional bits by performing a shift prior to the integer division.

Suppose you want to divide two numbers. Each of these numbers is represented by an 8-bit word, and each has a binary-point-only scaling of 2^{-4} . Additionally, the output is restricted to an 8-bit word with binary-point-only scaling of 2^{-4} .

The division of 9.1875 by 1.5000 is shown in the following model.



For this example,

$$\begin{aligned} Q_a &= 2^{-4 - (-4) - (-4)}(Q_b/Q_c) \\ &= 2^4(Q_b/Q_c). \end{aligned}$$

Assuming a large data type was available, this could be implemented as

$$Q_a = \frac{(2^4 Q_b)}{Q_c},$$

where the numerator uses the larger data type. If a larger data type was not available, integer division combined with four repeated subtractions would be used. Both approaches produce the same result, with the former being more efficient.

Shifts

Nearly all microprocessors and digital signal processors support well-defined *bit-shift* (or simply *shift*) operations for integers. For example, consider the 8-bit unsigned integer 00110101. The results of a 2-bit shift to the left and a 2-bit shift to the right are shown in the following table.

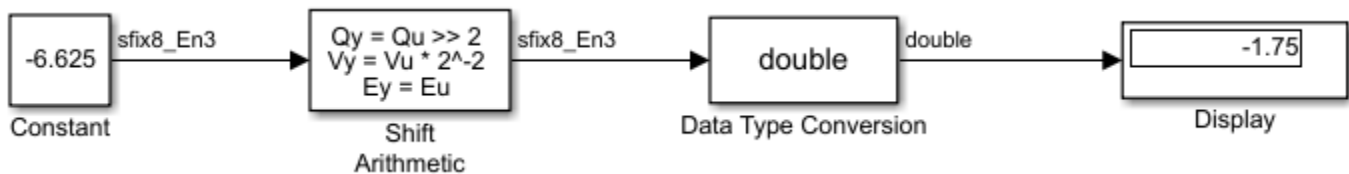
| Shift Operation | Binary Value | Decimal Value |
|----------------------------|--------------|---------------|
| No shift (original number) | 00110101 | 53 |
| Shift left by 2 bits | 11010100 | 212 |
| Shift right by 2 bits | 00001101 | 13 |

You can perform a shift using the Simulink Shift Arithmetic block. Use this block to perform a bit shift, a binary point shift, or both

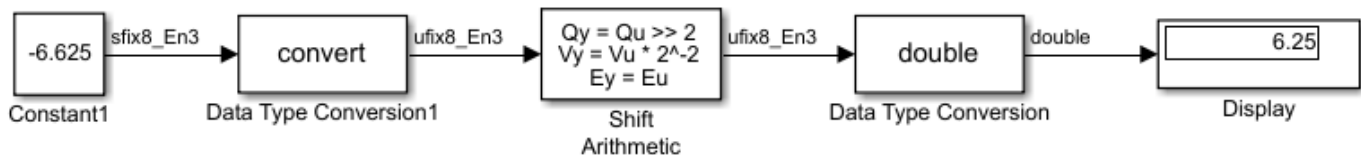
Shifting Bits to the Right

The special case of shifting bits to the right requires consideration of the treatment of the leftmost bit, which can contain sign information. A shift to the right can be classified either as a *logical* shift right or an *arithmetic* shift right. For a logical shift right, a 0 is incorporated into the most significant bit for each bit shift. For an arithmetic shift right, the most significant bit is recycled for each bit shift.

The Shift Arithmetic block performs an arithmetic shift right and, therefore, recycles the most significant bit for each bit shift right. For example, given the fixed-point number 11001.011 (-6.625), a bit shift two places to the right with the binary point unmoved yields the number 11110.010 (-1.75), as shown in the model below:



To perform a logical shift right on a signed number using the Shift Arithmetic block, use the Data Type Conversion block to cast the number as an unsigned number of equivalent length and scaling, as shown in the following model. The model shows that the fixed-point signed number 11001.001 (-6.625) becomes 00110.010 (6.25).



Conversions and Arithmetic Operations

This example uses the Discrete FIR Filter block to illustrate when parameters are converted from a double to a fixed-point number, when the input data type is converted to the output data type, and when the rules for addition, subtraction, and multiplication are applied.

Note If a block can perform all four arithmetic operations, then the rules for multiplication and division are applied first. The Discrete FIR Filter block is an example of this.

Suppose you configure the Discrete FIR Filter block for two outputs, where the first output is given by

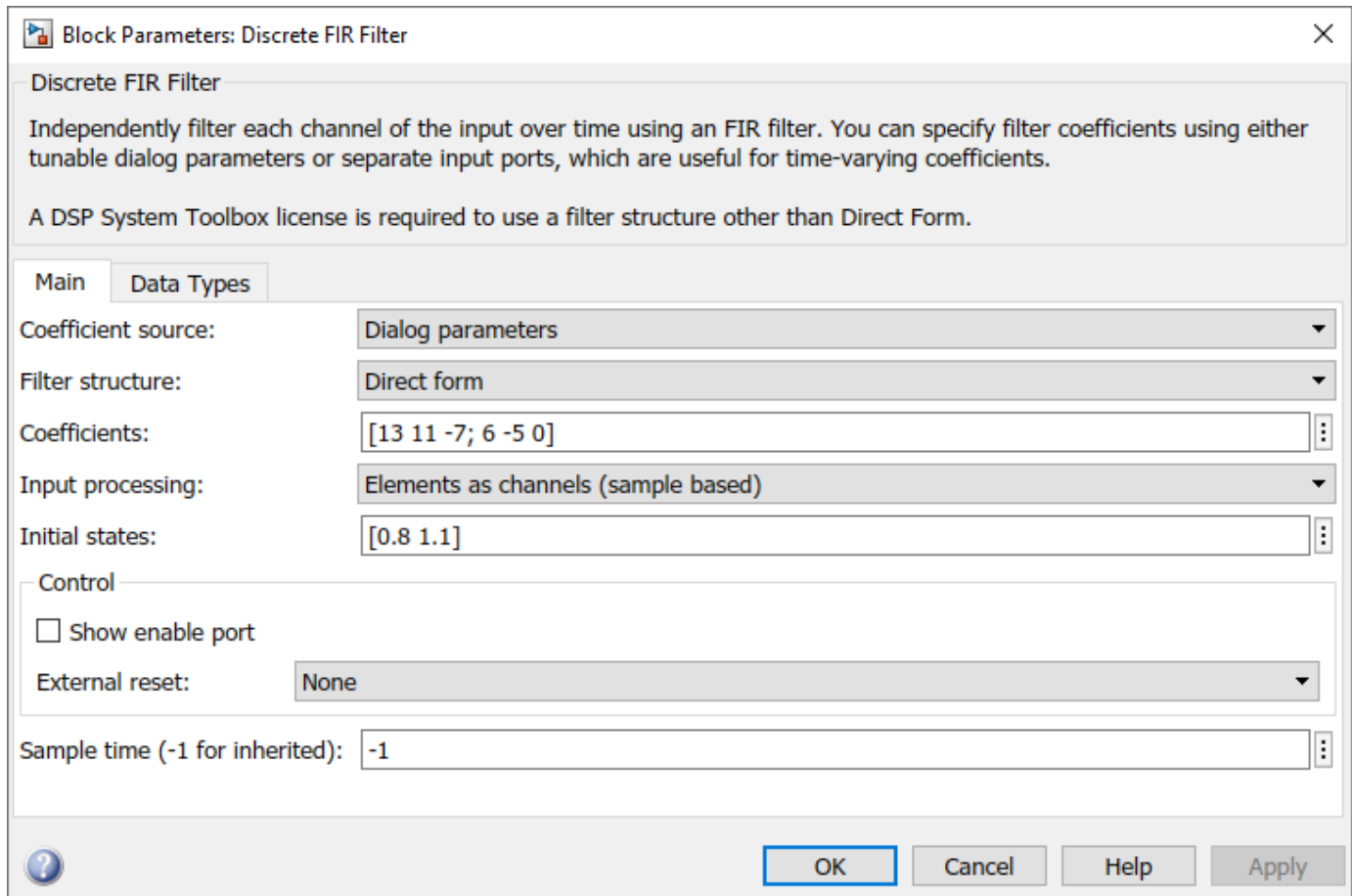
$$y_1(k) = 13 \cdot u(k) + 11 \cdot u(k - 1) - 7 \cdot u(k - 2),$$

and the second output is given by

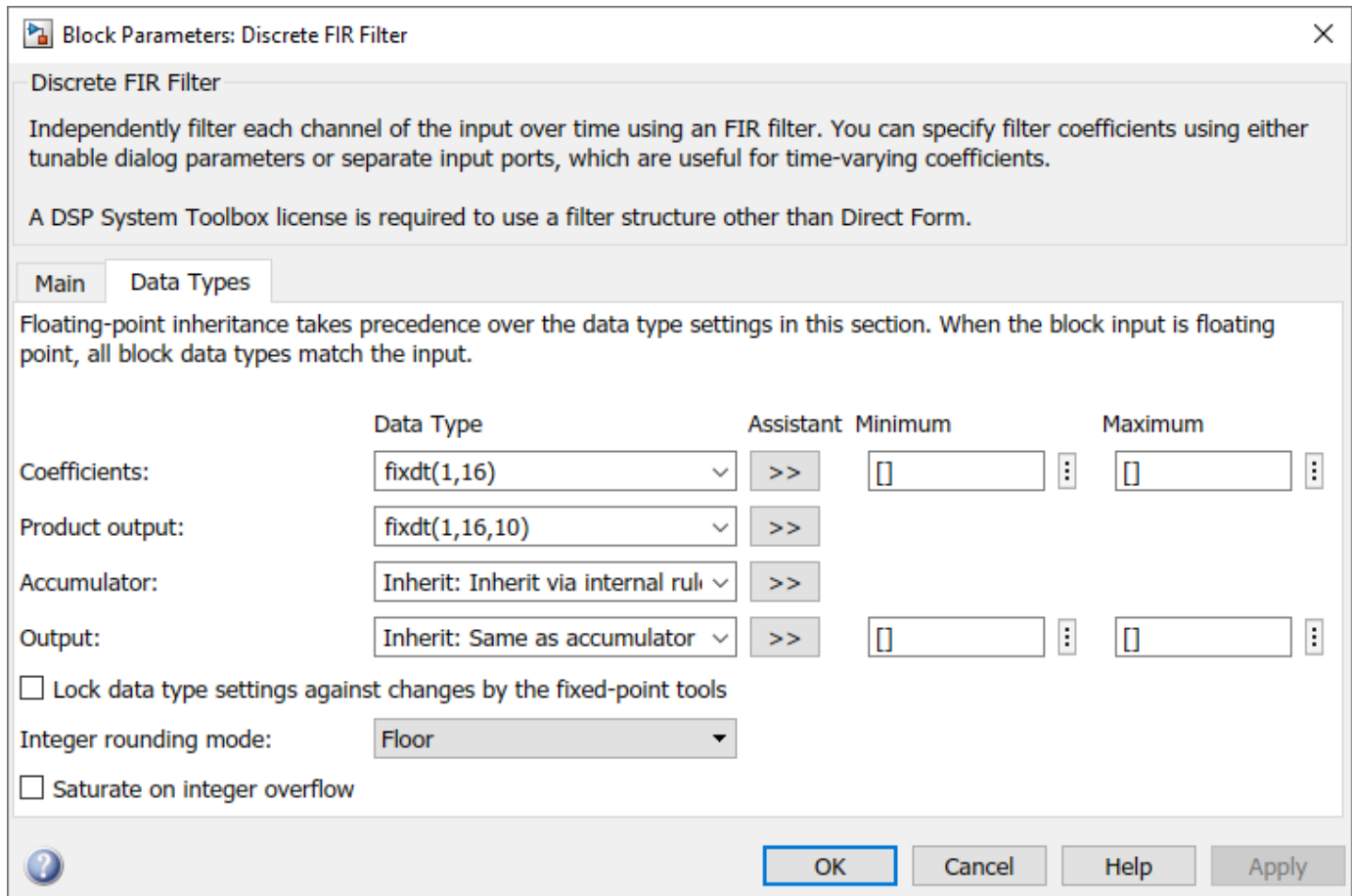
$$y_2(k) = 6 \cdot u(k) - 5 \cdot u(k - 1).$$

Additionally, the initial values of $u(k-1)$ and $u(k-2)$ are given by 0.8 and 1.1, respectively, and all inputs, parameters, and outputs have binary-point-only scaling.

To configure the Discrete FIR Filter block for this situation, on the **Main** pane of its dialog box, you must specify the **Coefficients** parameter as [13 11 -7; 6 -5 0] and the **Initial states** parameter as [0.8 1.1], as shown here.



Similarly, configure the options on the **Data Types** pane of the block dialog box to appear as follows:



The Discrete FIR Filter block performs parameter conversions and block operations in the following order:

- 1 The **Coefficients** parameter is converted offline from doubles to the **Coefficients** data type using round-to-nearest and saturation.

The **Initial states** parameter is converted offline from doubles to the input data type using round-to-nearest and saturation.

- 2 The coefficients and inputs are multiplied together for the initial time step for both outputs. For $y_1(0)$, the operations $13 \cdot u(0)$, $11 \cdot 0.8$, and $-7 \cdot 1.1$ are performed, while for $y_2(0)$, the operations $6 \cdot u(0)$ and $-5 \cdot 0.8$ are performed.

The results of these operations are stored as **Product output**.

- 3 The sum is carried out in **Accumulator**. The final summation result is then converted to **Output**.
- 4 Steps 2 and 3 repeat for subsequent time steps.

See Also

Discrete FIR Filter

More About

- “Parameter and Signal Conversions” on page 36-39
- “Rules for Arithmetic Operations” on page 36-42

Realization Structures

- “Realizing Fixed-Point Digital Filters” on page 37-2
- “Targeting an Embedded Processor” on page 37-3
- “Canonical Forms” on page 37-5

Realizing Fixed-Point Digital Filters

| |
|--|
| In this section... |
| “Introduction” on page 37-2 |
| “Realizations and Data Types” on page 37-2 |

Introduction

This chapter investigates how you can realize fixed-point digital filters using Simulink blocks and the Fixed-Point Designer software.

The Fixed-Point Designer software addresses the needs of the control system, signal processing, and other fields where algorithms are implemented on fixed-point hardware. In signal processing, a digital filter is a computational algorithm that converts a sequence of input numbers to a sequence of output numbers. The algorithm is designed such that the output signal meets frequency-domain or time-domain constraints (desirable frequency components are passed, undesirable components are rejected).

In general terms, a discrete transfer function controller is a form of a digital filter. However, a digital controller can contain nonlinear functions such as lookup tables in addition to a discrete transfer function. This guide uses the term *digital filter* when referring to discrete transfer functions.

Note To design and implement a wide variety of floating-point and fixed-point filters suitable for use in signal processing applications and for deployment on DSP chips, use the DSP System Toolbox software.

Realizations and Data Types

In an ideal world, where numbers, calculations, and storage of states have infinite precision and range, there are virtually an infinite number of realizations for the same system. In theory, these realizations are all identical.

In the more realistic world of double-precision numbers, calculations, and storage of states, small nonlinearities are introduced by the finite precision and range of floating-point data types. Therefore, each realization of a given system produces different results. In most cases however, these differences are small.

In the world of fixed-point numbers, where precision and range are limited, the differences in the realization results can be very large. Therefore, you must carefully select the data type, word size, and scaling for each realization element such that results are accurately represented. To assist you with this selection, design rules for modeling dynamic systems with fixed-point math are provided in “Targeting an Embedded Processor” on page 37-3.

Targeting an Embedded Processor

In this section...

“Introduction” on page 37-3

“Size Assumptions” on page 37-3

“Operation Assumptions” on page 37-3

“Design Rules” on page 37-4

Introduction

The sections that follow describe issues that often arise when targeting a fixed-point design for use on an embedded processor, such as some general assumptions about integer sizes and operations available on embedded processors. These assumptions lead to design issues and design rules that might be useful for your specific fixed-point design.

Size Assumptions

Embedded processors are typically characterized by a particular bit size. For example, the terms “8-bit micro,” “32-bit micro,” or “16-bit DSP” are common. It is generally safe to assume that the processor is predominantly geared to processing integers of the specified bit size. Integers of the specified bit size are referred to as the *base data type*. Additionally, the processor typically provides some support for integers that are twice as wide as the base data type. Integers consisting of double bits are referred to as the *accumulator data type*. For example a 16-bit micro has a 16-bit base data type and a 32-bit accumulator data type.

Although other data types may be supported by the embedded processor, this section describes only the base and accumulator data types.

Operation Assumptions

The embedded processor operations discussed in this section are limited to the needs of a basic simulation diagram. Basic simulations use multiplication, addition, subtraction, and delays. Fixed-point models also need shifts to do scaling conversions. For all these operations, the embedded processor should have native instructions that allow the base data type as inputs. For accumulator-type inputs, the processor typically supports addition, subtraction, and delay (storage/retrieval from memory), but not multiplication.

Multiplication is typically not supported for accumulator-type inputs because of complexity and size issues. A difficulty with multiplication is that the output needs to be twice as big as the inputs for full precision. For example, multiplying two 16-bit numbers requires a 32-bit output for full precision. The need to handle the outputs from a multiplication operation is one of the reasons embedded processors include accumulator-type support. However, if multiplication of accumulator-type inputs is also supported, then there is a need to support a data type that is twice as big as the accumulator type. To restrict this additional complexity, multiplication is typically not supported for inputs of the accumulator type.

Design Rules

The important design rules that you should be aware of when modeling dynamic systems with fixed-point math follow.

Design Rule 1: Only Multiply Base Data Types

It is best to multiply only inputs of the base data type. Embedded processors typically provide an instruction for the multiplication of base-type inputs, but not for the multiplication of accumulator-type inputs. If necessary, you can combine several instructions to handle multiplication of accumulator-type inputs. However, this can lead to large, slow embedded code.

You can insert blocks to convert inputs from the accumulator type to the base type prior to Product or Gain blocks, if necessary.

Design Rule 2: Delays Should Use the Base Data Type

There are two general reasons why a Unit Delay should use only base-type numbers:

- The Unit Delay essentially stores a variable's value to RAM and, one time step later, retrieves that value from RAM. Because the value must be in memory from one time step to the next, the RAM must be exclusively dedicated to the variable and can't be shared or used for another purpose. Using accumulator-type numbers instead of the base data type doubles the RAM requirements, which can significantly increase the cost of the embedded system.
- The Unit Delay typically feeds into a Gain block. The multiplication design rule requires that the input (the unit delay signal) use the base data type.

Design Rule 3: Temporary Variables Can Use the Accumulator Data Type

Except for unit delay signals, most signals are not needed from one time step to the next. This means that the signal values can be temporarily stored in shared and reused memory. This shared and reused memory can be RAM or it can simply be registers in the CPU. In either case, storing the value as an accumulator data type is not much more costly than storing it as a base data type.

Design Rule 4: Summation Can Use the Accumulator Data Type

Addition and subtraction can use the accumulator data type if there is justification. The typical justification is reducing the buildup of errors due to roundoff or overflow.

For example, a common filter operation is a weighted sum of several variables. Multiplying a variable by a weight naturally produces a product of the accumulator type. Before summing, each product can be converted back to the base data type. This approach introduces round-off error into each part of the sum.

Alternatively, the products can be summed using the accumulator data type, and the final sum can be converted to the base data type. Round-off error is introduced in just one point and the precision is generally better. The cost of doing an addition or subtraction using accumulator-type numbers is slightly more expensive, but if there is justification, it is usually worth the cost.

Canonical Forms

The Fixed-Point Designer software does not attempt to standardize on one particular fixed-point digital filter design method. For example, you can produce a design in continuous time and then obtain an “equivalent” discrete-time digital filter using one of many transformation methods. Alternatively, you can design digital filters directly in discrete time. After you obtain a digital filter, it can be realized for fixed-point hardware using any number of canonical forms. Typical canonical forms are the direct form, series form, and parallel form, each of which is outlined in the sections that follow.

For a given digital filter, the canonical forms describe a set of fundamental operations for the processor. Because there are an infinite number of ways to realize a given digital filter, you must make the best realization on a per-system basis. The canonical forms presented in this chapter optimize the implementation with respect to some factor, such as minimum number of delay elements.

In general, when choosing a realization method, you must take these factors into consideration:

- **Cost**

The cost of the realization might rely on minimal code and data size.

- **Timing constraints**

Real-time systems must complete their compute cycle within a fixed amount of time. Some realizations might yield faster execution speed on different processors.

- **Output signal quality**

The limited range and precision of the binary words used to represent real-world numbers will introduce errors. Some realizations are more sensitive to these errors than others.

The Fixed-Point Designer software allows you to evaluate various digital filter realization methods in a simulation environment. Following the development cycle outlined in “Developing and Testing Fixed-Point Systems” on page 34-11, you can fine-tune the realizations with the goal of reducing the cost (code and data size) or increasing signal quality. After you have achieved the desired performance, you can use the Simulink Coder product to generate rapid prototyping C code and evaluate its performance with respect to your system's real-time timing constraints. You can then modify the model based upon feedback from the rapid prototyping system.

The presentation of the various realization structures takes into account that a summing junction is a fundamental operator, thus you may find that the structures presented here look different from those in the fixed-point filter design literature. For each realization form, an example is provided using the transfer function shown here:

$$\begin{aligned}
 H_{ex}(z) &= \frac{1 + 2.2z^{-1} + 1.85z^{-2} + 0.5z^{-3}}{1 - 0.5z^{-1} + 0.84z^{-2} + 0.09z^{-3}} \\
 &= \frac{(1 + 0.5z^{-1})(1 + 1.7z^{-1} + z^{-2})}{(1 + 0.1z^{-1})(1 - 0.6z^{-1} + 0.9z^{-2})} \\
 &= 5.5556 - \frac{3.4639}{1 + 0.1z^{-1}} + \frac{-1.0916 + 3.0086z^{-1}}{1 - 0.6z^{-1} + 0.9z^{-2}}.
 \end{aligned}$$

Fixed-Point Advisor

Use the Fixed-Point Tool to Prepare a System for Conversion

Using the Fixed-Point Tool, you can prepare a model for conversion from a floating-point model or subsystem to an equivalent fixed-point representation. During the preparation stage, the Fixed-Point Tool checks the system under design for compatibility with the conversion process and reports any issues found in the model. When possible, the Fixed-Point Tool automatically changes settings that are not compatible. In cases where the tool is not able to automatically change the settings, the tool notifies you of the changes you must make manually to help the conversion process be successful.

To prepare a system for conversion:

- 1 Open the Fixed-Point Tool. In the **Apps** gallery of the model, select **Fixed-Point Tool**.
- 2 Under **New**, select the **Iterative Fixed-Point Conversion** workflow.
- 3 Under **System Under Design (SUD)**, select the system or subsystem you want to convert.
- 4 Under **Range Collection Mode**, select the method that you want to use to collect ranges. The Fixed-Point Tool uses these collected ranges to later generate data type proposals.

The preparation checks performed by the Fixed-Point Tool differ slightly between the range collection methods.

For more information on deciding which method of range collection is right for your application, see “Choosing a Range Collection Method” on page 42-2.

- 5 Under **Simulation Inputs**, you can specify `Simulink.SimulationInput` objects to exercise your design over its full operating range, or you can select to **Use default model inputs**.
- 6 To specify tolerances for the system, under **Signal Tolerances** in the table, specify tolerances for any signal in the model with signal logging enabled.
- 7 Click **Prepare**. The Fixed-Point Tool checks the system under design and the model containing the system under design for compatibility with the conversion process.

Selecting any of the checks displays additional information in the **Preparation Details** pane. This pane also contains details for resolving remaining issues.

- 8 After addressing any issues found by the Fixed-Point Tool, click **Prepare** to rerun the checks and verify that all issues are now resolved.

Preparation Checks

The following sections describe the checks performed by the Fixed-Point Tool during the preparation stage of the conversion.

Create Restore Point

The Fixed-Point Tool creates a restore point of your model at its current state. If after the conversion you want to restore your design to its state before converting the data types, click the **Restore Original Model** button.

| Status | Description |
|--------|--|
| Pass | This check passes when the Fixed-Point Tool is able to create a restore point for the model. |

| Status | Description |
|--------|---|
| Fail | <p>This check fails when one of the following occurs:</p> <ul style="list-style-type: none"> • The Fixed-Point Tool is unable to create a restore point for the model because the model is not in a writeable directory. • The Fixed-Point Tool is unable to create a restore point for the model because the model contains unsaved changes. |

Hardware Setting Consistency

Before converting your design to fixed point, you must specify the intended target hardware in the Configuration Parameters **Hardware Implementation** pane. These hardware implementation settings must be consistent throughout the model hierarchy of the model containing the system under design. For more information on how the Fixed-Point Tool uses these settings when proposing data types, see “How the Fixed-Point Tool Uses Target Hardware Information” on page 42-48.

| Status | Description |
|------------------|---|
| Pass | <p>This check passes when the intended target hardware is specified for the system under design and the settings do not conflict with the settings of any other systems in the model.</p> |
| Pass with change | <p>When the hardware implementation settings of the system under design are specified, but they do not match other systems within the model hierarchy, for example, if the model contains a referenced model that uses a different hardware configuration, the Fixed-Point Tool updates the hardware implementation settings of the other systems in your model so that they match the settings of the system under design.</p> |

| Status | Description |
|--------|---|
| Fail | <p>This check fails when one of the following two cases occurs.</p> <ul style="list-style-type: none"> The subsystem under design does not specify any target hardware information. <p>To fix this issue, specify target hardware information for the system under design in the Configuration Parameters Hardware Implementation pane.</p> <ul style="list-style-type: none"> The subsystem specifies target hardware information, but the settings do not match other systems in the model hierarchy and the Fixed-Point Tool is not able to change the settings. <p>To fix this issue, manually change the settings of other systems in the model hierarchy to match the settings of the system under design.</p> |

Check Diagnostic Settings

Certain diagnostics that alert you to numerical issues in your design cannot be set to none. This check passes only when the following diagnostic settings in the Configuration Parameters are set to either warning, or error.

- **Diagnostics > Data Validity > Signals > Wrap on overflow**
- **Diagnostics > Data Validity > Signals > Saturate on overflow**
- **Diagnostics > Data Validity > Signals > Simulation range checking**

| Status | Description |
|------------------|---|
| Pass | This check passes when the diagnostic settings of the model containing the system under design are set to either warning or error. |
| Pass with change | When the diagnostic settings are set to none, the Fixed-Point Tool changes these settings to warning. |
| Fail | This check fails when the Fixed-Point Tool is not able to set the diagnostic settings of the model containing the system under design to warning. This may be because the configuration parameters for the model specify a configuration set. |

Unsupported Constructs

The Fixed-Point Tool identifies any blocks or constructs in your system under design that do not support fixed-point types.

| Status | Description |
|------------------|---|
| Pass | This check passes when the system under design does not contain any unsupported constructs. |
| Pass with change | When the system under design contains unsupported constructs, the Fixed-Point Tool encapsulates any unsupported elements in a subsystem containing the unsupported block surrounded by Data Type Conversion blocks. After you complete the conversion process using the Fixed-Point Tool, you can replace the subsystem containing the unsupported block with a lookup table approximation. For more information, see “Convert Floating-Point Model to Fixed Point” on page 40-2. |
| Fail | This check fails when the Fixed-Point Tool is not able to isolate the unsupported constructs using Data Type Conversion blocks. |

System Under Design Boundary

When model objects within the system under design share a data type with objects outside of the system under design, data type propagation issues can occur after conversion to fixed point. You can prevent these propagation issues by isolating the system under design using Data Type Conversion blocks at the inputs and outputs of the system. The Data Type Conversion block converts an input signal of any Simulink software data type to the data type and scaling you specify for its **Output data type** parameter.

| Status | Description |
|------------------|--|
| Pass | This check passes when the system under design is isolated from the rest of the model by Data Type Conversion blocks. |
| Pass with change | When the system under design is not isolated from the rest of the system, the Fixed-Point Tool places data type conversion blocks at the ports of the system under design to isolate it during the conversion. |
| Fail | This check fails when the Fixed-Point Tool is not able to place Data Type Conversion blocks at the ports of the system under design. |

Design Ranges

When you select **Derived ranges** or **Simulation with Range Analysis** as your range collection method, the software performs a static range analysis of your model to derive minimum and maximum range values for signals in the model. This range analysis relies on specified design ranges. The Fixed-Point Tool checks that you have specified design ranges for all input and output ports of the system under design.

| Status | Description |
|---------------|--|
| Pass | This check passes when all input ports and output ports in the system under design have design range information specified. |
| Warn | This check warns when inputs to the subsystem specify design ranges, but outputs do not specify design ranges. To get the best results from range analysis, specify design ranges for both inputs and outputs to the system. |
| Fail | This check fails when inputs and outputs to the system under design are missing design range information. Specify design range information for all inputs of the system under design. |

See Also

More About

- “Convert Floating-Point Model to Fixed Point” on page 40-2
- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9

Fixed-Point Tool

- “Data Type Conversion Overview” on page 39-2
- “Run Management” on page 39-5
- “Convert a Referenced Model to Fixed Point” on page 39-7
- “Control Views in the Fixed-Point Tool” on page 39-13
- “Model Multiple Data Type Behaviors Using a Data Dictionary” on page 39-17
- “Compare Numerical Response of Sum Block and Sum in MATLAB Function Block” on page 39-21

Data Type Conversion Overview

Within digital hardware, numbers are represented as either fixed-point or floating-point data types. For both these data types, word sizes are fixed at a set number of bits. Fixed-point representation often offers advantages in terms of the power consumption, size, memory usage, speed, and cost of the final product. However, the dynamic range of fixed-point values is much smaller than floating-point values with equivalent word sizes. Therefore, in order to avoid overflow or unreasonable quantization errors, fixed-point values must be scaled.

The following summarizes the process of converting data types in a system from floating point to fixed point. After converting data types in your model to an embedded-efficient representation, you can further optimize your design for your intended hardware, generate code and then deploy code onto your target.

- 1 Identify system requirements.
- 2 Model ideal system.
- 3 Convert data types of system to data types that are efficient on target hardware.
- 4 Verify numeric behavior of the converted system.
- 5 Verify performance of the converted system. Optimize performance of the system based on target hardware.
- 6 Generate code.
- 7 Deploy code onto hardware.

Methods for Converting a System to Fixed Point

The Fixed-Point Designer software provides three methods for automatically specifying fixed-point data types for a system in your model. The following table summarizes the methods available for converting a floating-point system to fixed-point data types.

| Method | Description |
|---|---|
| Fixed-Point Tool | <p>The Fixed-Point Tool is a user interface that automates specifying fixed-point data types in a model.</p> <ul style="list-style-type: none"> • Iterative Fixed-Point Conversion workflow — The tool collects range data for model objects. Based on these values, the tool proposes fixed-point data types that maximize precision and cover the range. You can then review the data type proposals and apply them selectively to objects in your model. • Optimized Fixed-Point Conversion — If you know your system behavior tolerances, you can use <code>fxpopt</code> in the Fixed-Point Tool to find the optimal data types for your system which minimize total bit-width (sum of word lengths) of the system while staying within specified tolerances. <p>For an example, see “Convert Floating-Point Model to Fixed Point” on page 40-2.</p> |
| <code>DataTypeWorkflow.Converter</code> | <p>The <code>DataTypeWorkflow.Converter</code> object and its associated object functions are a command-line alternative to the Fixed-Point Tool. These functions offer the same functionality as the Fixed-Point Tool.</p> <p>For an example, see “Convert a Model to Fixed Point Using the Command Line” on page 47-4.</p> |
| <code>fxpopt</code> | <p>If you know your system behavior tolerances, then the command-line <code>fxpopt</code> function can find the optimal data types for your system which minimize total bit-width (sum of word lengths) of the system while staying within specified tolerances.</p> <p>For an example, see “Optimize Fixed-Point Data Types for a System” on page 40-14</p> |

After converting a system to fixed point, verify that the behavior of the fixed-point system meets your requirements. For more information, see “Verify New Settings” on page 42-30.

Optimize performance of the system based on target hardware. For example, replace trigonometric functions with equivalent CORDIC implementations, or use the **Lookup Table Optimizer** app to replace parts of your model with an embedded-efficient lookup table implementation.

See Also

More About

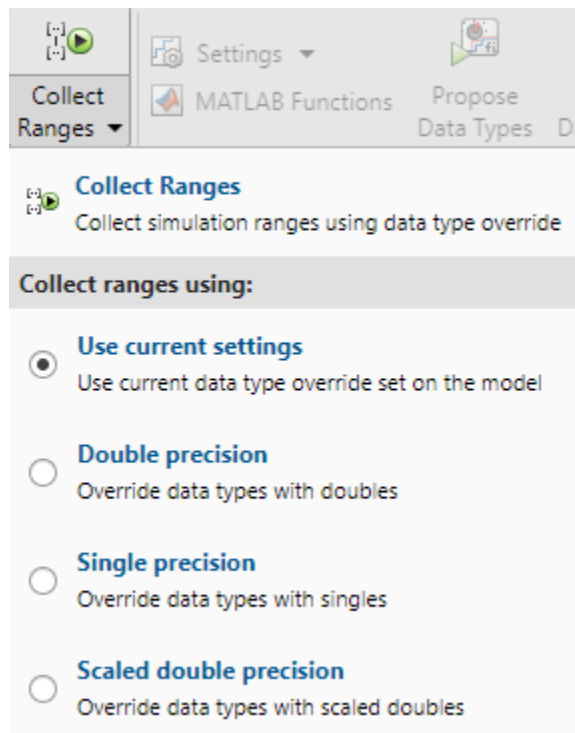
- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9
- “Optimize the Fixed-Point Data Types of a System Using the Fixed-Point Tool” on page 40-40
- “Optimize Generated Code with the Model Advisor” on page 48-21

Run Management



By default, the Fixed-Point Tool creates a run for range collection and a run for verification. By default, these runs are named `BaselineRun` and `EmbeddedRun` respectively.

The tool creates a range collection run when you click the **Collect Ranges** button; it creates a verification run when you click the **Simulate with Embedded Types** button. The default behavior of the range collection run is to use the current data type override settings specified on the model and collect ranges either through simulation, range analysis, or simulation with range analysis. To override your model with double-precision types to avoid quantization effects and collect idealized ranges, select `Double precision`. The verification run simulates your model using the currently specified data types.


The Fixed-Point Tool also provides additional configurations for the range collection and verification runs. You can edit which settings the tool uses by clicking the **Collect Ranges** button arrow and selecting a configuration. The tool overwrites previous range collection runs.



You can edit the default name for the embedded run. Under the **Simulate with Embedded Types** menu, type a new name in the **Run name** field.

 Run to compare in SDI 

Simulate with Embedded Types Compare Results

 **Simulate with Embedded Types**
Simulate the model using the currently specified data types

Verify using:

Specified data types
Use data types specified on model

Scaled double precision
Override data types with scaled doubles

Run name

Convert a Referenced Model to Fixed Point

In this section...

“Open ex_mdhref_controller Model” on page 39-7

“View Model Hierarchy in the Fixed-Point Tool” on page 39-7

“Viewing Simulation Ranges for Referenced Models” on page 39-8

“Propose Data Types for a Referenced Model” on page 39-10

When a system under design contains a referenced model, the Fixed-Point Tool proposes data types for the objects in the referenced model based on ranges collected through simulation or derived range analysis. If the system under design contains several instances of the same referenced model, the Fixed-Point Tool uses a union of the collected ranges for data type proposals.

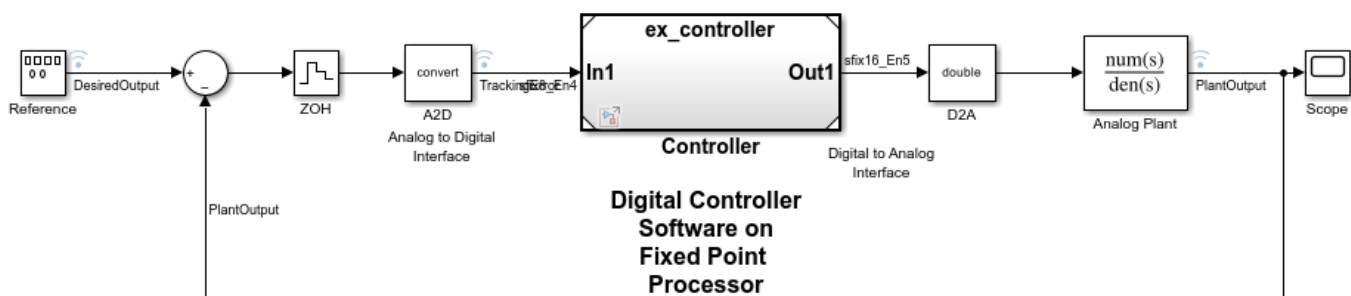
The Fixed-Point tool logs simulation minimum and maximum values only for instances of the referenced model that are in Normal mode. It does not log simulation minimum and maximum values for instances of the referenced model that are in non-Normal modes. If your model contains multiple instances of a referenced model and some are instances are in normal mode and some are not, the tool logs and displays data for those that are in normal mode.

Open ex_mdhref_controller Model

Open the ex_mdhref_controller model.

```
open_system('ex_mdhref_controller')
```

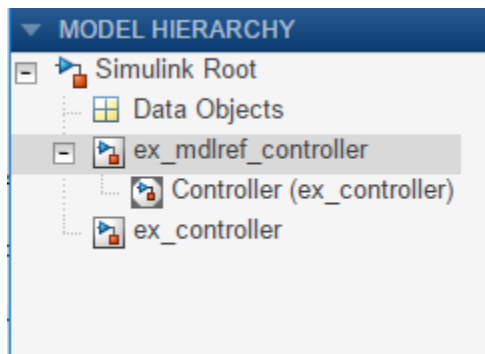
Scaling a Fixed-Point Control Design



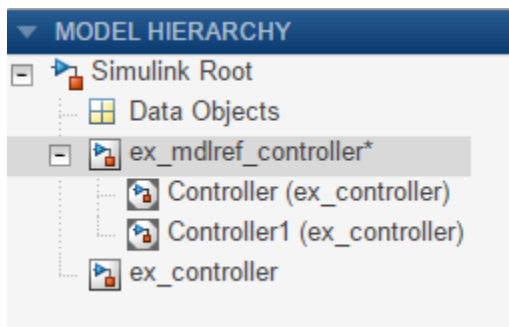
View Model Hierarchy in the Fixed-Point Tool

In the **Apps** gallery of the model, select **Fixed-Point Tool**.

When a model contains a referenced model, the Fixed-Point Tool **Model Hierarchy** pane displays a subnode for the instance of the referenced model and a node for the referenced model. For example, the ex_mdhref_controller model contains a Model block that references the ex_controller model. The Fixed-Point Tool shows both models in the model hierarchy.



If a model contains multiple instances of a referenced model, the tool displays each instance of the referenced model in this model and a node for the referenced model. For example, in the same model, if you duplicate the referenced model such that the `ex_mdref_controller` model contains two instances of the referenced model `ex_controller`. The Fixed-Point Tool displays both models and both instances of the referenced model in the model hierarchy.



Viewing Simulation Ranges for Referenced Models

- 1 In the Fixed-Point Tool, under **New**, select the Iterative Fixed-Point Conversion workflow.
- 2 Under **System Under Design (SUD)**, select the `ex_controller` model as the system you want to convert to fixed point.
- 3 Under **Range Collection Mode**, select **Simulation Ranges** as the range collection method.
- 4 In the toolstrip, click **Prepare**. The Fixed-Point Tool checks the system under design for compatibility with the conversion process and reports any issues found in the model. In this example, the tool reports that the model is ready for conversion.
- 5 Expand the **Collect Ranges** button arrow and select **Double precision**. Click the **Collect Ranges** button to start the simulation. The Fixed-Point Tool overrides the data types in the model with doubles and collects the minimum and maximum values for each object in your model that occur during the simulation. The Fixed-Point Tool stores this range information in a run titled **BaselineRun**.

The tool logs and displays the results for each instance of the referenced model. For example, here are the results for the first instance of the referenced model `ex_controller`.

The screenshot shows the Model Hierarchy on the left with 'ex_controller' selected. The Results table on the right displays the following data:

| Name | Run | CompiledDT | SpecifiedDT | SimMin | SimMax |
|-----------------------|----------------|------------|---------------------------|---------------------|--------------------|
| Combine Terms : A... | Ranges(Double) | double | Inherit: Inherit via i... | -7.009179391892659 | 4.864604758947099 |
| Combine Terms : O... | Ranges(Double) | double | fixdt(1,32,12) | -2.4143600572426305 | 4.864604758947099 |
| Denominator Terms... | Ranges(Double) | double | fixdt(1,32,12) | -9.574777620784765 | 6.066967386441957 |
| Denominator Terms... | Ranges(Double) | double | fixdt(1,32,12) | -7.009179391892659 | 3.4877081684371363 |
| Denominator Terms... | Ranges(Double) | double | fixdt(1,32,12) | -9.574777620784765 | 6.066967386441957 |
| Down Cast | Ranges(Double) | double | fixdt(1,16,5) | -2.4143600572426305 | 4.864604758947099 |
| Numerator Terms : ... | Ranges(Double) | double | fixdt(1,32,12) | -6.010590566449742 | 6.03517375895056 |
| Numerator Terms : ... | Ranges(Double) | double | fixdt(1,32,12) | -3.3599642646539447 | 3.546199714799863 |
| Numerator Terms : ... | Ranges(Double) | double | fixdt(1,32,12) | -6.010590566449742 | 6.03517375895056 |
| Up Cast | Ranges(Double) | double | fixdt(1,16,14) | -2.6853112577280545 | 4.23481996395132 |

Here are the results for the second instance of `ex_controller`.

The screenshot shows the Model Hierarchy on the left with 'Controller1 (ex_controller)' selected. The Results table on the right displays the following data:

| Name | Run | CompiledDT | SpecifiedDT | SimMin | SimMax |
|-----------------------|----------------|------------|---------------------------|----------------------|--------------------|
| Combine Terms : A... | Ranges(Double) | double | Inherit: Inherit via i... | -4.690670084764391 | 4.6853112577280545 |
| Combine Terms : O... | Ranges(Double) | double | fixdt(1,32,12) | -3.269951040092966 | 4.6853112577280545 |
| Denominator Terms... | Ranges(Double) | double | fixdt(1,32,12) | -9.221882475528233 | 6.4360941106591705 |
| Denominator Terms... | Ranges(Double) | double | fixdt(1,32,12) | -4.690670084764391 | 3.273064171832616 |
| Denominator Terms... | Ranges(Double) | double | fixdt(1,32,12) | -9.221882475528233 | 6.4360941106591705 |
| Down Cast | Ranges(Double) | double | fixdt(1,16,5) | -3.269951040092966 | 4.6853112577280545 |
| Numerator Terms : ... | Ranges(Double) | double | fixdt(1,32,12) | -9.436549978564978 | 9.475145275000827 |
| Numerator Terms : ... | Ranges(Double) | double | fixdt(1,32,12) | -0.08486130746613... | 0.1712008426003025 |
| Numerator Terms : ... | Ranges(Double) | double | fixdt(1,32,12) | -9.436549978564978 | 9.475145275000827 |
| Up Cast | Ranges(Double) | double | fixdt(1,16,14) | -4.359902080682029 | 6.648612943811886 |

In the referenced model node, the tool displays the union of the results for each instance of the referenced model.

The screenshot shows the Model Hierarchy on the left with 'ex_controller' selected. The Results table on the right displays the union of results from both instances:

| Name | Run | CompiledDT | SpecifiedDT | SimMin | SimMax |
|-----------------------|----------------|------------|---------------------------|---------------------|--------------------|
| Combine Terms : A... | Ranges(Double) | double | Inherit: Inherit via i... | -7.009179391892659 | 4.864604758947099 |
| Combine Terms : O... | Ranges(Double) | double | fixdt(1,32,12) | -3.269951040092966 | 4.864604758947099 |
| Denominator Terms... | Ranges(Double) | double | fixdt(1,32,12) | -9.574777620784765 | 6.4360941106591705 |
| Denominator Terms... | Ranges(Double) | double | fixdt(1,32,12) | -7.009179391892659 | 3.4877081684371363 |
| Denominator Terms... | Ranges(Double) | double | fixdt(1,32,12) | -9.574777620784765 | 6.4360941106591705 |
| Down Cast | Ranges(Double) | double | fixdt(1,16,5) | -3.269951040092966 | 4.864604758947099 |
| Numerator Terms : ... | Ranges(Double) | double | fixdt(1,32,12) | -9.436549978564978 | 9.475145275000827 |
| Numerator Terms : ... | Ranges(Double) | double | fixdt(1,32,12) | -3.3599642646539447 | 3.546199714799863 |
| Numerator Terms : ... | Ranges(Double) | double | fixdt(1,32,12) | -9.436549978564978 | 9.475145275000827 |
| Up Cast | Ranges(Double) | double | fixdt(1,16,14) | -4.359902080682029 | 6.648612943811886 |

Fixed-Point Instrumentation and Data Type Override Settings

When you simulate a model that contains referenced models, the data type override and fixed-point instrumentation settings for the top-level model do not control the settings for the referenced models. You must specify these settings separately for the referenced model. If the settings are inconsistent, for example, if you set the top-level model data type override setting to double and the referenced model to use local settings and the referenced model uses fixed-point data types, data type propagation issues might occur.

You can define custom data type override settings using `set_param`. For an example, see “Use Custom Data Type Override Settings for Range Collection” on page 44-9.

When you change the fixed-point instrumentation and data type override settings for any instance of a referenced model, the settings change on all instances of the model and on the referenced model itself.

Propose Data Types for a Referenced Model

- 1 In the **Convert** section of the toolstrip, click **Settings**. Specify the **Safety margin for simulation min/max (%)** parameter as 20.
- 2 Click **Propose Data Types**.

Because no design minimum and maximum information is supplied, the simulation minimum and maximum data that was collected during the simulation run is used to propose data types. The **Safety margin for simulation min/max (%)** parameter value multiplies the “raw” simulation values by a factor of 1.2. Setting the **Safety margin for simulation min/max (%)** parameter to a value greater than 1 decreases the likelihood that an overflow will occur when fixed-point data types are being used.

Because of the nonlinear effects of quantization, a fixed-point simulation produces results that are different from an idealized, doubles-based simulation. Signals in a fixed-point simulation can cover a larger or smaller range than in a doubles-based simulation. If the range increases enough, overflows or saturations could occur. A safety margin decreases the likelihood of this happening, but it might also decrease the precision of the simulation.

The Fixed-Point Tool analyzes the scaling of all fixed-point blocks whose **Lock output data type setting against changes by the fixed-point tools** parameter is not selected.

The Fixed-Point Tool uses the minimum and maximum values collected during simulation to propose a scaling for each block such that the precision is maximized while the full range of simulation values is spanned. The tool displays the proposed scaling in the spreadsheet.

The screenshot displays the MATLAB Fixed-Point Tool interface. The main window is titled "ITERATIVE FIXED-POINT CONVERSION" and "EXPLORE". The interface is divided into several panes:

- Workflow Browser:** Shows the current workflow steps: Setup, Preparation Results, and BaselineRun_2.
- Model Hierarchy:** Shows the Simulink Root, Data Objects, and the Controller (ex_control) block.
- Results Table:** A table listing simulation results for various blocks. The table has columns for Name, CompiledDT, SpecifiedDT, ProposedDT, Accept, SimMin, and SimMax. The "Accept" column contains checkmarks for most blocks, indicating that the proposed fixed-point scaling is accepted.
- Histograms of Simulation Data:** A plot showing the distribution of simulation data for various blocks. The plot is divided into regions for Overflows (red), Representable (white), In-Range (blue), and Underflows (yellow). The "In-Range" region is the largest, indicating that the proposed scaling is suitable for most of the data.
- Result Details:** A pane showing detailed information for a selected block, including the Proposed Data Type Summary and Ranges used for proposal.

| Name | CompiledDT | SpecifiedDT | ProposedDT | Accept | SimMin | SimMax |
|-------------------------------|------------|----------------|----------------|-------------------------------------|-----------------|-----------------|
| Combine Terms : Output | double | fixdt(1,32,12) | fixdt(1,32,28) | <input checked="" type="checkbox"/> | -2.413500903... | 4.3270018075... |
| Denominator Terms : Accum... | double | fixdt(1,32,12) | fixdt(1,32,27) | <input checked="" type="checkbox"/> | -8.516638478... | 5.3964875151... |
| Denominator Terms : Output | double | fixdt(1,32,12) | fixdt(1,32,28) | <input checked="" type="checkbox"/> | -6.475416336... | 3.4877081684... |
| Denominator Terms : Produ... | double | fixdt(1,32,12) | fixdt(1,32,27) | <input checked="" type="checkbox"/> | -8.516638478... | 5.3964875151... |
| Down Cast | double | fixdt(1,16,5) | fixdt(1,16,12) | <input checked="" type="checkbox"/> | -2.413500903... | 4.3270018075... |
| In1 | | fixdt(1,8,4) | fixdt(1,8,4) | <input type="checkbox"/> | | |
| Numerator Terms : Accumul... | double | fixdt(1,32,12) | fixdt(1,32,28) | <input checked="" type="checkbox"/> | -5.677304459... | 5.7005245184... |
| Numerator Terms : Output | double | fixdt(1,32,12) | fixdt(1,32,28) | <input checked="" type="checkbox"/> | -3.367372640... | 3.5439615259... |
| Numerator Terms : Product ... | double | fixdt(1,32,12) | fixdt(1,32,28) | <input checked="" type="checkbox"/> | -5.677304459... | 5.7005245184... |
| Out1 | | fixdt(1,16,5) | fixdt(1,16,12) | <input checked="" type="checkbox"/> | | |
| Up Cast | double | fixdt(1,16,14) | fixdt(1,16,12) | <input checked="" type="checkbox"/> | -2 | 3.9999999999... |

| Property | Proposed Data Type | Specified I |
|-----------|-----------------------|----------------|
| Data Type | fixdt(1,32,28) | fixdt(1,32,12) |
| Minimum | -8 | -524288 |
| Maximum | 7.99999999627471 | 524287.99975 |
| Precision | 3.725290298461914e... | 0.0002441406 |

| Property | Minimum | Maximum |
|------------|------------------|-----------------|
| Simulation | -5.6773044592... | 5.7005245184... |

| Values | Potential Overflows | In-Range | Potential Underflows |
|----------|---------------------|----------|----------------------|
| Positive | 0 | 825 | 170 |
| Negative | 0 | 835 | 155 |
| Zero | 0 | 10 | 0 |

- Review the scaling that the Fixed-Point Tool proposes. You can choose to accept the scaling proposal for each block by selecting the corresponding **Accept** check box. By default, the Fixed-Point Tool accepts all scaling proposals that differ from the current scaling. For this example, verify that the **Accept** check box is selected for each of the Controller system's blocks.

To view more information about a proposal, select the result and view the **Result Details** pane.

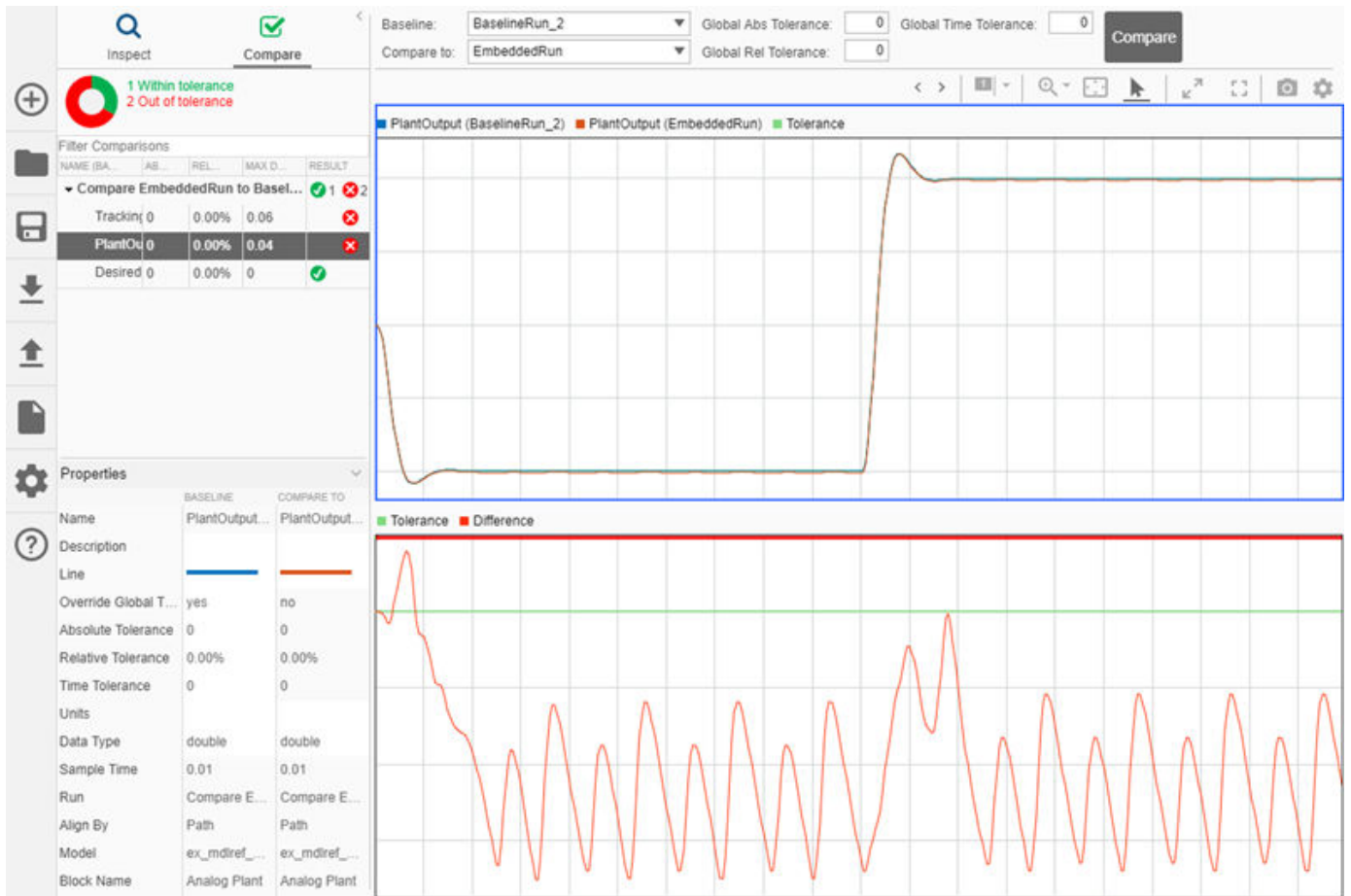
- In the Fixed-Point Tool, click the **Apply Data Types** button.

The Fixed-Point Tool applies the scaling proposals that you accepted in the previous step.

- In the **Verify** section of the toolbar, click the **Simulate with Embedded Types** button.

Simulink simulates the `ex_mdref_controller` model using the new scaling that you applied. Afterward, the Fixed-Point Tool displays information about blocks that logged fixed-point data.

- Click **Compare Results**. The Simulation Data Inspector plots the Analog Plant output for the floating-point and fixed-point runs and the difference between them.



See Also

Related Examples

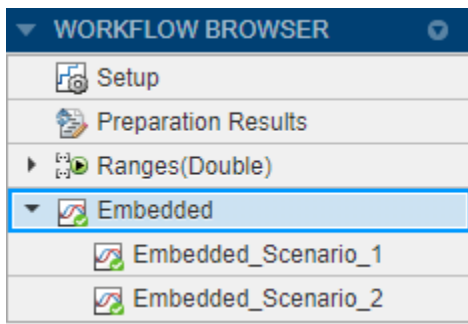
- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9

Control Views in the Fixed-Point Tool

The following sections describe how to control the amount of information that is shown in the Fixed-Point Tool. Within the Fixed-Point Tool, you can filter the results that are displayed at a given time by run, or by the subsystem in which the result belongs. You can also add or remove the columns that are shown in the spreadsheet of the Fixed-Point Tool.

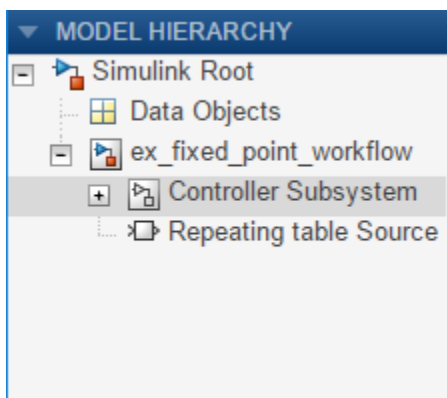
Filter Results by Run

In the Fixed-Point Tool, each time you collect ranges, either through simulation or derived range analysis, or run an embedded simulation, the tool stores the collected information in a run. Using the **Workflow Browser**, you can filter the results shown in the spreadsheet. Select the runs that you want to view in the spreadsheet.



Filter Results by Subsystem


By default, the Fixed-Point Tool displays only the results for model objects in your specified system under design. To filter the results shown or to see additional results, you can select a different node in the model hierarchy pane. The spreadsheet displays all model objects at and below the selected node.

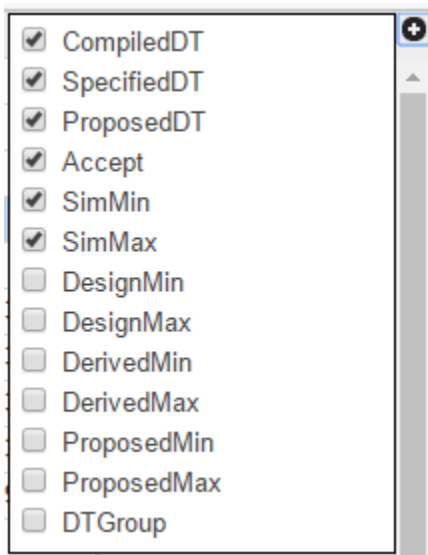


To view the data objects specified in your model, click the **Data Objects** node in the tree.

Control Column Views

As you follow the workflow for converting your model to fixed point using the Fixed-Point Tool, the tool displays the spreadsheet columns that are most pertinent to your current step in the workflow. If

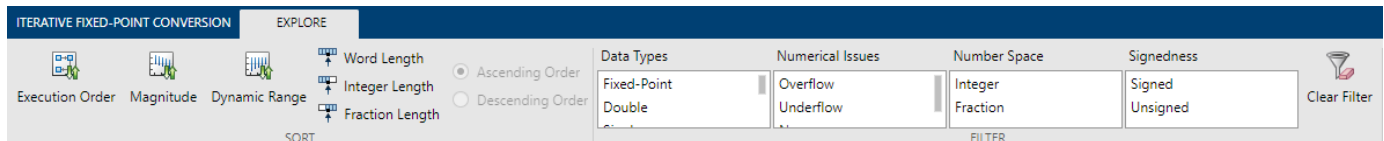
you want to view additional columns, you can add them to the spreadsheet using the  button in the top-right corner of the spreadsheet.



You can sort the results in the spreadsheet by clicking the header of the column on which you want to sort.

Use the Explore Tab to Sort and Filter Results

Using the **Explore** tab of the Fixed-Point Tool, you can sort and filter results in the tool based on additional criteria. The **Explore** tab is available whenever a single run that contains visualization data is selected in the **Run Browser** of the Fixed-Point Tool.



When there are unapplied proposals in the Fixed-Point Tool, the tool sorts and filters based on the **ProposedDT**. If there are no unapplied proposals, the tool sorts and filters the results based on the **SpecifiedDT**. If there is no **SpecifiedDT** available, for example if the result specifies an inherited data type, then the tool uses the **CompiledDT**.

Sort

You can sort results based on the following criteria.

| Sorting Criteria | Description |
|------------------------|--|
| Execution Order | Order in which the blocks are executed during simulation |
| Magnitude | The larger of the absolute values of the SimMin and SimMax |

| Sorting Criteria | Description |
|------------------------|--|
| Dynamic Range | Difference between the value representable by the largest bin and smallest bin in the histogram of logged values |
| Word Length | Total number of bits in the data type |
| Integer Length | Number of bits devoted to representing integer values in the data type The tool calculates the integer length as <i>word length</i> – <i>fraction length</i> . For example, for the data type <code>fixdt(1, 16, 12)</code> , the tool calculates the integer length as 4. The data type <code>fixdt(1, 16, -16)</code> would have a calculated integer length of 32. |
| Fraction Length | Number of bits in the fractional part of the data type |

By default, the tool sorts in ascending order. To have the tool sort in descending order, select **Descending**. You can select only one sorting criteria at a time.

Filter

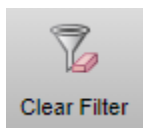
You can filter results based on the following criteria.

| Filter Criteria | Selections |
|-------------------------|--|
| Data Types | <ul style="list-style-type: none"> • Fixed-Point • Double • Single • Half • Boolean • Base Integer |
| Numerical Issues | <ul style="list-style-type: none"> • Overflow - Show only results containing an overflow • Underflow - Show only results containing an underflow • None - Show only results containing no underflows or overflows |
| Number Space | <ul style="list-style-type: none"> • Integer - Show only results for which the logged values are always integers. • Fraction - Show results for which the logged values contain a fractional part. |

| Filter Criteria | Selections |
|-------------------|--|
| Signedness | <ul style="list-style-type: none"> • Signed - Show only results for which the data type can represent signed values. For example, <code>int16</code>, <code>fixdt(1,16,12)</code>, and floating-point data types. • Unsigned - Show only results for which the data type can represent only unsigned values. For example, <code>uint16</code>, <code>Boolean</code>, <code>fixdt(0,16,12)</code> |

You can select multiple filtering selections from a criteria category by holding the **Ctrl** key while you select additional selections. When you select multiple filtering selections from the same criteria category, they are combined using OR logic. When you select multiple filtering selections across different criteria categories, they are combined using AND logic.

You can deselect a selections by holding the **Ctrl** key. To clear all filters, click the **Clear Filter**



button.

See Also

More About

- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9

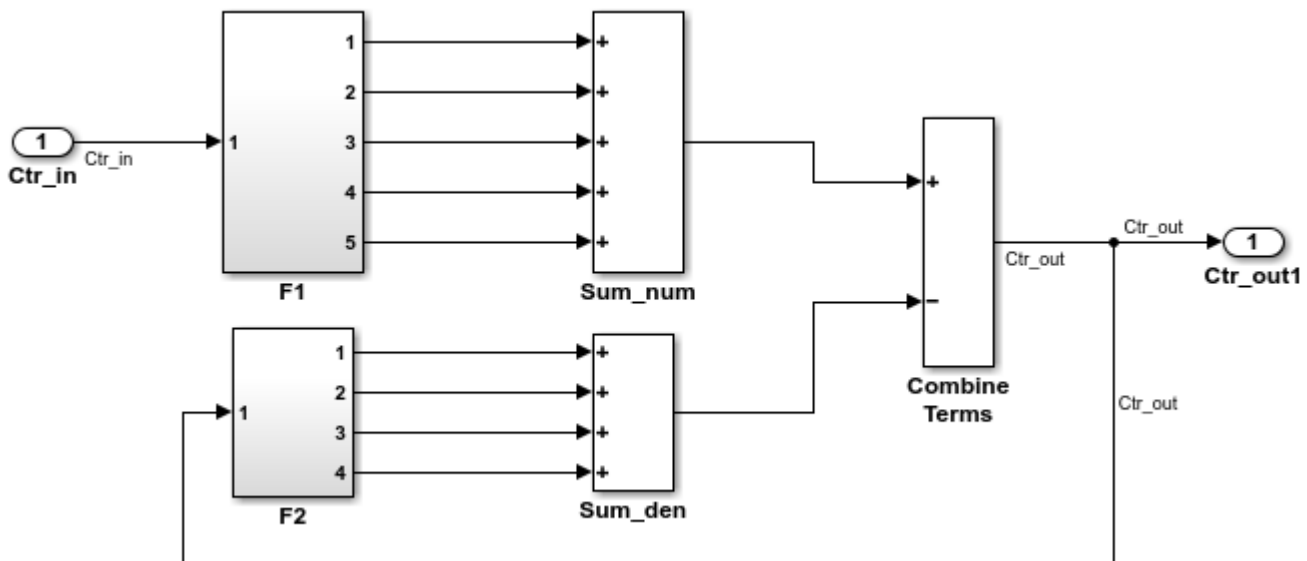
Model Multiple Data Type Behaviors Using a Data Dictionary

This example shows how to use referenced data dictionaries to store multiple sets of data types for a model. This example also shows how to change the data types by switching the referenced data dictionary.

Open the Model

Open the `ex_data_dictionary` model.

```
open_system('ex_data_dictionary')
```




The `ex_data_dictionary` model uses a data dictionary to store its data types.

- `mdl_dd.sldd` - Main data dictionary
- `flt_dd.sldd` - Referenced data dictionary using floating-point data types
- `fix_dd.sldd` - Referenced data dictionary using fixed-point data types

Explore How the Data Dictionary is Used in the Model

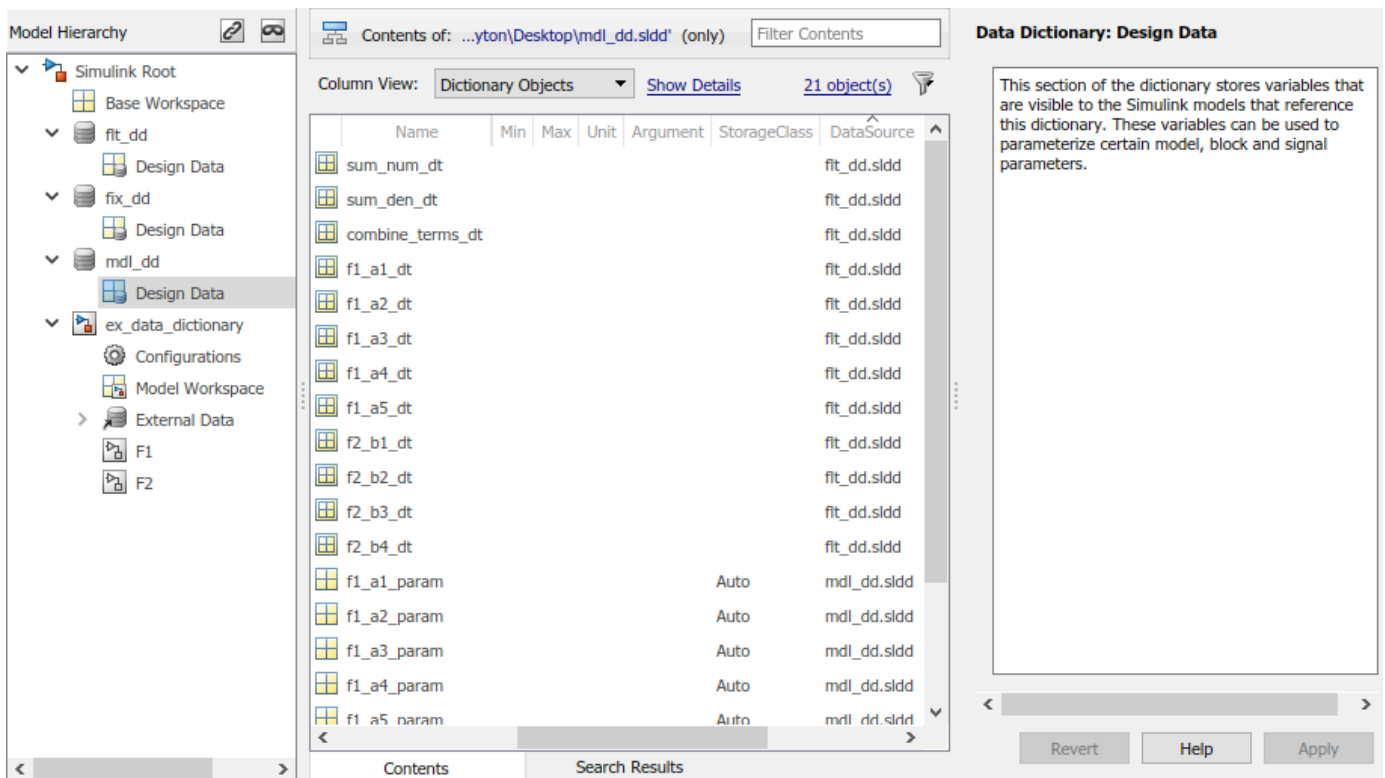
1 View the data dictionaries in the **Model Explorer**. On the **Modeling** tab, select **Model Explorer**.

2

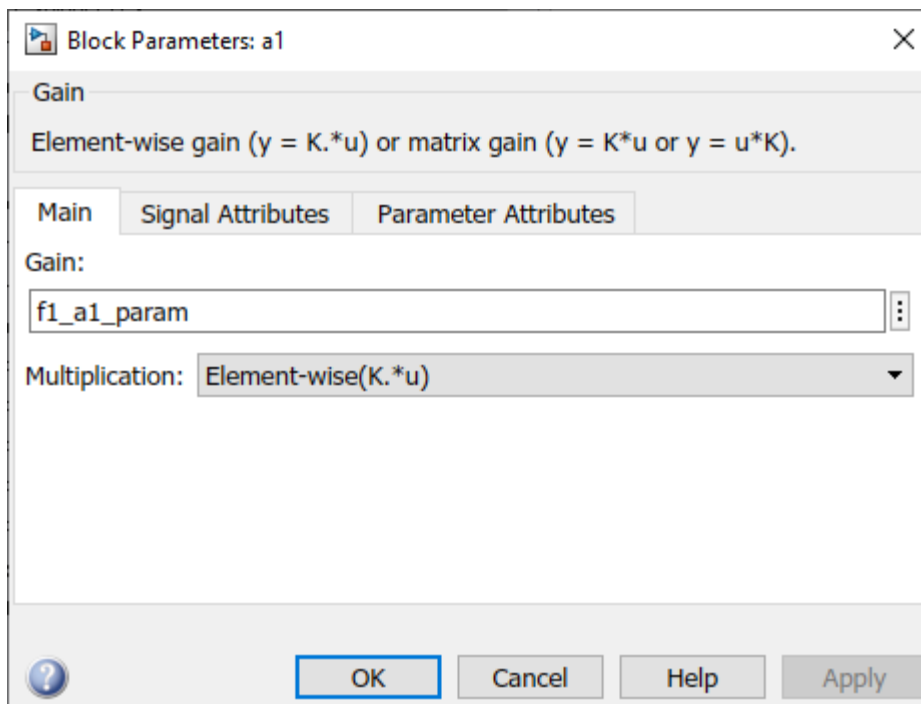
In the lower left corner of the Simulink Editor, click  to open the dictionary.

The data dictionary defines the parameters of the Gain blocks in the F1 and F2 subsystems. `mdl_dd` is associated with a referenced data dictionary, `flt_dd`, which defines the output data types of the gain blocks in the model's subsystems.

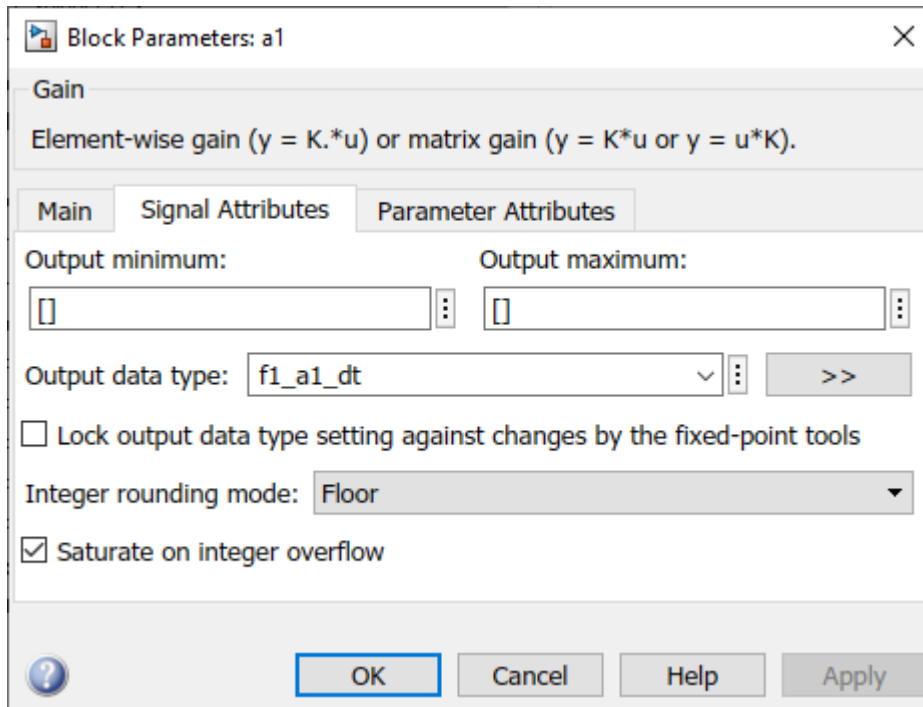
In the Model Explorer, in the **Contents** pane, the **Data Source** column shows the source data dictionary for each Gain block parameter.



- Return to the model. Open the F1 subsystem and double-click the a1 block. The block gain is specified as `f1_a1_param`, which is defined in the data dictionary.



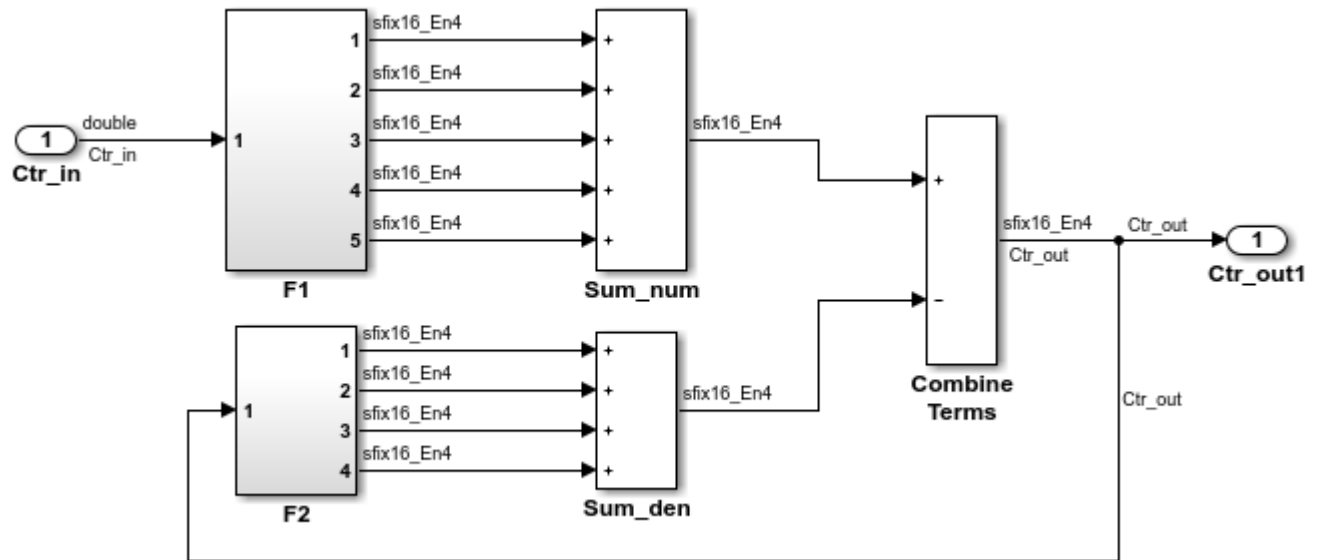
In the **Signal Attributes** tab, the block output data type is specified as `f1_a1_dt`. The data type of `f1_a1_dt` is defined in the referenced data dictionary, `flt_dd`.



Change Data Types of Model Parameters

The `fix_dd` data dictionary contains the same entries as `flt_dd`, but defines fixed-point data types instead of floating-point data types. To use the fixed-point data types without changing the model, replace `flt_dd` with `fix_dd` as the referenced data dictionary of `mdl_dd`.

- 1 In the Model Explorer, in the **Model Hierarchy** pane, right-click `mdl_dd` and select **Properties**.
- 2 Remove the referenced floating-point data dictionary. In the Data Dictionary dialog box, in the **Referenced Dictionaries** pane, select `flt_dd` and click **Remove**.
- 3 Add a reference to the fixed-point data dictionary. Click **Add** and select `fix_dd`. Click **OK** to close the dialog box.
- 4 In the Model Explorer, right-click `mdl_dd` and select **Save Changes**.
- 5 Return to the Simulink editor and update the model.



The model now uses fixed-point data types.

See Also

Related Examples

- “Migrate Single Model to Use Dictionary”

More About

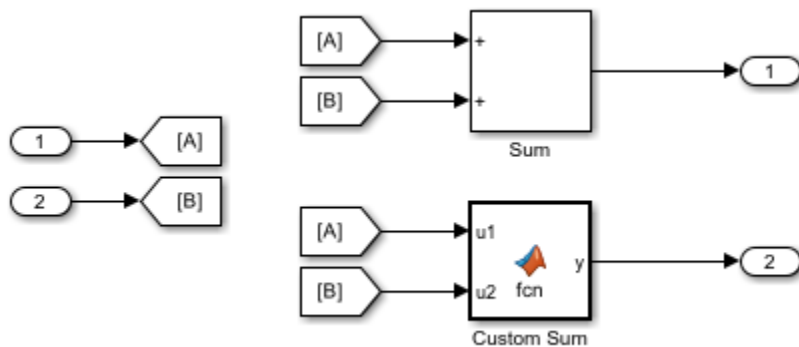
- “What Is a Data Dictionary?”

Compare Numerical Response of Sum Block and Sum in MATLAB Function Block

This example shows how to generate simulation inputs and use them to exercise a model over its full operating range. In this example, generate test data to simulate a model and compare the numerical response of the Sum block, and sum implemented in a MATLAB® Function block in the `ex_testsum` model.

Open the model.

```
model = 'ex_testsum';
open_system(model);
```



Specify Data Attributes and Generate Data

Use the `fixed.DataSpecification` object to specify input data attributes. In this example, create two `DataSpecification` objects, one with a double-precision data type, and the other using a single-precision data type. The interval of values generated by the first object is from 1 to 64, and the interval of values generated by the second is from 1 to 32.

```
dataspec1 = fixed.DataSpecification('double', 'Intervals', {1 64});
dataspec2 = fixed.DataSpecification('single', 'Intervals', {1 32});
```

The `DataGenerator` object generates combinations of numerically-rich values. To use the output data in a Simulink® model, set the format of the output to `'timeseries'`.

```
datagen = fixed.DataGenerator;
datagen.DataSpecifications = {dataspec1, dataspec2};
[tsdata1, tsdata2] = outputAllData(datagen, 'timeseries');
```

Set Up Model and Simulate

Apply the attributes of the `DataSpecification` objects to the Inport blocks in the model.

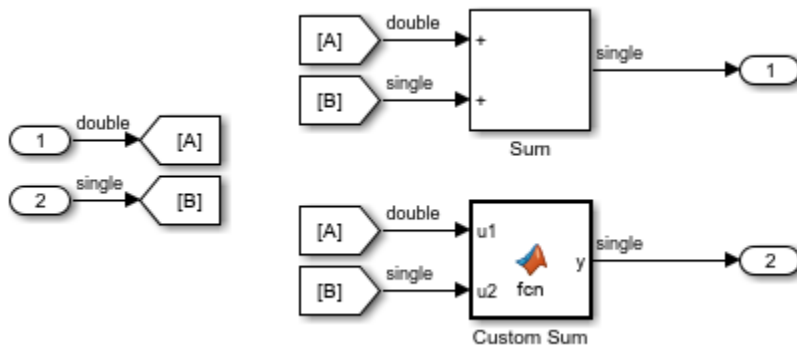
```
set_param('ex_testsum/In1', ...
    'OutDataTypeStr', dataspec1.DataTypeStr, ...
    'SignalType', dataspec1.Complexity, ...
    'PortDimensions', ['1' num2str(dataspec1.Dimensions) ''])
set_param('ex_testsum/In2', ...
    'OutDataTypeStr', dataspec2.DataTypeStr, ...
```

```
'SignalType',dataspec2.Complexity,...
'PortDimensions',{'[ ' num2str(dataspec2.Dimensions) ' ]']})
```

Load the generated timeseries data into the model and simulate.

```
set_param(model, 'LoadExternalInput', 'on',...
'ExternalInput', 'tsdata1, tsdata2',...
'StopTime', string(tsdata1.Time(end)));
```

```
simout = sim(model);
```



Visualize Output

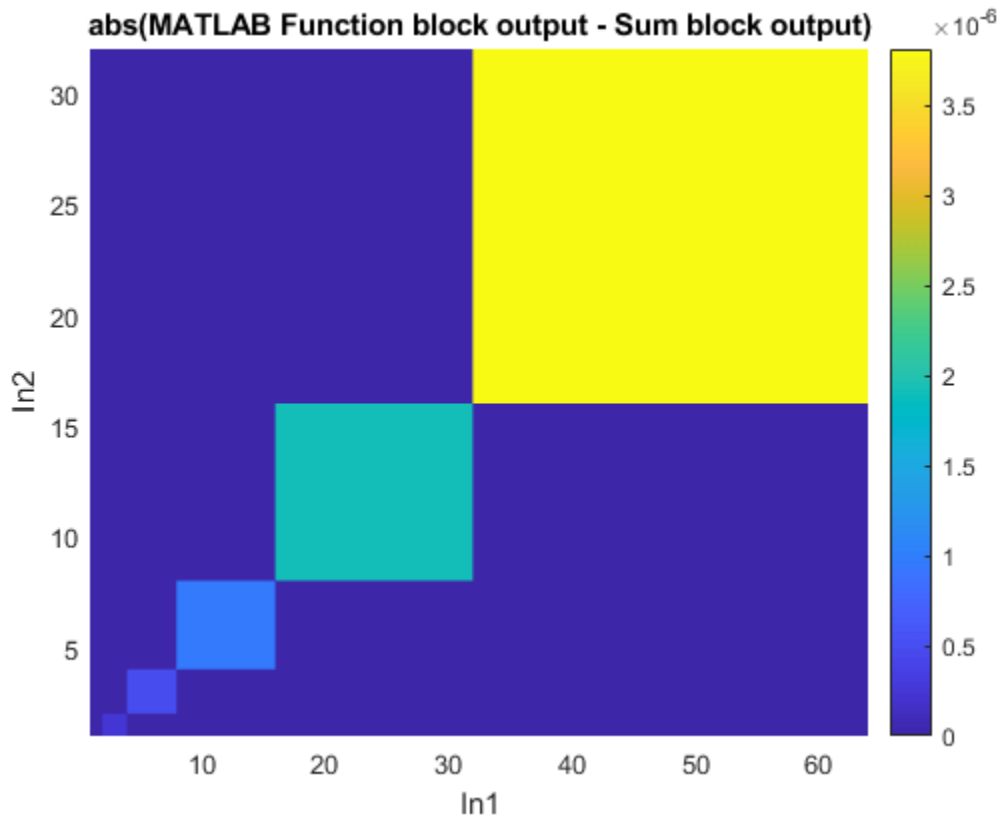
Visualize the output of the simulation, and compare the numerical behavior of the two implementations of the sum operation.

Get the unique values in the generated data for each set of data.

```
[x, y] = datagen.getUniqueValues;
d = abs(simout.yout{1}.Values.Data - simout.yout{2}.Values.Data);
X = reshape(tsdata1.Data, numel(x), []);
Y = reshape(tsdata2.Data, numel(x), []);
D = reshape(d, numel(x), []);
```

Plot the difference between outputs as a function of the input values.

```
figure;
surf(X, Y, D, 'EdgeColor', 'none');
grid on;
view(2);
axis tight;
xlabel('In1');
ylabel('In2');
colorbar;
title('abs(MATLAB Function block output - Sum block output)');
```

From the plot, you can see that the difference between the two implementations increases as the values of the numeric inputs get larger. This difference is due to the difference in the data type of the accumulator in the two implementations.

Compare Numerical Response with Single-Precision Accumulator

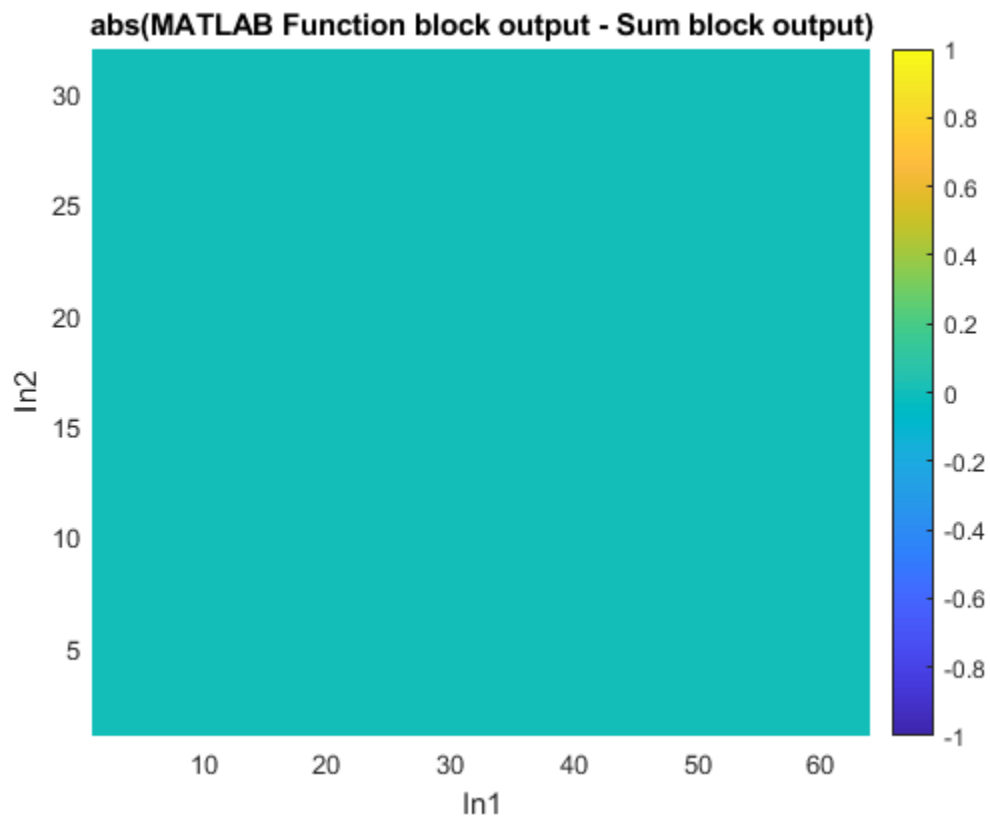
The **Accumulator Data Type** parameter of the Sum block is set to `Inherit`: `Inherit` via `internal` rule. In this case, the data type used for the accumulator is a double-precision floating-point type. Set the **Accumulator data type** to `single` and compare the output again.

```
set_param([model, '/Sum'], 'AccumDataTypeStr', 'single')
simout = sim(model);
```

Visualize the output. When the accumulator type of the Sum block is set to `single`, the implementations return the same result at all values.

```
[x, y] = datagen.getUniqueValues;
d = abs(simout.yout{1}.Values.Data - simout.yout{2}.Values.Data);
X = reshape(tpdata1.Data, numel(x), []);
Y = reshape(tpdata2.Data, numel(x), []);
D = reshape(d, numel(x), []);
figure;
surf(X, Y, D, 'EdgeColor', 'none');
grid on;
view(2);
axis tight;
xlabel('In1');
```

```
ylabel('In2');  
colorbar;  
title('abs(MATLAB Function block output - Sum block output)');
```



Convert Floating-Point Model to Fixed Point

- “Convert Floating-Point Model to Fixed Point” on page 40-2
- “Explore Multiple Floating-Point to Fixed-Point Conversions” on page 40-11
- “Optimize Fixed-Point Data Types for a System” on page 40-14
- “Optimize Data Types Using Multiple Simulation Scenarios” on page 40-20
- “Optimize Data Types for an FPGA with DSP Slices” on page 40-23
- “Use Data Type Optimization to Minimize Operator Counts” on page 40-30
- “Image Denoising Using Fixed-Point Quantized Restricted Boltzmann Machine Algorithm” on page 40-33
- “Optimize the Fixed-Point Data Types of a System Using the Fixed-Point Tool” on page 40-40
- “Perform Data Type Optimization with Custom Behavioral Constraints” on page 40-46

Convert Floating-Point Model to Fixed Point

In this section...

- “Set up the Model” on page 40-2
- “Prepare System for Conversion” on page 40-3
- “Collect Ranges” on page 40-5
- “Convert Data Types” on page 40-6
- “Verify New Settings” on page 40-7
- “Replace Unsupported Blocks with a Lookup Table Approximation” on page 40-8
- “Verify Behavior of System with Lookup Table Approximation” on page 40-10

In this example, learn how to:

- Convert a floating-point system to an equivalent fixed-point representation.

The Fixed-Point Tool automates the task of specifying fixed-point data types in a system. In this example, the tool collects range data for model objects, either from design minimum and maximum values that you specify explicitly for signals and parameters, or from logged minimum and maximum values that occur during simulation. Based on these values, the tool proposes fixed-point data types that maximize precision and cover the range. The tool allows you to review the data type proposals and then apply them selectively to objects in your model.

- Replace blocks that are not supported for conversion with a lookup table approximation.

During the preparation stage of the conversion, the Fixed-Point Tool isolates any blocks that do not support fixed-point conversion by placing these blocks inside a subsystem surrounded by Data Type Conversion blocks. You can use the Lookup Table Optimizer to replace the unsupported blocks with a lookup table approximation.

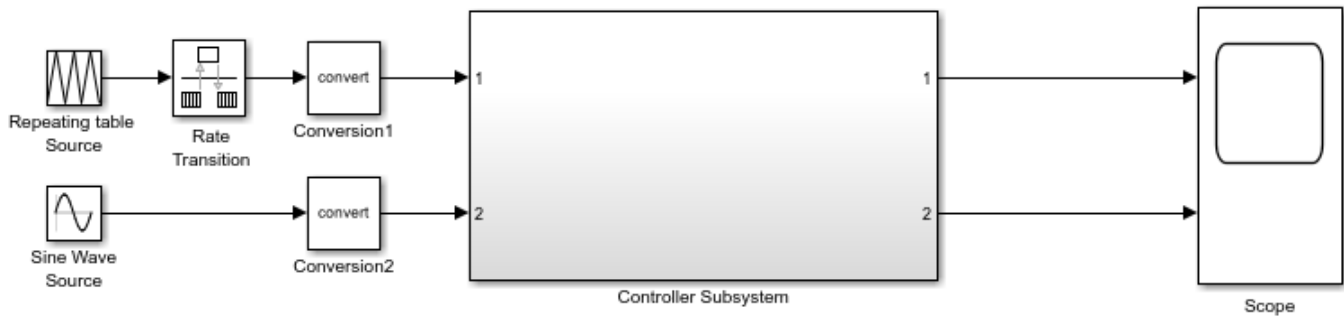
For a visual walk-through of these steps, watch the video:



Set up the Model

Open the model and configure it for fixed-point conversion.

```
open_system('ex_fixed_point_workflow')
```



The model consists of a source, a Controller Subsystem that you want to convert to fixed point, and a scope to visualize the subsystem outputs. Configuring a model in this way helps you to determine the effect of fixed-point data types on a system. Using this approach, you convert only the subsystem because this is the system of interest. There is no need to convert the source or scope to fixed point.

This configuration allows you to modify the inputs and collect simulation data for multiple stimuli. You can then examine the behavior of the subsystem with different input ranges and scale your fixed-point data types to provide maximum precision while accommodating the full simulation range.

To compare the behavior before and after conversion, enable signal logging at the outputs of the system under design.

```
ph = get_param('ex_fixed_point_workflow/Controller Subsystem', 'PortHandles');
set_param(ph.Outputport(1), 'DataLogging', 'on')
set_param(ph.Outputport(2), 'DataLogging', 'on')
```

Prepare System for Conversion

To convert the model to fixed point, use the Fixed-Point Tool.

- 1 In the **Apps** gallery of the `ex_fixed_point_workflow` model, select **Fixed-Point Tool**.
- 2 In the Fixed-Point Tool, under **New** workflow, select **Iterative Fixed-Point Conversion**.
- 3 Under **System Under Design (SUD)**, select the subsystem you want to convert to fixed point. In this example, select **Controller Subsystem**.
- 4 Under **Range Collection Mode**, select **Simulation Ranges** as the range collection method.
- 5 Under **Simulation Inputs**, you can specify `Simulink.SimulationInput` objects to exercise your design over its full operating range. In this example, set **Simulation inputs** to **Use default model inputs**.
- 6 To specify tolerances for the system, in the table under **Signal Tolerances**, specify tolerances for any signal in the model with signal logging enabled. For more information, see “Specify Signal Tolerances” on page 42-18.

Set the relative tolerance (**Rel Tol**) of the signals that you logged to 15%.

▼ Signal Tolerances

Specify tolerances for signals in your model that have signal logging enabled. After converting your system to fixed point, the Workflow Browser displays whether the embedded run meets the specified signal tolerances.

Filter signal list: Refresh Signals

| Signal Name | Abs Tol | Rel Tol | Time Tol (seconds) |
|------------------------|---------|---------|--------------------|
| Controller Subsystem:1 | | 0.15 | |
| Controller Subsystem:2 | | 0.15 | |

- 7 In the toolbar, click **Prepare**. The Fixed-Point Tool checks the system under design for compatibility with the conversion process and reports any issues found in the model. When possible, the Fixed-Point Tool automatically changes settings that are not compatible. For more information, see “Use the Fixed-Point Tool to Prepare a System for Conversion” on page 38-2.

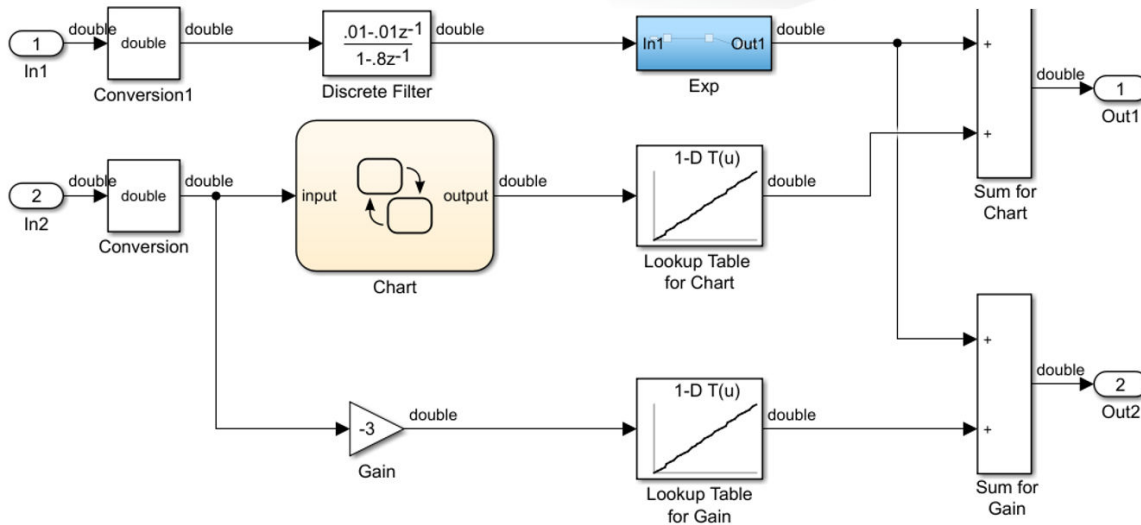
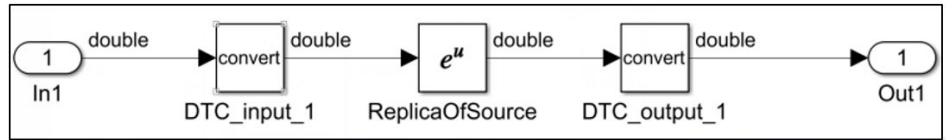
The screenshot shows the Fixed-Point Tool interface. The top toolbar includes buttons for New, Prepare, Collect Ranges, MATLAB Functions, Propose Data Types, Apply Data Types, Simulate with Embedded Types, and Restore Original Model. The main workspace displays the 'Selected system under design: ex_fixed_point_workflow/Controller Subsystem'. A 'Progress' indicator shows 100% completion. A table lists checks with their status:

| Selection | Check | Status |
|----------------------------------|-------------------------------------|--------|
| <input checked="" type="radio"/> | Create Restore Point | ✓ |
| <input type="radio"/> | Hardware Implementation Consistency | ✓ |
| <input type="radio"/> | Diagnostic Settings | ✓ |
| <input type="radio"/> | Unsupported Constructs | ✓ |
| <input type="radio"/> | System Under Design Boundary | ✓ |

Preparation is complete for the selected system under design.


The 'Preparation Details' panel on the right provides instructions: 'To ensure your original design is saved before making fixed-point data type changes, create a restore point for the model.' It also includes a 'Check Status' section stating: 'A restore point has previously been created for this model. To restore the model to this state, click the Restore Original Model button.'

The subsystem under design contains an Exp block, which does not support fixed-point data types. The Fixed-Point Tool surrounds this block with Data Type Conversion blocks and places it inside a subsystem. When you finish converting the rest of the subsystem to fixed point, you can replace the subsystem with a lookup table approximation of the exp function.



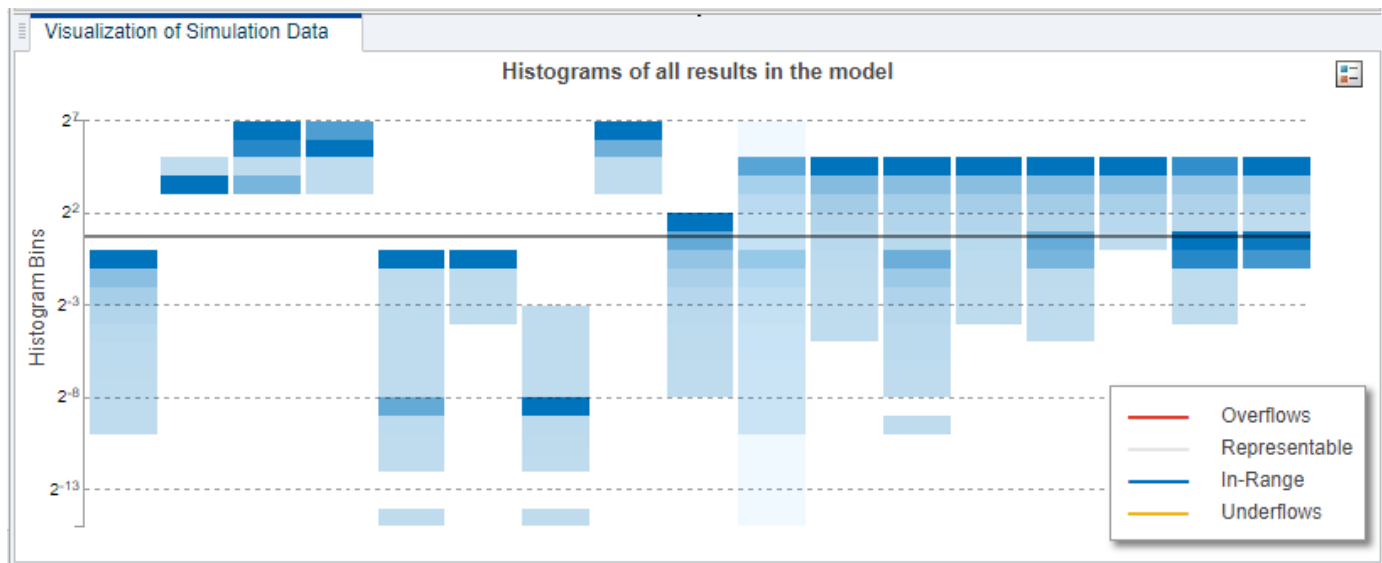
Collect Ranges

By default, the Fixed-Point Tool uses the current data type override set on the model. For this example, override data types in the model with double precision for the range collection run.

- 1 Expand the **Collect Ranges** button arrow and select **Double precision**.
- 2 Click **Collect Ranges**  to simulate the model.

The Fixed-Point Tool overrides the data types in the model with doubles and collects the minimum and maximum values for each object in your model that occur during the simulation. The Fixed-Point Tool stores this range information in a run titled **BaselineRun**. You can view the collected ranges in the **SimMin** and **SimMax** columns of the spreadsheet, or in the **Result Details** pane.

The **Visualization of Simulation Data** pane offers another view of the simulation results. Select the **Explore** tab of the Fixed-Point Tool for additional tools for sorting and filtering the data in the spreadsheet and the visualization.



Convert Data Types

Use the Fixed-Point Tool to propose fixed-point data types for appropriately configured blocks based on the double-precision simulation results stored in the run `BaselineRun`.

- 1 In the **Convert** section of the toolstrip, click the **Propose Data Types** button.

The Fixed-Point Tool analyzes the scaling of all fixed-point blocks whose **Lock output data type setting against changes by the fixed-point tools** parameter is not selected.


The Fixed-Point Tool uses the default proposal settings to propose data types with 16-bit word length and best-precision fraction length and updates the results in the spreadsheet.

You can edit the proposal settings by clicking the **Settings** button in the **Convert** section of the toolstrip before proposing types.

- 2 The tool displays the proposed data types in the **ProposedDT** column in the spreadsheet.

By default, it selects the **Accept** check box for each result where the proposed data type differs from the current data type of the object. If you apply data types, the tool applies these proposed data types to the system under design.

- 3 Examine the results to resolve any issues and to ensure that you want to accept the proposed data type for each result. The **Visualization of Simulation Data** pane indicates results that would contain overflows or underflows with a red or yellow triangle, respectively. Underflows can be sources of numerical issues, but can sometimes be safely ignored.

The Fixed-Point Tool indicates results whose proposed data type conflicts with another type with a red icon . In this example, no results contain conflicts. For more information, see "Examine Results to Resolve Conflicts" on page 42-26.

- 4 After reviewing the results and ensuring that there are no issues, you are ready to apply the proposed data types to the model. Click **Apply Data Types** to write the proposed data types to the model.

The Fixed-Point Tool applies the data type proposals to the blocks in the system under design.

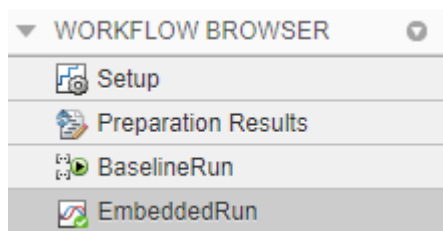
Verify New Settings

Next, simulate the model again using the new fixed-point settings. You then use the Simulation Data Inspector plotting capabilities to compare the results from the floating-point `BaselineRun` run with the fixed-point results.

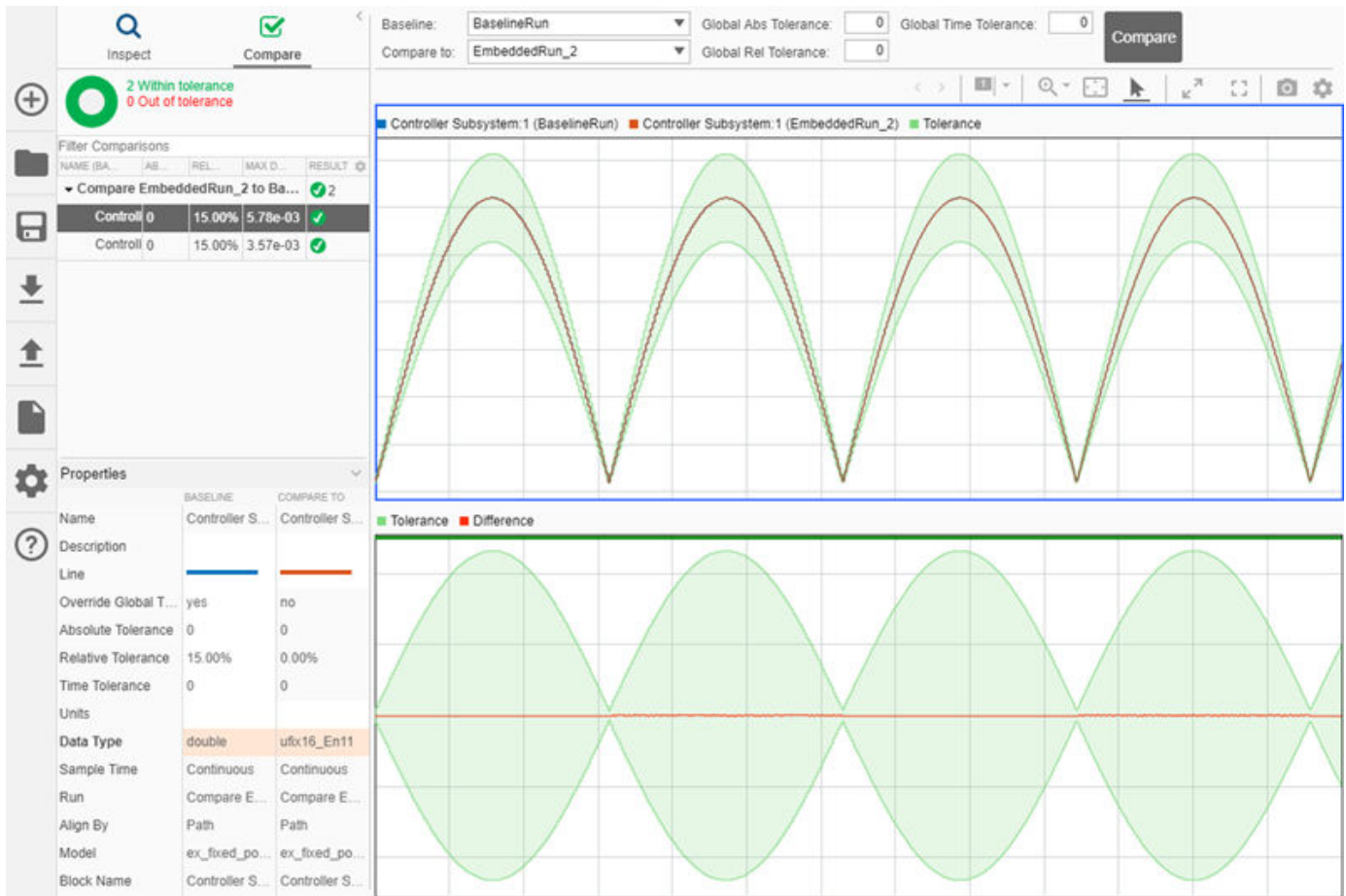
- 1 Click **Simulate with Embedded Types**. The Fixed-Point Tool simulates the model using the new fixed-point data types and stores the run information in a new run titled `EmbeddedRun`.

Afterward, the Fixed-Point Tool displays information about blocks that logged fixed-point data. The **CompiledDT** column for the run shows that the Controller Subsystem blocks use the new fixed-point data types.

- 2 Examine the histograms in the **Visualization of Simulation Data** pane to verify that there are no overflows or saturations. Overflows and saturations are marked with a red triangle ▲.
- 3 The workflow browser indicates that all signals for which you specified tolerances passed.



- 4 Click **Compare Results** to open the Simulation Data Inspector. In the Simulation Data Inspector, select one of the logged signals to view the fixed-point simulation behavior.



Replace Unsupported Blocks with a Lookup Table Approximation

In the “Prepare System for Conversion” on page 40-3 step of the workflow, the Fixed-Point Tool placed the Exp block, which is not supported for conversion, inside a subsystem surrounded with Data Type Conversion blocks. In this step, you replace the subsystem with a lookup table approximation.

- 1 To get a list of all of the subsystems the Fixed-Point Tool decoupled for conversion, at the command line enter:

```
decoupled = DataTypeWorkflow.findDecoupledSubsystems('ex_fixed_point_workflow')
decoupled =
```

```
1x2 table
```

| ID | BlockPath |
|----|--|
| 1 | {'ex_fixed_point_workflow/Controller Subsystem/Exp'} |

The `DataTypeWorkflow.findDecoupledSubsystems` function returns a table containing the block path of any subsystems that were created by the Fixed-Point Tool to isolate an unsupported block.

- 2 Open the **Lookup Table Optimizer**. In the **Apps** gallery, select **Lookup Table Optimizer**.
- 3 On the **Objective** page of the Lookup Table Optimizer, select **Simulink block or subsystem**. Click **Next**.
- 4 Under **Block Information**, copy from the command line and paste the path to the subsystem created by the Fixed-Point Tool.
- 5 Click the **Collect Current Values from Model** button to update the model diagram and allow the Lookup Table Optimizer to automatically gather information needed for the optimization process. Click **Next**.

Block Information

Simulink Block Path

'ex_fixed_point_workflow/Controller Subsystem/Exp'

Get Current Block

Attributes of Memory Efficient LUT

Collect Current Values from Model

Desired Output Data Type numerictype(0,16,15)

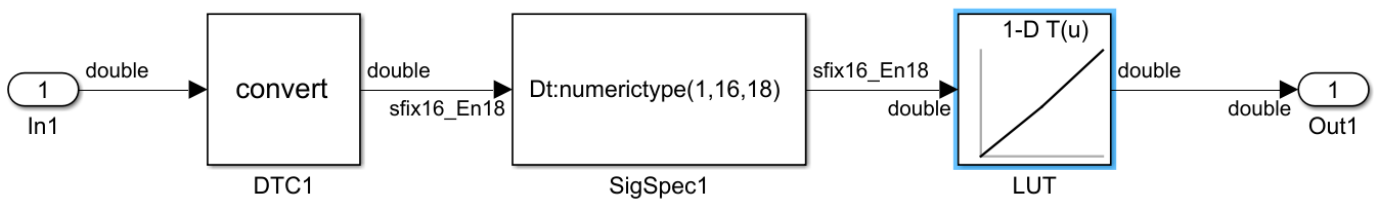
| Input | Desired Data Type | Minimum | Maximum |
|-------|----------------------|---------|---------------------|
| 1 | numerictype(1,16,18) | -0.125 | 0.12499618530273438 |

automatically gather information needed for the optimization process including current output data type, number and data type of inputs, and ranges of input values. You can manually edit all of these fields to specify ranges and data types other than those currently specified on the block.

- Specify the **Desired Output Data Type** of the generated lookup table in the form `numerictype(signedness, wordlength, fractionlength)`. For example, to specify a signed output data type of 16-bit word length and 8-bit fraction length, enter `numerictype(1,16,8)`. See `numerictype` for more information.
- Each input to the block being replaced represents a dimension of the replacement lookup table. Specify the minimum and maximum values of each dimension of the generated lookup table as scalars in the table.
- Specify the data type of each input to the block in the form `numerictype(signedness, wordlength, fractionlength)`.

Back Next

- 6 Specify the constraints to use in the optimization. For this example, use the default values. To create the lookup table, click **Optimize**. Click **Next**.
- 7 Click **Replace Original Function**. The Lookup Table Optimizer replaces the Math Function `exp` block with a new variant subsystem containing the lookup table approximation.



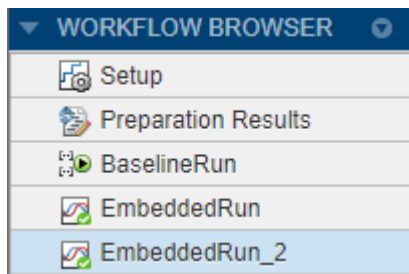
Verify Behavior of System with Lookup Table Approximation

Now that the system under design is fully converted, verify that the system still meets the tolerances you specified before conversion.

- 1 In the Fixed-Point Tool, in the **Verify** section of the toolbar, click **Simulate with Embedded Types**.

The Fixed-Point Tool simulates the model, which now contains the lookup table approximation, and saves the result as `EmbeddedRun_2`.

- 2 The **Workflow Browser** shows that the signals with specified tolerances pass in the model using the lookup table approximation.



See Also

More About

- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9
- “Optimize Lookup Tables for Memory-Efficiency” on page 41-15

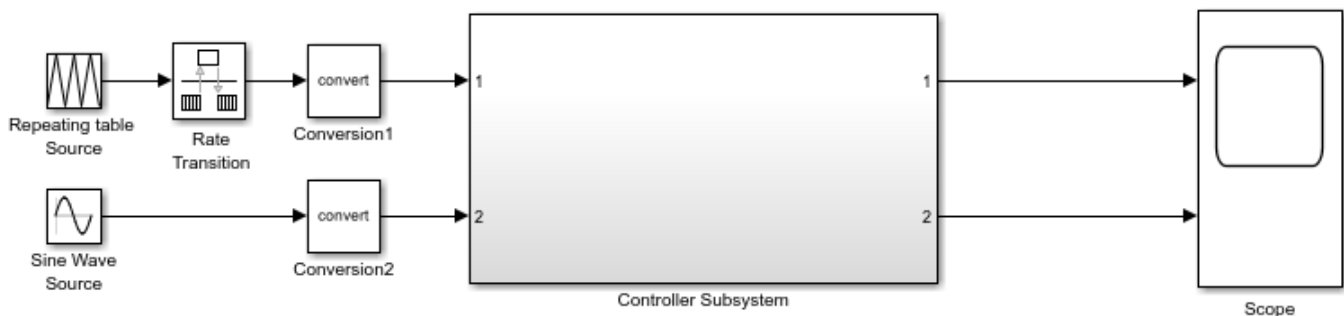
Explore Multiple Floating-Point to Fixed-Point Conversions

In this example, you use the Fixed-Point Tool and the `ex_fixed_point_workflow` model to explore different word length choices. After you simulate your model using embedded types, and compare the floating point and fixed-point behavior of your system, determine if the new behavior is satisfactory. If the behavior of the system using the newly applied fixed-point data types is not acceptable, you can iterate through the process until you find settings that work for your system.

Set up the Model

Open the model and configure it for fixed-point conversion.

```
open_system('ex_fixed_point_workflow')
```



The model consists of a source, a Controller Subsystem that you want to convert to fixed point, and a scope to visualize the subsystem outputs. Configuring a model in this way helps you to determine the effect of fixed-point data types on a system. Using this approach, you convert only the subsystem because this is the system of interest. There is no need to convert the source or scope to fixed point.

This configuration allows you to modify the inputs and collect simulation data for multiple stimuli. You can then examine the behavior of the subsystem with different input ranges and scale your fixed-point data types to provide maximum precision while accommodating the full simulation range.

To compare the behavior before and after conversion, enable signal logging at the outputs of the system under design.

```
ph = get_param('ex_fixed_point_workflow/Controller Subsystem', 'PortHandles');
set_param(ph.Outputport(1), 'DataLogging', 'on')
set_param(ph.Outputport(2), 'DataLogging', 'on')
```

Convert to Fixed-Point Using Default Proposal Settings

- 1 In the **Apps** gallery of the `ex_fixed_point_workflow` model, select **Fixed-Point Tool**.
- 2 In the Fixed-Point Tool, under **System Under Design**, select the subsystem you want to convert to fixed point. In this example, select **Controller Subsystem**.
- 3 Under **Range Collection Mode**, select **Simulation Ranges** as the range collection method.
- 4 Under **Simulation Inputs**, you can specify `Simulink.SimulationInput` objects to exercise your design over its full operating range. In this example, set **Simulation inputs** to Use default model inputs.

- To specify tolerances for the system, in the table under **Signal Tolerances**, specify tolerances for any signal in the model with signal logging enabled.


Set the relative tolerance (**Rel Tol**) of the signals that you logged to 15%.


▼ Signal Tolerances

Specify tolerances for signals in your model that have signal logging enabled. After converting your system to fixed point, the Workflow Browser displays whether the embedded run meets the specified signal tolerances.


Filter signal list: Refresh Signals


| Signal Name | Abs Tol | Rel Tol | Time Tol (seconds) |
|------------------------|---------|---------|--------------------|
| Controller Subsystem:1 | | 0.15 | |
| Controller Subsystem:2 | | 0.15 | |

- In the toolbar, click **Prepare**. The Fixed-Point Tool checks the system under design for compatibility with the conversion process and reports any issues found in the model. When possible, the Fixed-Point Tool automatically changes settings that are not compatible. For more information, see “Use the Fixed-Point Tool to Prepare a System for Conversion” on page 38-2.
- Expand the **Collect Ranges** button arrow and select **Double precision**. Click the **Collect Ranges** button  to override data types in the model with double precision and start the range collection simulation.

- In the **Convert** section, click the **Propose Data Types** button .

The Fixed-Point Tool uses the default proposal settings to propose data types with 16-bit word length and best-precision fraction length and updates the results in the spreadsheet.


- Click the **Apply Data Types** button  to write the proposed data types to the model.

- In the **Verify** section of the toolbar, click the **Simulate with Embedded Types** button . The Fixed-Point Tool simulates the model using the new fixed-point data types and stores the run information in a new run titled **EmbeddedRun**.

- Click **Compare Results** to open the Simulation Data Inspector and compare the floating-point and fixed-point behavior.

Return to the Fixed-Point Tool to update the proposal settings and generate new data type proposals.

Convert Using New Proposal Settings

- In the Fixed-Point Tool, in the **Convert** section of the toolbar, click the **Settings** button .
 - Edit the proposal settings to determine if a larger word length improves the fixed-point behavior of the system. Set the **Default Word Length** to 32.
 - To generate new proposals, click **Propose Data Types**.
 - Click the **Apply Data Types** to write the newly proposed data types to the model.
 - Click **Simulate with Embedded Types**. The Fixed-Point Tool simulates the model using the new fixed-point data types and stores the run information in a new run titled **EmbeddedRun_2**.
 - Click **Compare Results** to open the Simulation Data Inspector and compare the floating-point and fixed-point behavior.

You can continue to adjust the data type proposal settings, propose data types, and apply data types to your model until you find settings for which the fixed-point behavior of your system is acceptable.

See Also

More About

- “Explore Additional Data Types” on page 42-34

Optimize Fixed-Point Data Types for a System

Data type optimization seeks to minimize an objective function, such as the total bit-width or an estimated count of operators in generated code for a specified system, while maintaining original system behavior within a specified tolerance. During the optimization, the software establishes a baseline by simulating the original model. It then constructs different fixed-point versions of your model and runs simulations to determine the behavior using the new data types. The optimization selects the model that minimizes the objective function while also meeting the specified behavioral constraints.

The model containing the system you want to optimize must have the following characteristics:

- All blocks in the model must support fixed-point data types.
- The design ranges specified on blocks in the model must be consistent with the simulation ranges.
- If the model contains a MATLAB Function block, it must use MATLAB language features supported for fixed-point conversion. For more information, see “MATLAB Language Features Supported for Automated Fixed-Point Conversion” on page 7-35.
- The model must have finite simulation stop time.

During the optimization process, the software makes changes to several settings and model configuration parameters. The purpose of these changes include suppressing diagnostics, enabling logging with the Simulation Data Inspector, reducing the memory consumed by the result, ensuring validity of the model, accelerating the optimization process, and turning off data type override. For more information, see “Model Configuration Changes Made During Data Type Optimization” on page 42-63. You can restore these diagnostics after the optimization is complete.

Best Practices for Optimizing Data Types

Define Constraints

To determine if the behavior of a new fixed-point implementation is acceptable, the optimization requires well-defined behavioral constraints. To define a constraint, use the `addTolerance` method of the `fxpOptimizationOptions` object, or use one or more “Model Verification” blocks in your model. For more information, see “Specify Behavioral Constraints” on page 42-18.

Minimize Locked Data Types

When the **Lock data types against changes by the fixed-point tools** setting of a block within the system you want to optimize is enabled, it minimizes the freedom of the optimization process to find new solutions.

Model Management and Exploration

The `fxpopt` function returns an `OptimizationResult` object containing a series of fixed-point implementations called solutions. If the optimization process finds a fixed-point implementation that meets the specified behavioral constraints, the solutions are sorted by cost, giving the best solution with the smallest cost (bit width or operator counts) as the first element of the array.

In cases where the optimization is not able to find a fixed-point implementation meeting the behavioral constraints, the solutions are ordered by maximum absolute difference from the baseline model, with the smallest difference as the first element.

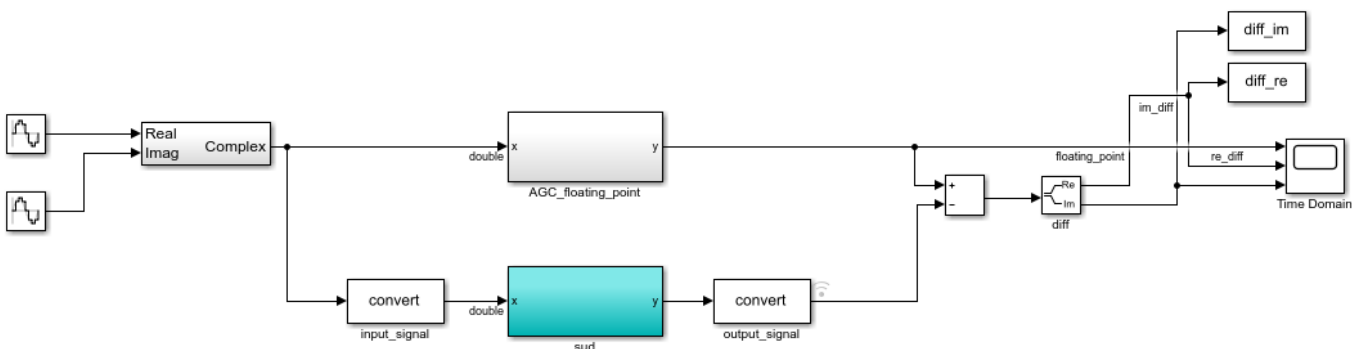
Explore the best found solution using the `explore` method of the `OptimizationResult` object. You can also explore any of the other found solutions in the same manner. Do not save and close the model until you select the solution you want to keep. Closing or saving the model inhibits further exploration of different solutions.

Optimize Fixed-Point Data Types

This example shows how to optimize the data types used by a system based on specified tolerances.

To begin, open the system for which you want to optimize the data types.

```
model = 'ex_auto_gain_controller';
sud = 'ex_auto_gain_controller/sud';
open_system(model)
```



Copyright 2017 The MathWorks, Inc.

Create an `fxpOptimizationOptions` object to define constraints and tolerances to meet your design goals. Set the `UseParallel` property of the `fxpOptimizationOptions` object to `true` to run iterations of the optimization in parallel. You can also specify word lengths to allow in your design through the `AllowableWordLengths` property.

```
opt = fxpOptimizationOptions('AllowableWordLengths', 10:24, 'UseParallel', true)
```

```
opt =
```

```
fxpOptimizationOptions with properties:
```

```

    MaxIterations: 50
         MaxTime: 600
        Patience: 10
       Verbosity: High
AllowableWordLengths: [10 11 12 13 14 15 16 17 18 19 20 21 22 23 24]
         UseParallel: 1
```

```
Advanced Options
  AdvancedOptions: [1x1 struct]
```

Use the `addTolerance` method to define tolerances for the differences between the original behavior of the system, and the behavior using the optimized fixed-point data types.

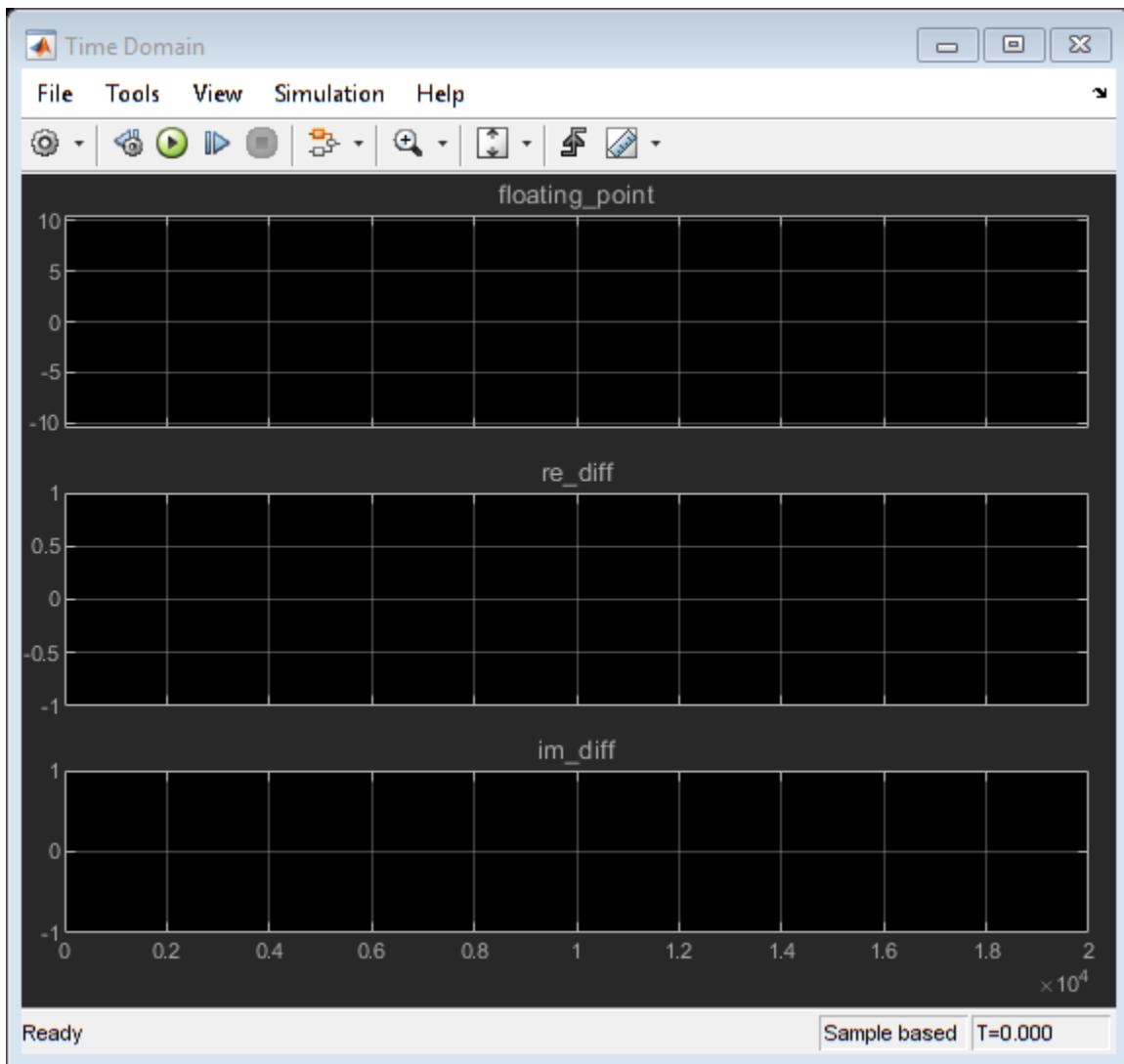
```
tol = 10e-2;
addTolerance(opt, [model '/output_signal'], 1, 'AbsTol', tol);
```

Use the `fxpopt` function to run the optimization. The software analyzes ranges of objects in your system under design and the constraints specified in the `fxpOptimizationOptions` object to apply heterogeneous data types to your system while minimizing total bit width.

```
result = fxpopt(model, sud, opt);
```

```
Starting parallel pool (parpool) using the 'local' profile ...
Connected to the parallel pool (number of workers: 4).
+ Preprocessing
+ Modeling the optimization problem
  - Constructing decision variables
+ Running the optimization solver
Analyzing and transferring files to the workers ...done.
- Evaluating new solution: cost 180, does not meet the tolerances.
- Evaluating new solution: cost 198, does not meet the tolerances.
- Evaluating new solution: cost 216, does not meet the tolerances.
- Evaluating new solution: cost 234, does not meet the tolerances.
- Evaluating new solution: cost 252, does not meet the tolerances.
- Evaluating new solution: cost 270, does not meet the tolerances.
- Evaluating new solution: cost 288, does not meet the tolerances.
- Evaluating new solution: cost 306, meets the tolerances.
- Evaluating new solution: cost 324, meets the tolerances.
- Evaluating new solution: cost 342, meets the tolerances.
- Evaluating new solution: cost 360, meets the tolerances.
- Evaluating new solution: cost 378, meets the tolerances.
- Evaluating new solution: cost 396, meets the tolerances.
- Evaluating new solution: cost 414, meets the tolerances.
- Evaluating new solution: cost 432, meets the tolerances.
- Updated best found solution, cost: 306
- Evaluating new solution: cost 304, meets the tolerances.
- Evaluating new solution: cost 304, meets the tolerances.
- Evaluating new solution: cost 301, meets the tolerances.
- Evaluating new solution: cost 305, does not meet the tolerances.
- Evaluating new solution: cost 305, meets the tolerances.
- Evaluating new solution: cost 301, meets the tolerances.
- Evaluating new solution: cost 299, meets the tolerances.
- Evaluating new solution: cost 299, meets the tolerances.
- Evaluating new solution: cost 296, meets the tolerances.
- Evaluating new solution: cost 299, meets the tolerances.
- Evaluating new solution: cost 291, meets the tolerances.
- Evaluating new solution: cost 296, does not meet the tolerances.
- Evaluating new solution: cost 299, meets the tolerances.
- Evaluating new solution: cost 300, meets the tolerances.
- Evaluating new solution: cost 296, does not meet the tolerances.
- Evaluating new solution: cost 301, meets the tolerances.
- Evaluating new solution: cost 303, meets the tolerances.
- Evaluating new solution: cost 299, meets the tolerances.
```

- Evaluating new solution: cost 304, does not meet the tolerances.
- Evaluating new solution: cost 300, meets the tolerances.
- Updated best found solution, cost: 304
- Updated best found solution, cost: 301
- Updated best found solution, cost: 299
- Updated best found solution, cost: 296
- Updated best found solution, cost: 291
- Evaluating new solution: cost 280, meets the tolerances.
- Evaluating new solution: cost 287, meets the tolerances.
- Evaluating new solution: cost 288, does not meet the tolerances.
- Evaluating new solution: cost 287, does not meet the tolerances.
- Evaluating new solution: cost 283, meets the tolerances.
- Evaluating new solution: cost 283, does not meet the tolerances.
- Evaluating new solution: cost 262, does not meet the tolerances.
- Evaluating new solution: cost 283, does not meet the tolerances.
- Evaluating new solution: cost 282, does not meet the tolerances.
- Evaluating new solution: cost 288, meets the tolerances.
- Evaluating new solution: cost 289, meets the tolerances.
- Evaluating new solution: cost 288, meets the tolerances.
- Evaluating new solution: cost 290, meets the tolerances.
- Evaluating new solution: cost 281, does not meet the tolerances.
- Evaluating new solution: cost 286, does not meet the tolerances.
- Evaluating new solution: cost 287, meets the tolerances.
- Evaluating new solution: cost 284, meets the tolerances.
- Evaluating new solution: cost 282, meets the tolerances.
- Evaluating new solution: cost 285, does not meet the tolerances.
- Evaluating new solution: cost 277, meets the tolerances.
- Updated best found solution, cost: 280
- Updated best found solution, cost: 277
- Evaluating new solution: cost 272, meets the tolerances.
- Evaluating new solution: cost 266, meets the tolerances.
- Evaluating new solution: cost 269, meets the tolerances.
- Evaluating new solution: cost 271, does not meet the tolerances.
- Evaluating new solution: cost 274, meets the tolerances.
- Evaluating new solution: cost 275, meets the tolerances.
- Evaluating new solution: cost 274, does not meet the tolerances.
- Evaluating new solution: cost 275, meets the tolerances.
- Evaluating new solution: cost 276, does not meet the tolerances.
- Evaluating new solution: cost 271, meets the tolerances.
- Evaluating new solution: cost 267, meets the tolerances.
- Evaluating new solution: cost 270, meets the tolerances.
- Evaluating new solution: cost 272, meets the tolerances.
- Evaluating new solution: cost 264, does not meet the tolerances.
- Evaluating new solution: cost 265, does not meet the tolerances.
- Evaluating new solution: cost 269, meets the tolerances.
- Evaluating new solution: cost 270, meets the tolerances.
- Evaluating new solution: cost 269, meets the tolerances.
- Evaluating new solution: cost 276, meets the tolerances.
- Evaluating new solution: cost 274, meets the tolerances.
- Updated best found solution, cost: 272
- Updated best found solution, cost: 266
- + Optimization has finished.
 - Neighborhood search complete.
 - Maximum number of iterations completed.
- + Fixed-point implementation that met the tolerances found.
 - Total cost: 266
 - Maximum absolute difference: 0.087035
 - Use the explore method of the result to explore the implementation.



Use the `explore` method of the `OptimizationResult` object, `result`, to launch Simulation Data Inspector and explore the design containing the smallest total number of bits while maintaining the numeric tolerances specified in the `opt` object.

```
explore(result);
```

You can revert your model back to its original state using the `revert` method of the `OptimizationResult` object.

```
revert(result);
```

See Also

Functions

`fxpopt`

Classes

`fxpOptimizationOptions`

More About

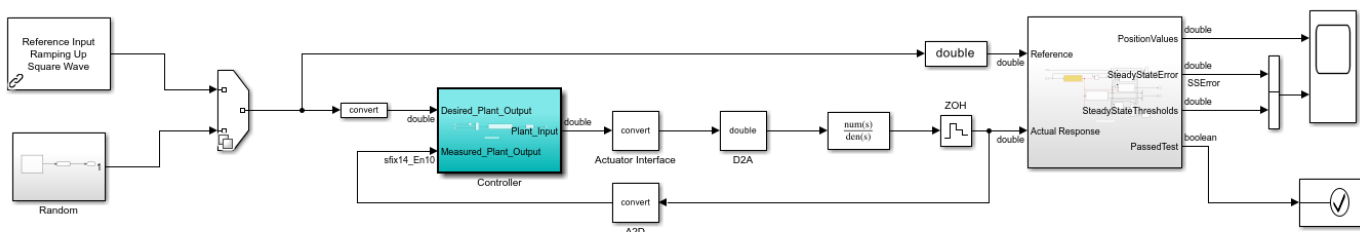
- “Data Type Optimization Not Successful” on page 49-31

Optimize Data Types Using Multiple Simulation Scenarios

This example shows how to create multiple simulation scenarios, and use the scenarios to optimize the fixed-point data types of a system.

Open the model. In this example, you optimize the data types of the Controller subsystem. The model is set up to use either a ramp input, or a random input. The model uses an Assertion block rather than using signal tolerances to verify the numerical behavior of the fixed-point implementation. For more information, see “Specify Behavioral Constraints” on page 42-18.

```
model = 'ex_controllerHarness';
open_system(model);
```



Copyright 2018 The MathWorks, Inc.

Create the Simulation Scenarios

Create a `Simulink.SimulationInput` object that contains the different scenarios. Use both the ramp input as well as four different seeds for the random input.

```
si = Simulink.SimulationInput.empty(5, 0);

% scan through 4 different seeds for the random input
rng(1);
seeds = randi(1e6, [1 4]);

for sIndex = 1:length(seeds)
    si(sIndex) = Simulink.SimulationInput(model);
    si(sIndex) = si(sIndex).setVariable('SOURCE', 2); % SOURCE == 2 corresponds to the random input
    si(sIndex) = si(sIndex).setBlockParameter([model '/Random/uniformRandom'], 'Seed', num2str(seeds(sIndex)));
    si(sIndex) = si(sIndex).setUserString(sprintf('random_%i', seeds(sIndex)));
end

% setting SOURCE == 1 corresponds to the ramp input
si(5) = Simulink.SimulationInput(model);
si(5) = si(5).setVariable('SOURCE', 1);
si(5) = si(5).setUserString('Ramp');
```

Specify Fixed-Point Optimization Options

To specify options for the optimization, such as the number of iterations and method for range collection, use the `fxpOptimizationOptions` object. This example uses derived range analysis to collect ranges for the system.

```
options = fxpOptimizationOptions('MaxIterations', 3e2, 'Patience', 50);
options.AdvancedOptions.PerformNeighborhoodSearch = false;
```

```

% use derived range analysis for range collection
options.AdvancedOptions.UseDerivedRangeAnalysis = true

options =

fxpOptimizationOptions with properties:

    MaxIterations: 300
    MaxTime: 600
    Patience: 50
    Verbosity: High
    AllowableWordLengths: [2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 ... ]
    ObjectiveFunction: BitWidthSum
    UseParallel: 0

Advanced Options
    AdvancedOptions: [1x1 DataTypeOptimization.AdvancedFxpOptimizationOptions]

```

Specify the simulation input objects as simulation scenarios in the advanced options.

```
options.AdvancedOptions.SimulationScenarios = si;
```

Run Optimization and Explore the Results

During the optimization, the software derives ranges for all simulation scenarios specified in the advanced options. The software verifies solutions against each simulation input scenario.

```
result = fxpopt(model, [model '/Controller'], options)
```

```

+ Starting data type optimization...
+ Checking for unsupported constructs.
+ Preprocessing
+ Modeling the optimization problem
  - Constructing decision variables
+ Running the optimization solver
  - Evaluating new solution: cost 496, does not meet the behavioral constraints.
  - Evaluating new solution: cost 976, does not meet the behavioral constraints.
  - Evaluating new solution: cost 1936, meets the behavioral constraints.
  - Updated best found solution, cost: 1936
+ Optimization has finished.
+ Fixed-point implementation that satisfies the behavioral constraints found. The best found
  - Total cost: 1936
  - Use the explore method of the result to explore the implementation.

```

```
result =
```

```
OptimizationResult with properties:
```

```

    Model: 'ex_controllerHarness'
    SystemUnderDesign: 'ex_controllerHarness/Controller'
    FinalOutcome: 'Fixed-point implementation that satisfies the behavioral constraints found'
    OptimizationOptions: [1x1 fxpOptimizationOptions]
    Solutions: [1x1 DataTypeOptimization.OptimizationSolution]

```

You can explore each solution as it compares to each simulation scenario you defined. Explore the best found solution and view it with the ramp simulation input. The ramp input is simulation scenario five.

```
solutionIndex = 1; % get the best found solution
scenarioIndex = 5; % get the 5th scenario (ramp)
solution = explore(result, solutionIndex, scenarioIndex);
```

See Also

Functions

fxpopt

Classes

fxpOptimizationOptions

More About

- “Data Type Optimization Not Successful” on page 49-31

Optimize Data Types for an FPGA with DSP Slices

This example shows how to use the `addSpecification` method of the `fxpOptimizationOptions` class to achieve better mapping for product blocks on DSP slices for FPGA targets. Use `addSpecification` to specify known data types in a system. After specifying these known parameters, when you optimize data types in the system, the optimization process does not change the specified block parameter data type.

Many FPGA boards have specific multiply-accumulate hardware accelerators, called DSP slices that speed up the execution of signal processing functions. DSP slices vary in size depending on the vendor. To gain the hardware acceleration benefits of DSP slices, it is common in FPGA design to map multiply and accumulate operations in the algorithm onto these slices.

In this example, optimize data types for 3 DSP families of Xilinx® boards, as well as for a generic 18x18 bit input. Use the `addSpecification` method to achieve a good mapping between the Product blocks in the design and the target DSP slices.

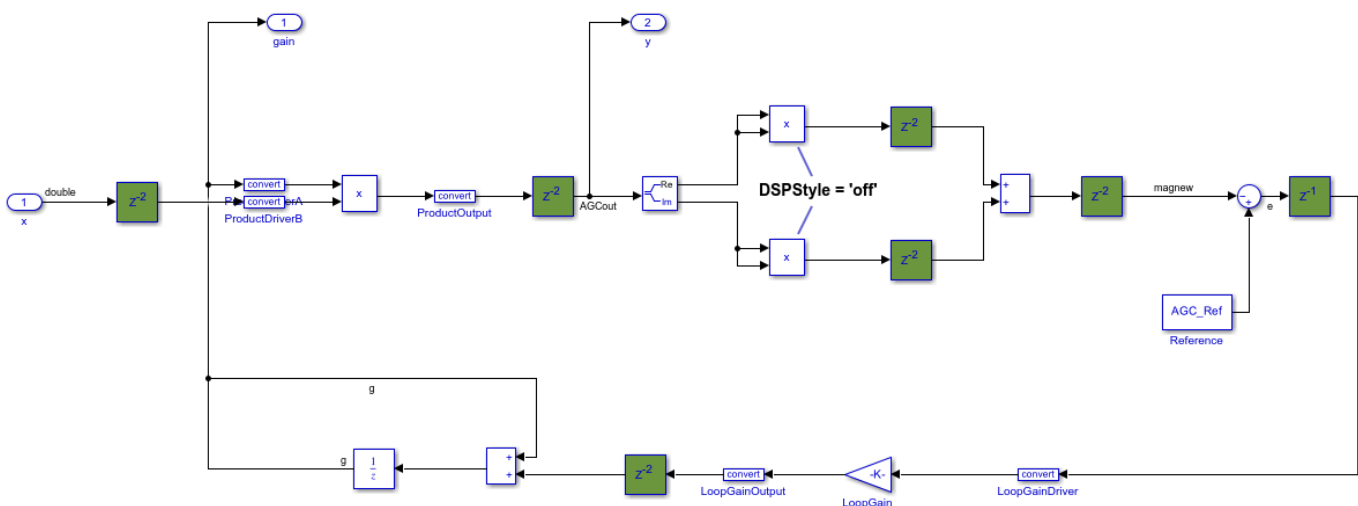
This example makes the following assumptions:

- 1 Only Product and Gain blocks are targeted for mapping to DSP slices. You can handle other blocks in a similar fashion.
- 2 Product blocks have only 2 inputs.
- 3 Driver blocks (blocks that precede Product or Gain blocks) have `OutDataTypeStr` as a parameter.
- 4 Only blocks that do not have the HDL property **DSPStyle** set to `off` are targeted.

Instrument the Model and Collect Ranges

To begin, open the system for which you want to optimize data types. In this example, data types are optimized for an automatic gain control algorithm.

```
model = 'mQAMAGC';
sud = [model '/Automatic Gain Control'];
open_system(model);
```



Initialize the QAM Tx subsystem.

```
initQAM;
```

Create a structure array that describes the properties of the target DSP slice. The `dspStyle` function included in this example provides information about common Xilinx® DSP slices including DSP48A1 (18x18 bit signed), DSP48E1 (18x25 bit signed), and DSP48E2 (18x27 bit signed). It also provides an example of a generic 18x18 signed DSP slice.

```
dspStyle = getDSPStyle('DSP48E2');
```

In this example, only Product and Gain blocks are targeted for mapping to DSP slices. Find all the Product and Gain blocks in the system under design that do not have the HDL block property **DSPStyle** set to off. For more information, see “DSPStyle” (HDL Coder).

```
productBlocks = find_system(sud, 'LookUnderMasks', 'on', 'BlockType', 'Product');
hasDSP0ff = cellfun(@(x) (isequal(hdlget_param(x, 'DSPStyle'), 'off')), productBlocks);
productDSP = productBlocks(~hasDSP0ff);
```

```
gainBlocks = find_system(sud, 'LookUnderMasks', 'on', 'BlockType', 'Gain');
hasDSP0ff = cellfun(@(x) (isequal(hdlget_param(x, 'DSPStyle'), 'off')), productBlocks);
gainDSP = gainBlocks(~hasDSP0ff);
```

Enable instrumentation to log minimum, maximum, and overflow data during simulation for Product blocks and driver blocks that precede Product blocks. Simulate the model to collect ranges.

```
c = DataTypeWorkflow.Converter(sud, 'TopModel', model);
c.CurrentRunName = 'RangeCollection';
c.simulateSystem('MinMaxOverflowLogging', 'MinMaxAndOverflow');
```

Get Specifications for Product Blocks

For efficient mapping of Product blocks to available DSP slices, consider the range requirements of the Product blocks. Create a structure array to store simulation minimum and maximum values for Product blocks collected during the range collection run.

```
specs = struct('Block', productDSP{1}, 'Drivers', [], 'Min', [], 'Max', []); %#ok<*SAGROW>
r = c.results(c.CurrentRunName, @(x) (strcmp(x.ResultName, productDSP{1})));
specs.Min = r.SimMin;
specs.Max = r.SimMax;
predecessorBlocks = predecessors(productDSP{1});
for pIndex = 1:numel(predecessorBlocks)
    pBlkObj = get_param(predecessorBlocks{pIndex}, 'Object');
    specs.Drivers(pIndex) = pBlkObj.Handle;
    r = c.results(c.CurrentRunName, @(x) (strcmp(x.ResultName, pBlkObj.getFullName())));
    specs.Min(pIndex+1) = r.SimMin;
    specs.Max(pIndex+1) = r.SimMax;
end
```

Store these known parameter specifications for the Product blocks in a `Simulink.Simulation.BlockParameter` object.

```
bpProductBlock = Simulink.Simulation.BlockParameter.empty(0,3);
```

```
fout = fi(max([abs(specs.Min(1)) abs(specs.Max(1))]), dspStyle.sout, dspStyle.wout);
bpProductBlock(1) = Simulink.Simulation.BlockParameter(specs.Block, ...
    'OutDataTypeStr', ...
    sprintf('fixdt(%i,%i,%i)', dspStyle.sout, dspStyle.wout, fout.FractionLength));
```

Assign the largest type to the largest range of the driver blocks. This ensures that the largest data type available in the target hardware is applied to the block with the largest range requirement.

```
dMax1 = max([abs(specs.Min(2)) abs(specs.Max(2))]);
dMax2 = max([abs(specs.Min(3)) abs(specs.Max(3))]);
if dMax1 < dMax2
    win_1 = dspStyle.win_1;
    sin_1 = dspStyle.sin_1;
    win_2 = dspStyle.win_2;
    sin_2 = dspStyle.sin_2;
    if dspStyle.win_1 >= dspStyle.win_2
        win_1 = dspStyle.win_2;
        sin_1 = dspStyle.sin_2;
        win_2 = dspStyle.win_1;
        sin_2 = dspStyle.sin_1;
    end
else
    win_1 = dspStyle.win_2;
    sin_1 = dspStyle.sin_2;
    win_2 = dspStyle.win_1;
    sin_2 = dspStyle.sin_1;
    if dspStyle.win_1 >= dspStyle.win_2
        win_1 = dspStyle.win_1;
        sin_1 = dspStyle.sin_1;
        win_2 = dspStyle.win_2;
        sin_2 = dspStyle.sin_2;
    end
end
```

Get specifications for blocks preceding Product blocks. Note that this example assumes that Product blocks have two inputs.

```
fin1 = fi(dMax1, sin_1, win_1);
blkObj = get_param(specs.Drivers(1), 'Object');
bpProductBlock(2) = Simulink.Simulation.BlockParameter(blkObj.getFullName, ...
    'OutDataTypeStr', ...
    sprintf('fixdt(%i, %i, %i)', sin_1, win_1, fin1.FractionLength));

fin2 = fi(dMax2, sin_2, win_2);
blkObj = get_param(specs.Drivers(2), 'Object');
bpProductBlock(3) = Simulink.Simulation.BlockParameter(blkObj.getFullName, ...
    'OutDataTypeStr', ...
    sprintf('fixdt(%i, %i, %i)', sin_2, win_2, fin2.FractionLength));
```

Get Specifications for Gain Blocks

Store known parameter specifications for the Gain blocks in a `Simulink.Simulation.BlockParameter` object.

```
bpGainBlock = Simulink.Simulation.BlockParameter.empty(0,3);
specs = struct('Block',gainDSP{1},'Drivers',[],'Min',[],'Max',[]); %#ok<*SAGROW>
r = c.results(c.CurrentRunName,@(x)(strcmp(x.ResultName,gainDSP{1})));
specs.Min = r.SimMin;
specs.Max = r.SimMax;
predecessorBlocks = predecessors(gainDSP{1});
pBlkObj = get_param(predecessorBlocks{1}, 'Object');
specs.Drivers(1) = pBlkObj.Handle;
r = c.results(c.CurrentRunName,@(x)(strcmp(x.ResultName,pBlkObj.getFullName())));
```

```
specs.Min(2) = r.SimMin;
specs.Max(2) = r.SimMax;
```

Get specifications for the output of Gain blocks.

```
fout = fi(max(abs([specs.Min(1) specs.Max(1)])),dspStyle.sout,dspStyle.wout);
bpGainBlock(1) = Simulink.Simulation.BlockParameter(gainDSP{1}, ...
    'OutDataTypeStr', ...
    sprintf('fixdt(%i, %i, %i)',dspStyle.sout,dspStyle.wout,fout.FractionLength));
```

Get specifications for the blocks preceding Gains blocks and assign this to the first configuration of the Simulink.Simulation.BlockParameter object bpGainBlock.

```
blkObj = get_param(specs.Drivers(1),'Object');
fin = fi(max(abs([specs.Min(2) specs.Max(2)])),dspStyle.sin_1,dspStyle.win_1);
bpGainBlock(2) = Simulink.Simulation.BlockParameter(blkObj.getFullName, ...
    'OutDataTypeStr', ...
    sprintf('fixdt(%i,%i,%i)',dspStyle.sin_1,dspStyle.win_1,fin.FractionLength));
```

Get specifications for the **Gain** parameter of the system under design and assign this value to the second configuration of bpGainBlock.

```
paramValue = str2double(get_param(sud,'AGC_Gain'));
fParam = fi(paramValue,dspStyle.sin_2,dspStyle.win_2);
bpGainBlock(3) = Simulink.Simulation.BlockParameter(gainDSP{1}, ...
    'ParamDataTypeStr', ...
    sprintf('fixdt(%i,%i,%i)',dspStyle.sin_2,dspStyle.win_2,fParam.FractionLength));
```

Define Constraints and Tolerances

Create an fxpOptimizationOptions object to define constraints and tolerances. Specify allowed word lengths of 8 bits to 32 bits.

```
options = fxpOptimizationOptions('AllowableWordLengths',8:2:32);
```

Use the addTolerance method to define tolerances for the differences between the original behavior of the system and the behavior using the optimized fixed-point data types.

```
addTolerance(options,sud,1,'RelTol',1e-2);
addTolerance(options,sud,2,'RelTol',1e-2);
addTolerance(options,sud,1,'AbsTol',1e-3);
addTolerance(options,sud,2,'AbsTol',1e-3);
```

Use the addSpecification method to define specifications for the Product and Gain blocks.

```
addSpecification(options,'BlockParameter',bpProductBlock); % set the specifications for the product
addSpecification(options,'BlockParameter',bpGainBlock); % set the specifications for the gain block
showSpecifications(options);
```

| Index | Name | BlockPath | Value |
|-------|------------------|---|------------------|
| 1 | OutDataTypeStr | mQAMAGC/Automatic Gain Control/LoopGain | 'fixdt(1, 45, 5) |
| 2 | ParamDataTypeStr | mQAMAGC/Automatic Gain Control/LoopGain | 'fixdt(1,27,35) |
| 3 | OutDataTypeStr | mQAMAGC/Automatic Gain Control/LoopGainDriver | 'fixdt(1,18,16) |
| 4 | OutDataTypeStr | mQAMAGC/Automatic Gain Control/Product | 'fixdt(1,45,43) |
| 5 | OutDataTypeStr | mQAMAGC/Automatic Gain Control/ProductDriverA | 'fixdt(1, 27, 2) |
| 6 | OutDataTypeStr | mQAMAGC/Automatic Gain Control/ProductDriverB | 'fixdt(1, 18, 2) |

Optimize Fixed-Point Data Types

Use the `fxpopt` function to run the optimization. The software analyzes ranges of objects in the system under design and the constraints specified in the `fxpOptimizationOptions` object to apply heterogeneous data types to your system while minimizing the total bit width. Known parameter specifications included using the `addSpecification` method is not affected by the optimization process.

```
result = fxpopt(model,sud,options);

+ Starting data type optimization...
+ Checking for unsupported constructs.
+ Preprocessing
+ Modeling the optimization problem
  - Constructing decision variables
+ Running the optimization solver
Exporting logged dataset prior to deleting run...done.
- Evaluating new solution: cost 200, does not meet the behavioral constraints.
- Evaluating new solution: cost 250, does not meet the behavioral constraints.
- Evaluating new solution: cost 300, does not meet the behavioral constraints.
- Evaluating new solution: cost 350, does not meet the behavioral constraints.
- Evaluating new solution: cost 400, does not meet the behavioral constraints.
- Evaluating new solution: cost 450, does not meet the behavioral constraints.
- Evaluating new solution: cost 500, meets the behavioral constraints.
- Updated best found solution, cost: 500
- Evaluating new solution: cost 494, meets the behavioral constraints.
- Updated best found solution, cost: 494
- Evaluating new solution: cost 492, meets the behavioral constraints.
- Updated best found solution, cost: 492
- Evaluating new solution: cost 490, does not meet the behavioral constraints.
- Evaluating new solution: cost 486, does not meet the behavioral constraints.
- Evaluating new solution: cost 490, meets the behavioral constraints.
- Updated best found solution, cost: 490
- Evaluating new solution: cost 488, meets the behavioral constraints.
- Updated best found solution, cost: 488
- Evaluating new solution: cost 478, meets the behavioral constraints.
- Updated best found solution, cost: 478
- Evaluating new solution: cost 474, meets the behavioral constraints.
- Updated best found solution, cost: 474
- Evaluating new solution: cost 470, meets the behavioral constraints.
- Updated best found solution, cost: 470
- Evaluating new solution: cost 466, meets the behavioral constraints.
- Updated best found solution, cost: 466
- Evaluating new solution: cost 462, meets the behavioral constraints.
- Updated best found solution, cost: 462
- Evaluating new solution: cost 458, meets the behavioral constraints.
- Updated best found solution, cost: 458
- Evaluating new solution: cost 452, meets the behavioral constraints.
- Updated best found solution, cost: 452
- Evaluating new solution: cost 450, meets the behavioral constraints.
- Updated best found solution, cost: 450
- Evaluating new solution: cost 448, does not meet the behavioral constraints.
- Evaluating new solution: cost 444, does not meet the behavioral constraints.
- Evaluating new solution: cost 448, meets the behavioral constraints.
- Updated best found solution, cost: 448
- Evaluating new solution: cost 446, meets the behavioral constraints.
- Updated best found solution, cost: 446
- Evaluating new solution: cost 436, meets the behavioral constraints.
```

- Updated best found solution, cost: 436
- Evaluating new solution: cost 432, meets the behavioral constraints.
- Updated best found solution, cost: 432
- Evaluating new solution: cost 428, meets the behavioral constraints.
- Updated best found solution, cost: 428
- Evaluating new solution: cost 424, meets the behavioral constraints.
- Updated best found solution, cost: 424
- Evaluating new solution: cost 420, meets the behavioral constraints.
- Updated best found solution, cost: 420
- Evaluating new solution: cost 416, meets the behavioral constraints.
- Updated best found solution, cost: 416
- Evaluating new solution: cost 410, meets the behavioral constraints.
- Updated best found solution, cost: 410
- Evaluating new solution: cost 408, meets the behavioral constraints.
- Updated best found solution, cost: 408
- Evaluating new solution: cost 406, does not meet the behavioral constraints.
- Evaluating new solution: cost 402, does not meet the behavioral constraints.
- Evaluating new solution: cost 406, meets the behavioral constraints.
- Updated best found solution, cost: 406
- Evaluating new solution: cost 404, meets the behavioral constraints.
- Updated best found solution, cost: 404
- Evaluating new solution: cost 394, meets the behavioral constraints.
- Updated best found solution, cost: 394
- Evaluating new solution: cost 390, meets the behavioral constraints.
- Updated best found solution, cost: 390
- Evaluating new solution: cost 386, meets the behavioral constraints.
- Updated best found solution, cost: 386
- Evaluating new solution: cost 382, meets the behavioral constraints.
- Updated best found solution, cost: 382
- Evaluating new solution: cost 378, meets the behavioral constraints.
- Updated best found solution, cost: 378
- Evaluating new solution: cost 374, meets the behavioral constraints.
- Updated best found solution, cost: 374
- Evaluating new solution: cost 368, does not meet the behavioral constraints.
- Evaluating new solution: cost 372, meets the behavioral constraints.
- Updated best found solution, cost: 372
- Evaluating new solution: cost 370, does not meet the behavioral constraints.
- Evaluating new solution: cost 366, does not meet the behavioral constraints.
- Evaluating new solution: cost 370, meets the behavioral constraints.
- Updated best found solution, cost: 370
- Evaluating new solution: cost 368, meets the behavioral constraints.
- Updated best found solution, cost: 368
- Evaluating new solution: cost 358, does not meet the behavioral constraints.
- Evaluating new solution: cost 364, meets the behavioral constraints.
- Updated best found solution, cost: 364
- Evaluating new solution: cost 360, meets the behavioral constraints.
- Updated best found solution, cost: 360
- Evaluating new solution: cost 356, meets the behavioral constraints.
- Updated best found solution, cost: 356
- Evaluating new solution: cost 352, meets the behavioral constraints.
- Updated best found solution, cost: 352
- Evaluating new solution: cost 348, meets the behavioral constraints.
- Updated best found solution, cost: 348
- Evaluating new solution: cost 342, does not meet the behavioral constraints.
- + Optimization has finished.
 - Neighborhood search complete.
 - Maximum number of iterations completed.
- + Fixed-point implementation that satisfies the behavioral constraints found. The best found

- Total cost: 348
- Maximum absolute difference: 0.006126
- Use the explore method of the result to explore the implementation.

Use the `explore` method of the `OptimizationResult` object, `result`, to launch the Simulation Data Inspector and explore the design.

```
explore(result)
```

```
ans =
```

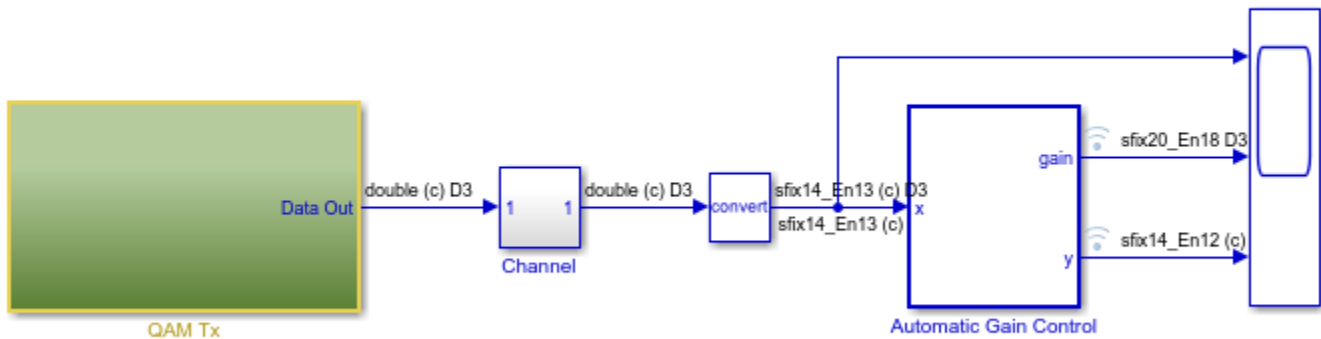
```
OptimizationSolution with properties:
```

```

    Cost: 348
    Pass: 1
    MaxDifference: 0.0061
    RunID: 23545
    RunName: {'solution_61d4a5478350738e21dfd31a47760c2cbf2e4794_1'}

```

Automatic Gain Control Sample



Copyright 2020 The MathWorks, Inc.

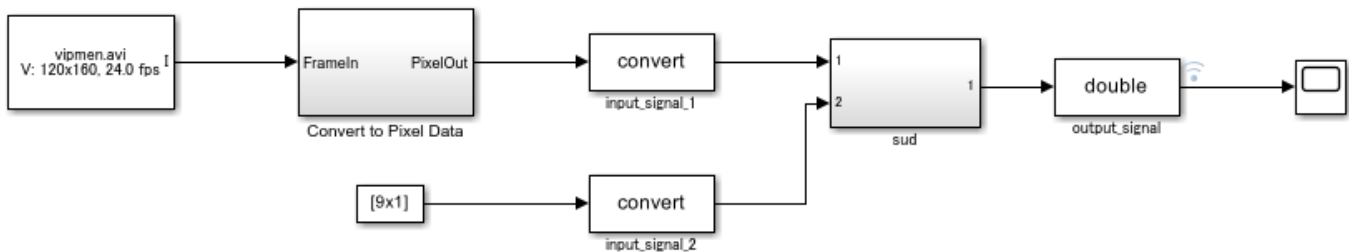
Use Data Type Optimization to Minimize Operator Counts

This example shows how to use the `ObjectiveFunction` property of the `fxpOptimizationOptions` class to minimize an estimated count of operators in the generated code. This results in a lower program memory size for C code generated from Simulink models.

Open the Model

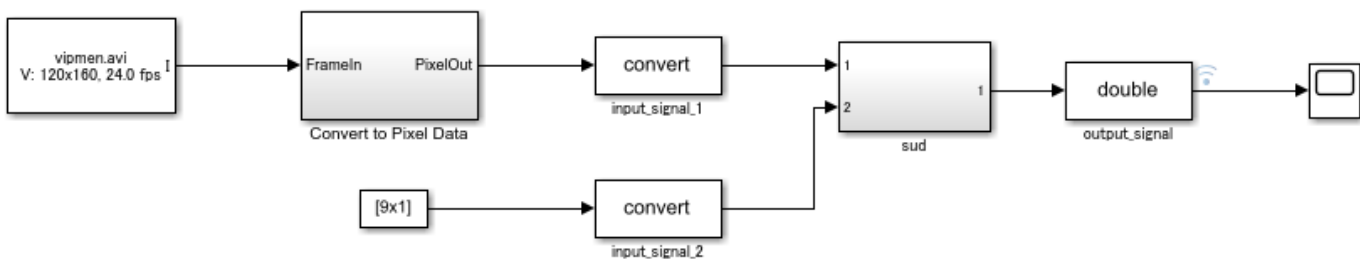
To begin, open the system for which you want to optimize the data types.

```
model = 'mFilter2D';
sud = [model '/sud'];
openExample('fixedpoint/FxpoptWithOperatorCountsExample')
```



Copyright 2017 The MathWorks, Inc.

Define Tolerances and Settings



Copyright 2017 The MathWorks, Inc.

Create an `fxpOptimizationOptions` object to define tolerances and optimization settings.

```
opt = fxpOptimizationOptions();
```

Use the `addTolerance` method to define tolerances for the differences between the original behavior of the system and the behavior using the optimized fixed-point data types.


```
tol = 1e-4;
opt.addTolerance([model '/output_signal'],1,'AbsTol',tol);
```

Set the advanced option `PerformNeighborhoodSearch` to `true` to perform a neighborhood search for the optimized solution.

```
opt.AdvancedOptions.PerformNeighborhoodSearch = true;
```

Set the `ObjectiveFunction` property to `OperatorCount` to instruct the optimization to minimize an estimated count of operators in the generated code.

```
opt.ObjectiveFunction = 'OperatorCount';
```

Use the `fxpopt` function to run the optimization. The software analyzes ranges of objects in the system under design and the constraints specified in the `fxpOptimizationOptions` object to apply heterogeneous data types to your system while minimizing the total operator count.

```
optRes = fxpopt(model,sud,opt);
```

```
+ Checking for unsupported constructs.
  + Preprocessing
  + Modeling the optimization problem
    - Constructing decision variables
  + Running the optimization solver
    - Evaluating new solution: cost 87, does not meet the tolerances.
    - Evaluating new solution: cost 87, does not meet the tolerances.
    - Evaluating new solution: cost 87, does not meet the tolerances.
    - Evaluating new solution: cost 87, does not meet the tolerances.
    - Evaluating new solution: cost 87, does not meet the tolerances.
    - Evaluating new solution: cost 85, does not meet the tolerances.
    - Evaluating new solution: cost 81, does not meet the tolerances.
    - Evaluating new solution: cost 83, does not meet the tolerances.
    - Evaluating new solution: cost 83, does not meet the tolerances.
    - Evaluating new solution: cost 83, does not meet the tolerances.
    - Evaluating new solution: cost 83, does not meet the tolerances.
    - Evaluating new solution: cost 83, does not meet the tolerances.
    - Evaluating new solution: cost 83, does not meet the tolerances.
    - Evaluating new solution: cost 83, does not meet the tolerances.
    - Evaluating new solution: cost 81, does not meet the tolerances.
    - Evaluating new solution: cost 83, meets the tolerances.
    - Updated best found solution, cost: 83
    - Evaluating new solution: cost 90, meets the tolerances.
    - Evaluating new solution: cost 83, does not meet the tolerances.
    - Evaluating new solution: cost 85, does not meet the tolerances.
    - Evaluating new solution: cost 89, meets the tolerances.
    - Evaluating new solution: cost 85, does not meet the tolerances.
    - Evaluating new solution: cost 88, meets the tolerances.
    - Evaluating new solution: cost 92, meets the tolerances.
    - Evaluating new solution: cost 90, does not meet the tolerances.
    - Evaluating new solution: cost 83, does not meet the tolerances.
    - Evaluating new solution: cost 89, does not meet the tolerances.
  + Optimization has finished.
    - Neighborhood search complete.
    - Reached limit of number of iterations without updates to the current best solution.
  + Fixed-point implementation that met the tolerances found.
    - Total cost: 83
    - Maximum absolute difference: 0.000091
    - Use the explore method of the result to explore the implementation.
```

Use the `explore` method of the `OptimizationResult` object, `optRes`, to launch the Simulation Data Inspector and explore the design containing the minimum number of operator counts while maintaining the numeric tolerances specified in the `opt` object.

```
explore(optRes)
```

Use the `openSimulationManager` method of the `OptimizationResult` object, `optRes`, to inspect the simulations run during optimization in Simulation Manager.

```
openSimulationManager(optRes)
```

Image Denoising Using Fixed-Point Quantized Restricted Boltzmann Machine Algorithm

This example highlights two workflows that can help you arrive at a fully embedded-efficient fixed-point design. This example shows how to:

- Use multiple simulation scenarios in data type optimization.
- Use ranges derived from design ranges for data-type optimization.
- Use different benchmarks of numerical behavior for each scenario using blocks from the Model Verification library.
- Replace math operations that do not support fixed-point data types with efficient lookup tables.

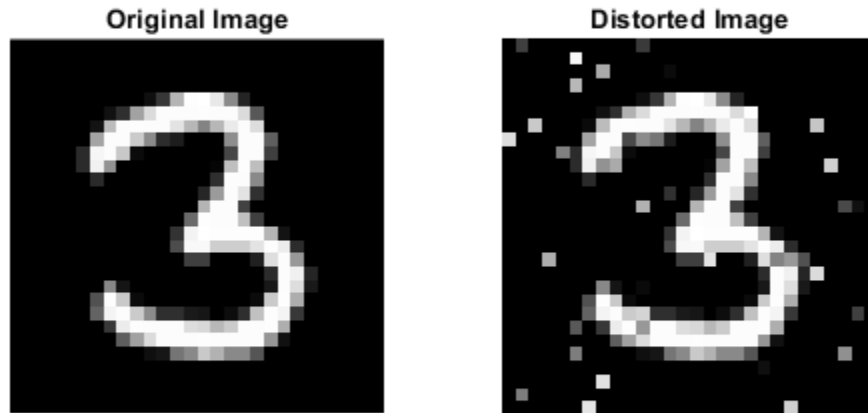
Convert Model to Use Optimal Fixed-Point Data Types

The model in this example uses a Restricted Boltzmann Machine (RBM) algorithm to denoise images. Load the image data and RBM algorithm weights. The original and distorted images are stored in the `imgOriginal` and `imgDistorted` variables. Each row of each matrix is a test image from the MNIST data set.

```
load RBMData;
```

Open and view the first set of test images.

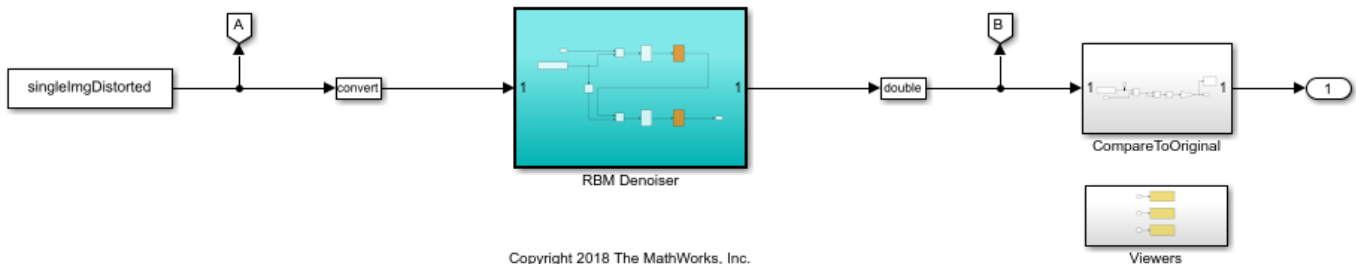
```
singleImgDistorted = imgDistorted(1,:);  
singleImgOriginal = imgOriginal(1,:);  
imgSize = length(singleImgOriginal);  
  
subplot(1,2,1)  
imshow(reshape(singleImgOriginal,[28,28]))'  
title('Original Image');  
subplot(1,2,2)  
imshow(reshape(singleImgDistorted,[28,28]))'  
title('Distorted Image')
```



Open the model. The model loads a distorted test image, uses the RBM algorithm to denoise the image, and then compares the denoised image to the original image without added noise. To improve simulation speed, the video display is turned off in this model. To turn on the video display, set the `DISPLAY_VIEWER` variable to 1.

```
model = 'ex_rbmDenoiser01';
open_system(model);

DISPLAY_VIEWER = 0;
```



Copyright 2018 The MathWorks, Inc.

RBM Denoiser as made by Andrej Karpathy
<https://code.google.com/archive/p/matrbm/>

When converting a model to use fixed-point data types, it is important to collect ranges while exercising the model over its full operating range. You can do this by defining multiple simulation

scenarios. In this example, each of the five simulation scenarios defines a new set of test images to denoise and compare to the original image.

```

IMGN = 5;
si = Simulink.SimulationInput.empty(0, IMGN);

for indx = 1:IMGN
    si(indx) = Simulink.SimulationInput(model);
    si(indx) = si(indx).setVariable('singleImgDistorted', imgDistorted(indx,:));
    si(indx) = si(indx).setVariable('singleImgOriginal', imgOriginal(indx,:));
end

```

In each simulation scenario, verify that the mean-squared error between the original image and the denoised image is less than 0.02.

```

si(1) = si(1).setBlockParameter([model '/CompareToOriginal/check'], 'max', '0.02');
si(2) = si(2).setBlockParameter([model '/CompareToOriginal/check'], 'max', '0.02');
si(3) = si(3).setBlockParameter([model '/CompareToOriginal/check'], 'max', '0.02');
si(4) = si(4).setBlockParameter([model '/CompareToOriginal/check'], 'max', '0.02');
si(5) = si(5).setBlockParameter([model '/CompareToOriginal/check'], 'max', '0.03');

```

Define the options to use during optimization. For this example, restrict the word lengths in the converted model to be between 8 and 16 bits. You can also restrict the number of iterations the optimization algorithm performs.

```

options = fxpOptimizationOptions(...
    'AllowableWordLengths', [8 16], ...
    'MaxIterations', 50, ...
    'Patience', 50);

```

To collect derived ranges in the model in addition to using the simulation scenarios to collect simulation ranges, set the `UseDerivedRangeAnalysis` option to `true`. Derived range analysis often returns a more conservative estimate of the dynamic ranges in the system than ranges collected through simulations.

```
options.AdvancedOptions.UseDerivedRangeAnalysis = true;
```

Specify the simulation scenarios to use during the optimization.

```
options.AdvancedOptions.SimulationScenarios = si;
```

Use the `fxpopt` function to optimize the data types in the RBM Denoiser subsystem according to the options specified in the `fxpOptimizationOptions` object, `options`

```

result = fxpopt(model, [model '/RBM Denoiser'], options);

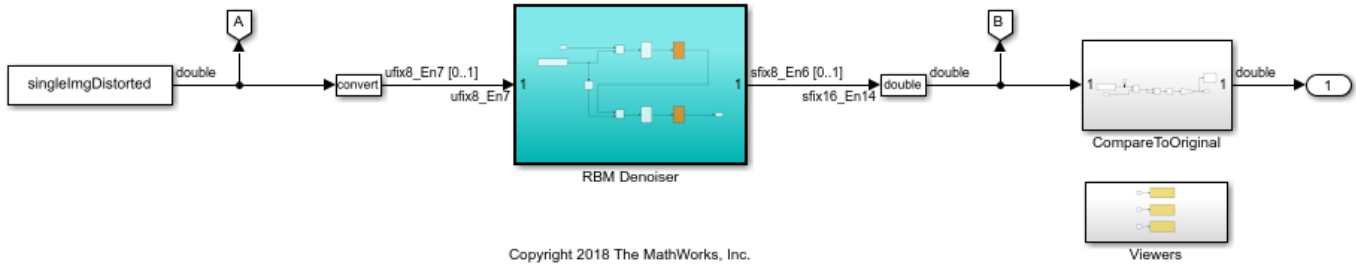
+ Starting data type optimization...
+ Checking for unsupported constructs.
+ Preprocessing
+ Modeling the optimization problem
  - Constructing decision variables
+ Running the optimization solver
  - Evaluating new solution: cost 344, does not meet the behavioral constraints.
  - Evaluating new solution: cost 656, meets the behavioral constraints.
  - Updated best found solution, cost: 656
  - Evaluating new solution: cost 640, meets the behavioral constraints.
  - Updated best found solution, cost: 640
  - Evaluating new solution: cost 624, does not meet the behavioral constraints.

```

```
- Evaluating new solution: cost 632, meets the behavioral constraints.
- Updated best found solution, cost: 632
- Evaluating new solution: cost 608, meets the behavioral constraints.
- Updated best found solution, cost: 608
- Evaluating new solution: cost 600, meets the behavioral constraints.
- Updated best found solution, cost: 600
- Evaluating new solution: cost 584, does not meet the behavioral constraints.
- Evaluating new solution: cost 592, meets the behavioral constraints.
- Updated best found solution, cost: 592
- Evaluating new solution: cost 568, meets the behavioral constraints.
- Updated best found solution, cost: 568
- Evaluating new solution: cost 560, meets the behavioral constraints.
- Updated best found solution, cost: 560
- Evaluating new solution: cost 544, meets the behavioral constraints.
- Updated best found solution, cost: 544
- Evaluating new solution: cost 504, meets the behavioral constraints.
- Updated best found solution, cost: 504
- Evaluating new solution: cost 440, meets the behavioral constraints.
- Updated best found solution, cost: 440
- Evaluating new solution: cost 432, meets the behavioral constraints.
- Updated best found solution, cost: 432
- Evaluating new solution: cost 424, meets the behavioral constraints.
- Updated best found solution, cost: 424
- Evaluating new solution: cost 408, meets the behavioral constraints.
- Updated best found solution, cost: 408
- Evaluating new solution: cost 400, meets the behavioral constraints.
- Updated best found solution, cost: 400
- Evaluating new solution: cost 392, meets the behavioral constraints.
- Updated best found solution, cost: 392
- Evaluating new solution: cost 376, meets the behavioral constraints.
- Updated best found solution, cost: 376
- Evaluating new solution: cost 360, does not meet the behavioral constraints.
- Evaluating new solution: cost 360, does not meet the behavioral constraints.
- Evaluating new solution: cost 392, meets the behavioral constraints.
- Evaluating new solution: cost 416, meets the behavioral constraints.
- Evaluating new solution: cost 424, meets the behavioral constraints.
- Evaluating new solution: cost 480, meets the behavioral constraints.
- Evaluating new solution: cost 472, meets the behavioral constraints.
- Evaluating new solution: cost 512, meets the behavioral constraints.
+ Optimization has finished.
  - Neighborhood search complete.
  - Maximum number of iterations completed.
+ Fixed-point implementation that satisfies the behavioral constraints found. The best found
  - Total cost: 376
  - Use the explore method of the result to explore the implementation.
```

To compare the behavior of the double-precision baseline model against the model that uses fixed-point data types, save the optimized, fixed-point model with a new name.

```
modelAfterFxpopt = 'ex_rbmDenoiser02';
save_system(model, modelAfterFxpopt);
```



Copyright 2018 The MathWorks, Inc.

RBM Denoiser as made by Andrej Karpathy
<https://code.google.com/archive/p/matrbm/>

Replace Logistic Regression with a Lookup Table

The LogisticExpression subsystems contain operations that do not support fixed-point data types. Replace these subsystems with lookup tables that closely approximate the original behavior.

```
functionToApproximate = [modelAfterFxpopt '/RBM Denoiser/Logistic/LogisticExpression'];
problem = FunctionApproximation.Problem(functionToApproximate);
problem.Options.AbsTol = 2^-6;
problem.Options.RelTol = 2^-7;
solution = solve(problem);
replaceWithApproximate(solution);
```

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLs | TableData WL | BreakpointSpec: |
|----|---------------|----------|------------|-----------------|--------------|-----------------|
| 0 | 32 | 0 | 2 | 8 | 8 | EvenPow |
| 1 | 160 | 0 | 18 | 8 | 8 | EvenPow |
| 2 | 312 | 0 | 37 | 8 | 8 | EvenPow |
| 3 | 704 | 0 | 86 | 8 | 8 | EvenPow |
| 4 | 2064 | 1 | 256 | 8 | 8 | EvenPow |
| 5 | 128 | 0 | 14 | 8 | 8 | EvenPow |
| 6 | 120 | 0 | 13 | 8 | 8 | EvenPow |
| 7 | 248 | 0 | 29 | 8 | 8 | EvenPow |
| 8 | 224 | 0 | 26 | 8 | 8 | EvenPow |
| 9 | 528 | 0 | 64 | 8 | 8 | EvenPow |
| 10 | 432 | 0 | 52 | 8 | 8 | EvenPow |
| 11 | 1040 | 1 | 128 | 8 | 8 | EvenPow |
| 12 | 96 | 0 | 10 | 8 | 8 | EvenPow |
| 13 | 88 | 0 | 9 | 8 | 8 | EvenPow |
| 14 | 168 | 0 | 19 | 8 | 8 | EvenPow |
| 15 | 128 | 1 | 8 | 8 | 8 | Explicit |
| 16 | 128 | 1 | 8 | 8 | 8 | Explicit |
| 17 | 2064 | 1 | 256 | 8 | 8 | EvenPow |
| 18 | 1040 | 1 | 128 | 8 | 8 | EvenPow |

Best Solution

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLs | TableData WL | BreakpointSpec: |
|----|---------------|----------|------------|-----------------|--------------|-----------------|
| 15 | 128 | 1 | 8 | 8 | 8 | Explicit |

Because both of the LogisticExpression subsystems implement the same algorithm, you can replace the second LogisticExpression subsystem with the same lookup table created in the previous step.

```
lutBlockPath = functionToApproximate;
subsystemToReplace = [modelAfterFxpopt '/RBM Denoiser/Logistic1/LogisticExpression'];
pos = get_param(subsystemToReplace, 'Position');
```

```
delete_block(subsystemToReplace);
add_block(lutBlockPath, subsystemToReplace, 'Position', pos);
set_param(subsystemToReplace, 'Commented', 'off');
```

Compare Behavior of Original Model with the Embedded-Efficient Version

Compare the simulation behavior of the fixed-point model with the lookup table approximations against the original double-precision baseline version. Define the same simulation scenarios for the updated model.

```
siFA = Simulink.SimulationInput.empty(0, IMG_N);
for indx = 1:IMG_N
    siFA(indx) = Simulink.SimulationInput(modelAfterFxpopt);
    siFA(indx) = siFA(indx).setVariable('singleImgDistorted', imgDistorted(indx,:));
    siFA(indx) = siFA(indx).setVariable('singleImgOriginal', imgOriginal(indx,:));
end
siFA(1) = siFA(1).setBlockParameter([modelAfterFxpopt '/CompareToOriginal/check'], 'max', '0');
siFA(2) = siFA(2).setBlockParameter([modelAfterFxpopt '/CompareToOriginal/check'], 'max', '0');
siFA(3) = siFA(3).setBlockParameter([modelAfterFxpopt '/CompareToOriginal/check'], 'max', '0');
siFA(4) = siFA(4).setBlockParameter([modelAfterFxpopt '/CompareToOriginal/check'], 'max', '0');
siFA(5) = siFA(5).setBlockParameter([modelAfterFxpopt '/CompareToOriginal/check'], 'max', '0');
```

Simulate and observe the simulation behavior of the model that contains the lookup table replacements. The model throws an error if the mean-square error between the original image, and the denoised image is greater than 0.02.

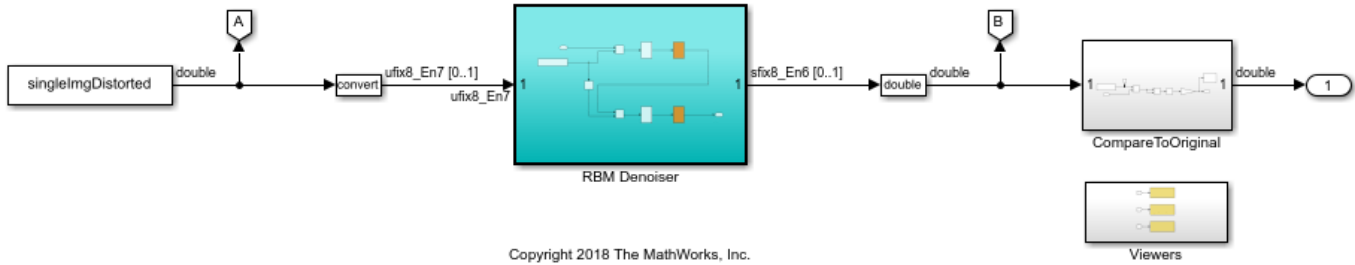
```
simOutAfterFA = sim(siFA);
assert(all(arrayfun(@(x) isempty(x.ErrorMessage), simOutAfterFA)), 'Final model does not meet tl

[01-Oct-2020 10:50:29] Running simulations...
[01-Oct-2020 10:50:30] Completed 1 of 5 simulation runs
[01-Oct-2020 10:50:31] Completed 2 of 5 simulation runs
[01-Oct-2020 10:50:32] Completed 3 of 5 simulation runs
[01-Oct-2020 10:50:32] Completed 4 of 5 simulation runs
[01-Oct-2020 10:50:33] Completed 5 of 5 simulation runs
```

Save the Model

Save the model after replacing the unsupported subsystems with lookup table approximations.

```
modelAfterFunctionApproximation = 'ex_rbmDenoiser03';
save_system(modelAfterFxpopt, modelAfterFunctionApproximation);
```

Copyright 2018 The MathWorks, Inc.

RBM Denoiser as made by Andrej Karpathy
<https://code.google.com/archive/p/matrbm/>

See Also

Classes

fxpOptimizationOptions

Functions

fxpopt

More About

- “Optimize Fixed-Point Data Types for a System” on page 40-14
- “Propose Data Types For Merged Simulation Ranges” on page 42-54
- “Approximate Functions with Lookup Tables”

Optimize the Fixed-Point Data Types of a System Using the Fixed-Point Tool

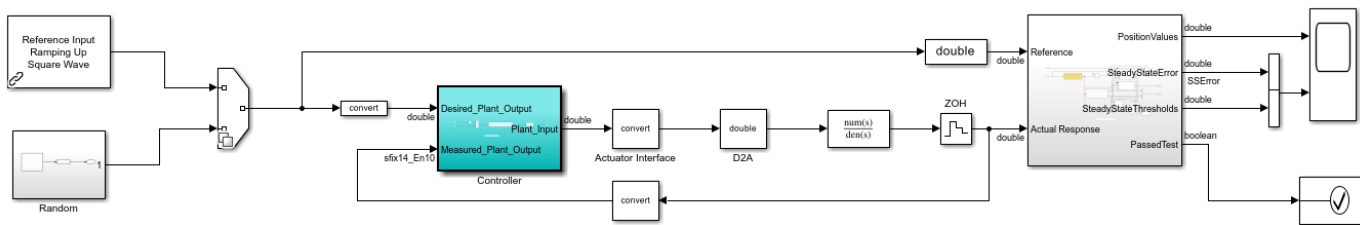
This example shows how to define simulation scenarios and use the Fixed-Point Tool to collect ranges by running simulations using these scenarios. You can then use the Fixed-Point Tool to optimize the fixed-point data types of the system.

During the optimization, the software establishes a baseline by simulating the original model. It then constructs different fixed-point versions of your model and runs simulations to determine the behavior using the new data types. The optimization selects the model that minimizes the objective function while also meeting the specified behavioral constraints. Including a `Simulink.SimulationInput` object in the setup allows you to define additional simulation scenarios to consider during the optimization. A comprehensive set of input signals can help to ensure that the full operating range of your design is exercised during the optimization process.

Open Model and Define Simulation Scenarios

Open the model. In this example, you optimize the data types of the Controller subsystem. The model is set up to use either a ramp input, or a random input.

```
model = 'ex_controllerHarness';
open_system(model)
```



Copyright 2018 The MathWorks, Inc.

Create a `Simulink.SimulationInput` object that contains the different scenarios. Use both the ramp input as well as four different seeds for the random input.

```
si = Simulink.SimulationInput.empty(5, 0);

% scan through 4 different seeds for the random input
rng(1);
seeds = randi(1e6, [1 4]);

for sIndex = 1:length(seeds)
    si(sIndex) = Simulink.SimulationInput(model);
    si(sIndex) = si(sIndex).setVariable('SOURCE', 2); % SOURCE == 2 corresponds to the random input
    si(sIndex) = si(sIndex).setBlockParameter([model '/Random/uniformRandom'], 'Seed', num2str(seeds(sIndex)));
    si(sIndex) = si(sIndex).setUserString(sprintf('random_%i', seeds(sIndex)));
end

% setting SOURCE == 1 corresponds to the ramp input
si(5) = Simulink.SimulationInput(model);
```

```
si(5) = si(5).setVariable('SOURCE', 1);
si(5) = si(5).setUserString('Ramp');
```

Prepare System for Conversion

To optimize the data types in the mode, use the Fixed-Point Tool.

- 1 In the **Apps** gallery of the `ex_controllerHarness` model, select **Fixed-Point Tool**.
- 2 In the Fixed-Point Tool, under **New** workflow, select **Optimized Fixed-Point Conversion**.
- 3 Under **System Under Design (SUD)**, select the subsystem for which you want to optimize the data types. In this example, select `Controller`.
- 4 Under **Range Collection Mode**, select **Simulation Ranges** as the range collection method.
- 5 Under **Simulation Inputs**, you can specify `Simulink.SimulationInput` objects to exercise your design over its full operating range. In this example, use the simulation scenarios you defined. Set **Simulation Inputs** to `si`.
- 6 You can specify tolerances for any signal in the model with signal logging enabled in the table under **Signal Tolerances**. In this example, the **Signal Tolerances** section indicates that the model contains no logged signals. Because this model uses an Assertion block from the Model Verification library to verify the numerical behavior of the system during optimization, specifying signal tolerances is optional. For more information, see “Specify Behavioral Constraints” on page 42-18.

If you have an `fxpOptimizationOptions` object saved from a prior command line optimization, for example if you previously optimized data types in this model following the example “Optimize Data Types Using Multiple Simulation Scenarios” on page 40-20, you can import the `fxpOptimizationOptions` object under the **Advanced Options** section. The **Setup** pane and toolbar **Settings** menu will be updated with the imported values. If an optimization options object is not imported, the default values are used unless otherwise changed manually in the Fixed-Point Tool. For this example, you will specify the settings manually in the app.

- 7 In the toolbar, click **Prepare**. The Fixed-Point Tool checks the system under design for compatibility with the conversion process and reports any issues found in the model. When possible, the Fixed-Point Tool automatically changes settings that are not compatible. For more information, see “Use the Fixed-Point Tool to Prepare a System for Conversion” on page 38-2.

Optimize Data Types in the Fixed-Point Tool

- 1 To specify settings to use during the optimization, in the toolbar, click **Settings**.

In this example, use the following settings.

- Set **Allowable Word Lengths** to `[2:32]`.

This setting defines the word lengths that can be used in your optimized system. Use this setting to target the neighborhood search of the optimization process. The final result of the optimization uses word lengths in the intersection of this setting and word lengths compatible with hardware constraints specified in the **Hardware Implementation** pane of your model.

- Set **Max Iterations** to `3e2`.

This setting specifies the maximum number of iterations to perform in the optimization. The optimization process iterates through different solutions until it finds an ideal solution, reaches the maximum number of iterations, or reaches another stopping criteria.

- Set **Patience** to `50`.

This setting defines the maximum number of iterations where no new best solution is found. The optimization continues as long as the algorithm continues to find new best solutions.

- Set **Objective Function** to **Bit Width Sum**. Using this setting instructs the optimization to minimize the total bit width of the final design while meeting the specified constraints.

For more information about optimization settings, see `fxpOptimizationOptions`.

- 2 To optimize the data types in the model according to the specified settings, click **Optimize Data Types**.

During the optimization process, the software analyzes ranges of objects in your system under design and the constraints specified in the settings to apply heterogeneous data types to your system while minimizing the objective function. Details about the optimization process are printed to the **Optimization Details** pane in the Fixed-Point Tool.

You can stop the optimization solver before the optimization search is complete by clicking **Stop** in the toolbar of the Fixed-Point Tool. To resume optimization from where you left off, click **Optimize Data Types** to restart the optimization solver. **Perform Neighborhood Search** must be activated to resume optimization. You can also choose to modify the stopping criteria for the optimization solver, including **Max Iterations**, **Max Time (sec)**, and **Patience (iterations)**. This can help to quickly narrow down a large optimization search space.

Examine Results

When the optimization completes, the Fixed-Point Tool displays a table that contains all of the solutions found during the optimization process. The first solution in the table corresponds to the solution with the lowest cost (smallest total bit width).

OPTIMIZED FIXED-POINT CONVERSION

Fixed-point implementation that satisfies the behavioral constraints found. The best found solution is applied on the model.

Best solution: Solution 1

Solution currently applied: Solution 1

Cost (Bit Width Sum): 1160

Max difference: 0

Scenarios passed/total: 5/5

Stopping criteria:

- Reached limit of number of iterations without updates to the current best solution.

[View optimization settings](#)

Hide Solutions Table

To apply an optimization solution to the system, select a solution from the table and click Apply and Compare.

| Name | Status | Cost (Bit Width Sum) | Max Difference | Scenarios Passed/Total |
|----------------------------------|--------|----------------------|----------------|------------------------|
| ▶ Solution 1 (currently applied) | ✓ | 1160 | 0 | 5/5 |
| ▶ Solution 2 | ✓ | 1168 | 0 | 5/5 |
| ▶ Solution 3 | ✓ | 1176 | 0 | 5/5 |
| ▶ Solution 4 | ✓ | 1176 | 0 | 5/5 |
| ▶ Solution 5 | ✓ | 1184 | 0 | 5/5 |
| ▶ Solution 6 | ✓ | 1240 | 0 | 5/5 |
| ▶ Solution 7 | ✓ | 1256 | 0 | 5/5 |
| ▶ Solution 8 | ✓ | 1264 | 0 | 5/5 |
| ▶ Solution 9 | ✓ | 1280 | 0 | 5/5 |
| ▶ Solution 10 | ✓ | 1312 | 0 | 5/5 |

To view the settings that were used for optimization, in the **Result** pane, click **View optimization settings**.

To inspect the ranges that were collected for objects in your model during the optimization process, in the **Workflow Browser** pane, select **BaselineRun**.

The Fixed-Point Tool displays a summary of the ranges of objects in your model and histograms of the bits used by each object. Each column in the **Visualization of Simulation Data** pane represents a histogram for one object in your model. Each bin in a histogram corresponds to a bit in the binary word.

Selecting a column highlights the corresponding model object in the **Results** spreadsheet of the Fixed-Point Tool and populates the **Result Details** pane with more detailed information about the selected result.

You can use the data type visualization to see a summary of the ranges of objects in your model and to spot sources of overflow, underflows, and inefficient data types. Using the **Explore** tab of the Fixed-Point Tool, you can sort and filter results in the tool based on additional criteria.

The screenshot displays the Fixed-Point Tool interface. The top toolbar includes options like New, Prepare, Settings, Stop, Apply and Compare, Compare, Export Script, and Restore Original Model. The main workspace is divided into several panes:

- Workflow Browser:** Shows a tree view with 'BaselineRun' selected.
- Results:** A table listing model objects with their compiled data types and simulation ranges.

| Name | CompiledDT | SimMin | SimMax |
|--|------------|---------------------|----------------------|
| PIController/AntiWindUp/Switch | double | -0.0519866943359375 | 0.032543182373046875 |
| PIController/AntiWindUp/WouldIncreaseNeg... | boolean | | |
| PIController/AntiWindUp/WouldIncreasePos... | boolean | | |
| PIController/AntiWindUp/Zero | | | |
| PIController/IntegralActionSum : Accumulator | double | -1.9370651245117188 | 2.3435935974121094 |
| PIController/IntegralActionSum : Output | double | -1.9370651245117188 | 2.3435935974121094 |
| PIController/lxTsGain | double | -0.0519866943359375 | 0.032543182373046875 |
| PIController/lxTsGain : Gain | | | |
| PIController/Out1 | | | |
| PIController/PGain | double | -1.356852722167969 | 0.8493770599365235 |
| PIController/PGain : Gain | | | |
| PIController/PISum : Accumulator | double | -1.9371608734130858 | 2.5681831359863283 |
| PIController/PISum : Output | double | 1.0274608724420858 | 2.5681831359863283 |
- Visualization of Simulation Data:** A histogram showing the distribution of simulation data for all results in the model. The x-axis represents simulation values, and the y-axis represents the number of histogram bins. A legend indicates Overflows (red), Representable (grey), In-Range (blue), and Underflows (yellow).
- Result Details:** Provides detailed information for the selected object 'ex_controllerHarness/Controller/PIController/AntiWindUp/Switch'.

| Property | Compiled Data Type |
|-----------|--------------------------|
| Data Type | double |
| Minimum | -1.7976931348623157e+308 |
| Maximum | 1.7976931348623157e+308 |
| Precision | 4.94065645841247e-324 |

| Property | Minimum | Maximum |
|-----------------|------------------|-----------------|
| Shared Simul... | -0.0519866943... | 0.0325431823... |
| Simulation | -0.0519866943... | 0.0325431823... |
- Visualization of Simulation Data using double:** A smaller histogram showing the distribution of simulation data using double precision. Below it is a table summarizing overflow and underflow statistics.

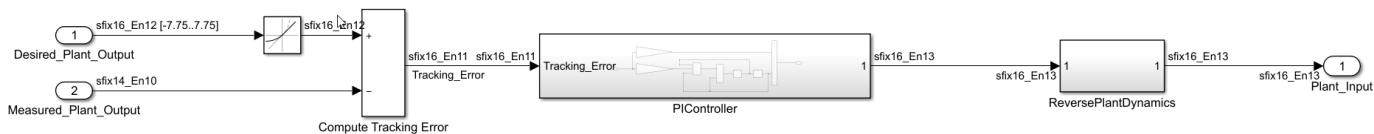
| Values | Potential Overflows | In-Range | Potential Underflows |
|----------|---------------------|----------|----------------------|
| Positive | 0 | 22706 | 0 |
| Negative | 0 | 22094 | 0 |
| Zero | 0 | 6400 | 0 |

Apply Optimized Data Types to the Model

- To apply the optimized data types to the model, in the solutions table, select the solution that you want to apply. In the **Explore** section of the toolbar, click **Apply and Compare**. The Fixed-Point Tool applies the selected solution that contains optimized fixed-point data types to the model and opens the Simulation Data Inspector.

In this example, select **Solution 1**, then click **Apply and Compare**.

2 In the Controller subsystem, you can see the applied, optimized fixed-point data types.



Tip The Fixed-Point Tool uses the Simulation Data Inspector tool plotting capabilities that enable you to plot logged signals for graphical analysis. Because this model does not contain any logged signals, the **Compare** button remains disabled in this example. To compare results of data type optimization using the Simulation Data Inspector, log one or more signals in your model.

Export Optimization Workflow Steps to a MATLAB Script

After optimizing data types in the Fixed-Point Tool, you can choose to export optimization workflow steps to a MATLAB script. This allows you to save the current optimization workflow steps and continue data type optimization using `fxpopt` at the command line.

In the toolstrip, click **Export Script**. The Fixed-Point Tool exports a script called `fxpOptimizationScript.m` to the current working directory:

```
model = 'ex_controllerHarness';
sud = 'ex_controllerHarness/Controller';
options = fxpOptimizationOptions();
options.MaxIterations = 300; % Maximum number of iterations to perform.
options.Patience = 50; % Maximum number of iterations where no new best solution is found.
options.AllowableWordLengths = 2:32; % Word lengths that can be used in your optimized system un
savedOptions = load('fxpOptimizationScript');
options.AdvancedOptions.SimulationScenarios = savedOptions.simulationScenarios;
result = fxpopt(model, sud, options);
explore(result);
```

The `Simulink.SimulationInput` object, `si`, used during optimization is exported to a MAT-file called `fxpOptimizationScript.mat`.

See Also

`fxpopt` | **Fixed-Point Tool**

More About

- “Image Denoising Using Fixed-Point Quantized Restricted Boltzmann Machine Algorithm” on page 40-33

- “Propose Data Types For Merged Simulation Ranges” on page 42-54

Perform Data Type Optimization with Custom Behavioral Constraints

This example shows how to optimize fixed-point data types with `fxpopt` using custom behavioral constraints in the frequency domain. The Simulink® models in this example demonstrate two ways of using blocks from the Model Verification library to author custom behavioral constraints on the SNR (signal-to-noise ratio) of a low-pass filter. In the first example, calculate an approximation of the noise floor; in the second example, calculate the SNR difference between the original model and the quantized model that uses fixed-point data types.

Behavioral Constraints for Data Type Optimization

Data type optimization seeks to minimize an objective function, such as the total bit width, while maintaining original system behavior within a specified tolerance. During the optimization, the software establishes a baseline by simulating the original model. The software then constructs different fixed-point versions of your model and runs simulations to determine the behavior using the new data types. The optimization selects the model that minimizes the objective function while also meeting the specified behavioral constraints.

To determine if the behavior of a new fixed-point implementation is acceptable, the optimization requires well-defined behavioral constraints. You must specify at least one behavioral constraint. There are two ways to specify behavioral constraints for use with `fxpopt`:

1. Tolerance-based simulation comparisons to the baseline reference - You can add tolerances for signals that have signal logging enabled using the `addTolerance` method of `fxpOptimizationOptions`. A tolerance specifies an envelope of passing behavior with respect to the baseline simulation of the model.
2. Simulation-based assertion checks - You can use blocks in the Model Verification block library to author custom verification expressions. Enabled model verification blocks in your model are interpreted as behavioral constraints by the optimization solver.

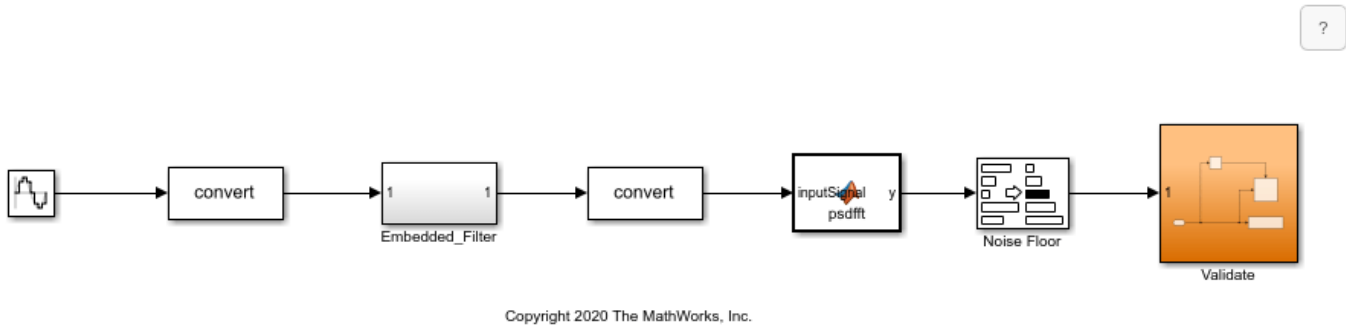
You can also use a combination of signal tolerances and model verification blocks to specify behavioral constraints for your model. For more information, see “Specify Behavioral Constraints” on page 42-18.

Example 1: Calculate an Approximation of the Noise Floor

In this example, you convert a low-pass filter to use optimized fixed-point data types using `fxpopt`. To ensure that the embedded version still meets the SNR requirements, the model `mLowPass_NoiseFloor` contains additional logic to compute the noise floor of the signal at the output of the filter. This signal is then routed to the `Validate` subsystem, which uses a `Check State Range` block from the Model Verification block library. If the noise floor values calculated during simulation are outside the specified static range, `fxpopt` will interpret the model as not passing the behavioral constraints.

To begin, open the system for which you want to optimize the data types.

```
model = 'mLowPass_NoiseFloor';  
sud = [model '/Embedded_Filter'];  
open_system(model);
```

Use the `fxpopt` function to run the optimization. Because a model verification block is used to specify the constraints, it is not required to specify signal tolerances using the `addTolerance` method of `fxpOptimizationOptions`.

```
result = fxpopt(model,sud)
```

```
+ Starting data type optimization...
+ Checking for unsupported constructs.
+ Preprocessing
+ Modeling the optimization problem
  - Constructing decision variables
+ Running the optimization solver
  - Evaluating new solution: cost 110, meets the behavioral constraints.
  - Updated best found solution, cost: 110
+ Optimization has finished.
  - Neighborhood search complete.
  - Reached limit of number of iterations without updates to the current best solution.
+ Fixed-point implementation that satisfies the behavioral constraints found. The best found
  - Total cost: 110
  - Use the explore method of the result to explore the implementation.
```

```
result =
```

```
OptimizationResult with properties:
```

```
    Model: 'mLowPass_NoiseFloor'
 SystemUnderDesign: 'mLowPass_NoiseFloor/Embedded_Filter'
  FinalOutcome: 'Fixed-point implementation that satisfies the behavioral constraints f
OptimizationOptions: [1x1 fxpOptimizationOptions]
    Solutions: [1x1 DataTypeOptimization.OptimizationSolution]
```

Use the `explore` method to explore the design containing the optimized fixed-point data types.

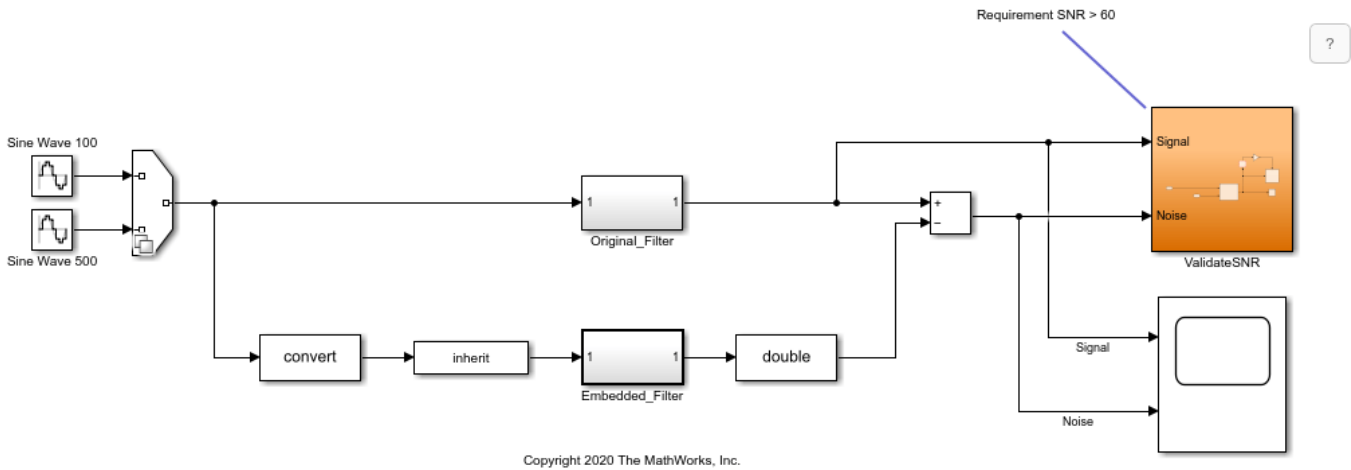
```
explore(result);
```

Example 2: Calculate SNR Difference Between Original and Quantized Models

In this example, optimize fixed-point data types for the low-pass filter model, `mLowPass_SNR`, using a different set of behavioral constraints that compare the SNR of the embedded filter with the SNR of the original filter, which uses double-precision data types. The `ValidateSNR` subsystem first computes the SNR then uses a Check Static Lower Bound block from the Model Verification library to assert that the SNR is greater than 60. If the SNR drops below this specified value, then `fxpopt` will interpret the model as not passing the behavioral constraints.

To begin, open the system for which you want to optimize the data types.

```
model = 'mLowPass_SNR';
sud = [model '/Embedded_Filter'];
Source = 1;
open_system(model);
```



For this example, specify additional simulation scenarios by creating a `Simulink.SimulationInput` object. Use one sine wave input with a frequency of 100 rad/sec and a second sine wave input with a frequency of 500 rad/sec. `fxpopt` will consider both simulation scenarios during the optimization. A comprehensive set of input signals can help to ensure that the full operating range of your design is exercised during the optimization process.

```
si(1) = Simulink.SimulationInput(model);
si(1) = si(1).setVariable('Source',1);
si(2) = Simulink.SimulationInput(model);
si(2) = si(2).setVariable('Source',2);
```

Create an `fxpOptimizationOptions` object. Use the advanced options to define simulation scenarios to consider during optimization.

```
options = fxpOptimizationOptions();
options.AdvancedOptions.SimulationScenarios = si;
```

Use the `fxpopt` function to run the optimization. Because a model verification block is used to specify the constraints, it is not required to specify signal tolerances using the `addTolerance` method of `fxpOptimizationOptions`.

```
result = fxpopt(model,sud,options);
```

```
+ Starting data type optimization...
+ Checking for unsupported constructs.
+ Preprocessing
+ Modeling the optimization problem
  - Constructing decision variables
+ Running the optimization solver
  - Evaluating new solution: cost 112, does not meet the behavioral constraints.
  - Evaluating new solution: cost 168, does not meet the behavioral constraints.
  - Evaluating new solution: cost 224, does not meet the behavioral constraints.
  - Evaluating new solution: cost 280, does not meet the behavioral constraints.
```

- Evaluating new solution: cost 336, does not meet the behavioral constraints.
- Evaluating new solution: cost 392, does not meet the behavioral constraints.
- Evaluating new solution: cost 448, does not meet the behavioral constraints.
- Evaluating new solution: cost 504, does not meet the behavioral constraints.
- Evaluating new solution: cost 560, does not meet the behavioral constraints.
- Evaluating new solution: cost 616, does not meet the behavioral constraints.
- Evaluating new solution: cost 672, does not meet the behavioral constraints.
- Evaluating new solution: cost 728, does not meet the behavioral constraints.
- Evaluating new solution: cost 784, does not meet the behavioral constraints.
- Evaluating new solution: cost 840, does not meet the behavioral constraints.
- Evaluating new solution: cost 896, does not meet the behavioral constraints.
- Evaluating new solution: cost 952, does not meet the behavioral constraints.
- Evaluating new solution: cost 1008, meets the behavioral constraints.
- Updated best found solution, cost: 1008
- Evaluating new solution: cost 995, meets the behavioral constraints.
- Updated best found solution, cost: 995
- Evaluating new solution: cost 994, meets the behavioral constraints.
- Updated best found solution, cost: 994
- Evaluating new solution: cost 993, meets the behavioral constraints.
- Updated best found solution, cost: 993
- Evaluating new solution: cost 992, meets the behavioral constraints.
- Updated best found solution, cost: 992
- Evaluating new solution: cost 991, meets the behavioral constraints.
- Updated best found solution, cost: 991
- Evaluating new solution: cost 990, meets the behavioral constraints.
- Updated best found solution, cost: 990
- Evaluating new solution: cost 989, meets the behavioral constraints.
- Updated best found solution, cost: 989
- Evaluating new solution: cost 988, meets the behavioral constraints.
- Updated best found solution, cost: 988
- Evaluating new solution: cost 987, meets the behavioral constraints.
- Updated best found solution, cost: 987
- Evaluating new solution: cost 986, meets the behavioral constraints.
- Updated best found solution, cost: 986
- Evaluating new solution: cost 985, meets the behavioral constraints.
- Updated best found solution, cost: 985
- Evaluating new solution: cost 984, meets the behavioral constraints.
- Updated best found solution, cost: 984
- Evaluating new solution: cost 982, meets the behavioral constraints.
- Updated best found solution, cost: 982
- Evaluating new solution: cost 981, does not meet the behavioral constraints.
- Evaluating new solution: cost 981, does not meet the behavioral constraints.
- Evaluating new solution: cost 981, does not meet the behavioral constraints.
- Evaluating new solution: cost 981, does not meet the behavioral constraints.
- Evaluating new solution: cost 981, does not meet the behavioral constraints.
- Evaluating new solution: cost 981, meets the behavioral constraints.
- Updated best found solution, cost: 981
- Evaluating new solution: cost 980, meets the behavioral constraints.
- Updated best found solution, cost: 980
- Evaluating new solution: cost 979, does not meet the behavioral constraints.
- Evaluating new solution: cost 979, meets the behavioral constraints.
- Updated best found solution, cost: 979
- Evaluating new solution: cost 978, does not meet the behavioral constraints.
- Evaluating new solution: cost 978, meets the behavioral constraints.
- Updated best found solution, cost: 978
- Evaluating new solution: cost 977, does not meet the behavioral constraints.
- Evaluating new solution: cost 977, meets the behavioral constraints.

- Updated best found solution, cost: 977
- Evaluating new solution: cost 976, meets the behavioral constraints.
- Updated best found solution, cost: 976
- Evaluating new solution: cost 975, meets the behavioral constraints.
- Updated best found solution, cost: 975
- Evaluating new solution: cost 974, does not meet the behavioral constraints.
- Evaluating new solution: cost 974, meets the behavioral constraints.
- Updated best found solution, cost: 974
- Evaluating new solution: cost 973, does not meet the behavioral constraints.
- Evaluating new solution: cost 973, meets the behavioral constraints.
- Updated best found solution, cost: 973
- Evaluating new solution: cost 972, does not meet the behavioral constraints.
- Evaluating new solution: cost 972, meets the behavioral constraints.
- Updated best found solution, cost: 972
- Evaluating new solution: cost 971, meets the behavioral constraints.
- Updated best found solution, cost: 971
- Evaluating new solution: cost 970, meets the behavioral constraints.
- Updated best found solution, cost: 970
- Evaluating new solution: cost 969, meets the behavioral constraints.
- Updated best found solution, cost: 969
- Evaluating new solution: cost 968, meets the behavioral constraints.
- Updated best found solution, cost: 968
- Evaluating new solution: cost 967, meets the behavioral constraints.
- Updated best found solution, cost: 967
- Evaluating new solution: cost 966, meets the behavioral constraints.
- Updated best found solution, cost: 966
- Evaluating new solution: cost 965, meets the behavioral constraints.
- Updated best found solution, cost: 965
- Evaluating new solution: cost 964, meets the behavioral constraints.
- Updated best found solution, cost: 964
- Evaluating new solution: cost 951, does not meet the behavioral constraints.
- Evaluating new solution: cost 963, does not meet the behavioral constraints.
- Evaluating new solution: cost 963, does not meet the behavioral constraints.
- Evaluating new solution: cost 963, does not meet the behavioral constraints.
- Evaluating new solution: cost 963, does not meet the behavioral constraints.
- Evaluating new solution: cost 963, does not meet the behavioral constraints.
- + Optimization has finished.
 - Neighborhood search complete.
 - Maximum number of iterations completed.
- + Fixed-point implementation that satisfies the behavioral constraints found. The best found
 - Total cost: 964
 - Use the explore method of the result to explore the implementation.

Use the explore method to explore the design containing the optimized fixed-point data types.

```
explore(result);
```

Producing Lookup Table Data

- “Producing Lookup Table Data” on page 41-2
- “Worst-Case Error for a Lookup Table” on page 41-3
- “Create Lookup Tables for a Sine Function” on page 41-5
- “Use Lookup Table Approximation Functions” on page 41-14
- “Optimize Lookup Tables for Memory-Efficiency” on page 41-15
- “Optimize Lookup Tables for Memory-Efficiency Programmatically” on page 41-19
- “Generate an Optimized Lookup Table as a MATLAB Function” on page 41-36
- “Generate an Optimized Lookup Table as a MATLAB Function Programmatically” on page 41-38
- “Convert Neural Network Algorithms to Fixed-Point Using fxpopt and Generate HDL Code” on page 41-41
- “Convert Neural Network Algorithms to Fixed Point and Generate C Code” on page 41-52
- “Effects of Spacing on Speed, Error, and Memory Usage” on page 41-59
- “Approximate Functions with a Direct Lookup Table” on page 41-65
- “Convert Digit Recognition Neural Network to Fixed Point and Generate C Code” on page 41-70
- “Calculate Complex dB Using a Direct Lookup Table” on page 41-79
- “Optimize Lookup Tables for Periodic Functions” on page 41-82
- “Replace Fitted Curve with Optimized Lookup Table” on page 41-89

Producing Lookup Table Data

A function lookup table is a method by which you can approximate a function by a table with a finite number of points (X,Y). Function lookup tables are essential to many fixed-point applications. The function you want to approximate is called the *ideal function*. The X values of the lookup table are called the *breakpoints*. You approximate the value of the ideal function at a point by linearly interpolating between the two breakpoints closest to the point.

In creating the points for a function lookup table, you generally want to achieve one or both of the following goals:

- Minimize the worst-case error for a specified maximum number of breakpoints
- Minimize the number of breakpoints for a specified maximum allowed error

“Create Lookup Tables for a Sine Function” on page 41-5 shows you how to create function lookup tables using the function `fixpt_look1_func_approx`. You can optimize the lookup table to minimize the number of data points, the error, or both. You can also restrict the spacing of the breakpoints to be even or even powers of two to speed up computations using the table.

“Worst-Case Error for a Lookup Table” on page 41-3 explains how to use the function `fixpt_look1_func_plot` to find the worst-case error of a lookup table and plot the errors at all points.

Worst-Case Error for a Lookup Table

The error at any point of a function lookup table is the absolute value of the difference between the ideal function at the point and the corresponding Y value found by linearly interpolating between the adjacent breakpoints. The *worst-case error*, or *maximum absolute error*, of a lookup table is the maximum absolute value of all errors in the interval containing the breakpoints.

For example, if the ideal function is the square root, and the breakpoints of the lookup table are 0, 0.25, and 1, then in a perfect implementation of the lookup table, the worst-case error is $1/8 = 0.125$, which occurs at the point $1/16 = 0.0625$. In practice, the error could be greater, depending on the fixed-point quantization and other factors.

The section that follows shows how to use the function `fixpt_look1_func_plot` to find the worst-case error of a lookup table for the square root function.

Approximate the Square Root Function

This example shows how to use the function `fixpt_look1_func_plot` to find the maximum absolute error for the simple lookup table whose breakpoints are 0, 0.25, and 1. The corresponding Y data points of the lookup table, which you find by taking the square roots of the breakpoints, are 0, 0.5, and 1.

To use the function `fixpt_look1_func_plot`, you need to define its parameters first. To do so, type the following at the MATLAB prompt:

```
funcstr = 'sqrt(x)'; % Define the square root function
xdata = [0;.25;1]; % Set the breakpoints
ydata = sqrt(xdata); % Find the square root of the breakpoints
xmin = 0; % Set the minimum breakpoint
xmax = 1; % Set the maximum breakpoint
xdt = ufix(16); % Set the x data type
xscale = 2^-16; % Set the x data scaling
ydt = sfix(16); % Set the y data type
yscale = 2^-14; % Set the y data scaling
rndmeth = 'Floor'; % Set the rounding method
```

To get the worst-case error of the lookup table and a plot of the error, type:

```
errworst = fixpt_look1_func_plot(xdata,ydata,funcstr, ...
xmin,xmax,xdt,xscale,ydt,yscale,rndmeth)
```

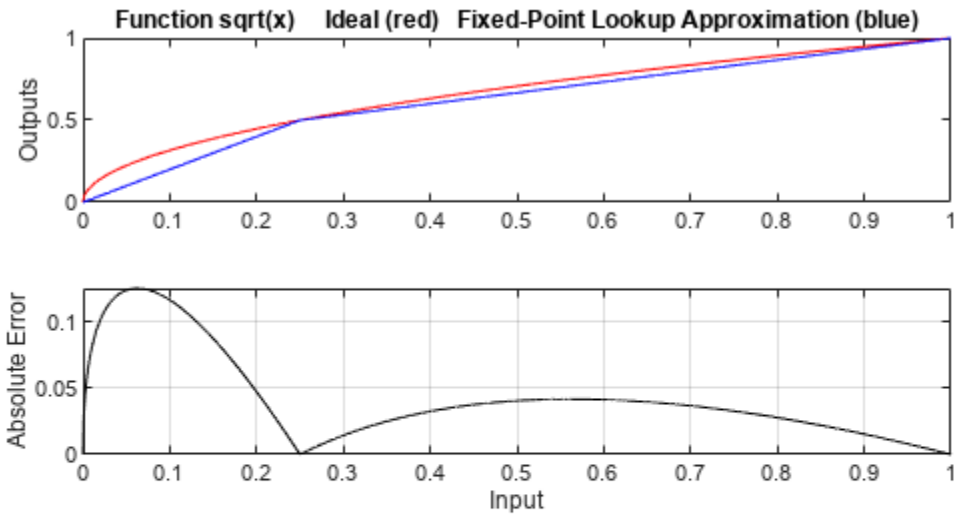


Table uses 3 unevenly spaced data points.
 The input is Unsigned 16 Bit with 16 bits right of binary point
 The output is Signed 16 Bit with 14 bits right of binary point
 Maximum Absolute Error 0.12503 $\log_2(\text{MAE}) = -2.9996$ MAE/yBit = 2048.5
 The least significant 12 bits of the output can be inaccurate.
 The most significant nonsign bit of the output is used.
 The remaining 3 nonsign bits of the output are used and always accurate.
 The sign bit of the output is not used.
 The rounding mode is to Floor

errworst = 0.1250

The upper box (Outputs) displays a plot of the square root function with a plot of the fixed-point lookup approximation underneath. The approximation is found by linear interpolation between the breakpoints. The lower box (Absolute Error) displays the errors at all points in the interval from 0 to 1. Notice that the maximum absolute error occurs at 0.0625. The error at the breakpoints is 0.

Create Lookup Tables for a Sine Function

In this section...

“Introduction” on page 41-5

“Set Function Parameters for the Lookup Table” on page 41-5

“Specifying Both `errmax` and `nptsmax`” on page 41-12

“Comparison of Example Results” on page 41-13

Introduction

The sections that follow explain how to use the function `fixpt_look1_func_approx` to create lookup tables. It gives examples that show how to create lookup tables for the function $\sin(2\pi x)$ on the interval from 0 to 0.25.

Set Function Parameters for the Lookup Table

First, define parameter values for the `fixpt_look1_func_approx` function.

```
% Required parameters
funcstr = 'sin(2*pi*x)'; % Ideal function
xmin = 0; % Minimum input of interest
xmax = 0.25; % Maximum input of interest
xdt = ufix(16); % x data type
xscale = 2^-16; % x data scaling
ydt = sfix(16); % y data type
yscale = 2^-14; % y data scaling
rndmeth = 'Floor'; % Rounding method
% Optional parameters
errmax = 2^-10; % Maximum allowed error of the lookup table
nptsmax = 21; % Maximum number of points of the lookup table
```

The parameters `errmax`, `nptsmax`, and `spacing` are optional. You must use at least one of the parameters `errmax` and `nptsmax`. If you use only the `errmax` parameter without `nptsmax`, the function creates a lookup table with the fewest points for which the worst-case error is at most `errmax`. If you use only the `nptsmax` parameter without `errmax`, the function creates a lookup table with at most `nptsmax` points which has the smallest worst-case error. You can use the optional `spacing` parameter to restrict the spacing between breakpoints of the lookup table.

Use `errmax` with Unrestricted Spacing

Create a lookup table that has the fewest data points for a specified worst-case error, with unrestricted spacing.

```
[xdata,ydata,errworst] = fixpt_look1_func_approx(funcstr, ...
xmin,xmax,xdt,xscale,ydt,yscale,rndmeth,errmax,[]);
```

Note that the `nptsmax` and `spacing` parameters are not specified.

```
length(xdata)
```

```
ans = 16
```

16 points are required to approximate $\sin(2\pi x)$ to within the tolerance specified by `errmax`. The maximum error is:

```
errworst
```

```
errworst = 9.7656e-04
```

The value of `errworst` is less than or equal to the value of `errmax`.

Plot the results:

```
figure(1)
fixpt_look1_func_plot(xdata,ydata,funcstr,xmin,xmax,xdt, ...
xscale,ydt,yscale,rndmeth);
```

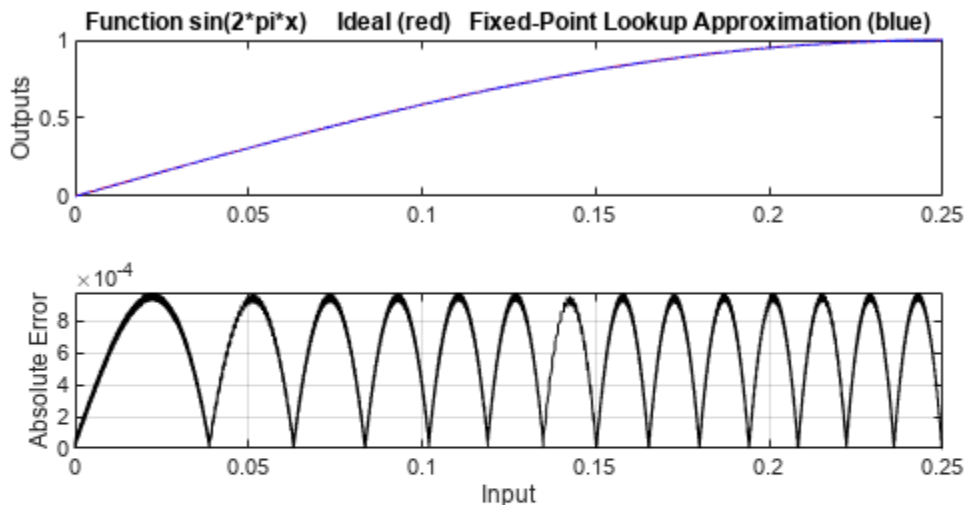


Table uses 16 unevenly spaced data points.

The input is Unsigned 16 Bit with 16 bits right of binary point

The output is Signed 16 Bit with 14 bits right of binary point

Maximum Absolute Error 0.00097656 $\log_2(\text{MAE}) = -10$ MAE/yBit = 16

The least significant 4 bits of the output can be inaccurate.

The most significant nonsign bit of the output is used.

The remaining 11 nonsign bits of the output are used and always accurate.

The sign bit of the output is not used.

The rounding mode is to Floor

The upper plot shows the ideal function $\sin(2\pi x)$ and the fixed-point lookup table approximation between the breakpoints. In this example, the ideal function and the approximation are so close together that the two graphs appear to coincide. The lower plot displays the errors.

In this example, the Y data points, `ydata`, are equal to the ideal function applied to the point in `xdata`. However, you can define a different set of values for `ydata` after running `fixpt_look1_func_approx`. This can sometimes reduce the maximum error.

You can also change the values of `xmin` and `xmax` to evaluate the lookup table on a subset of the original interval.

To find the new maximum error after changing `ydata`, `xmin`, or `xmax`, type

```
errworst = fixpt_look1_func_plot(xdata,ydata,funcstr,xmin,xmax, ...
xdt,xscale,ydt,yscale,rndmeth)
```

Use nptsmax with Unrestricted Spacing

Create a lookup table that minimizes the worst-case error for a specified maximum number of data points, with unrestricted spacing.

```
[xdata, ydata, errworst] = fixpt_look1_func_approx(funcstr, ...
xmin,xmax,xdt,xscale,ydt,yscale,rndmeth,[],nptsmax);
```

The empty brackets tell the function to ignore the parameter `errmax`. Omitting `errmax` causes the function to return a lookup table of size specified by `nptsmax`, with the smallest worst-case error.

```
length(xdata)
```

```
ans = 21
```

The function returns a vector `xdata` with 21 points. The maximum error is:

```
errworst
```

```
errworst = 5.1139e-04
```

The value of `errworst` is less than or equal to the value of `errmax`.

Plot the results:

```
figure(2)
fixpt_look1_func_plot(xdata,ydata,funcstr,xmin,xmax,xdt, ...
xscale,ydt,yscale,rndmeth);
```

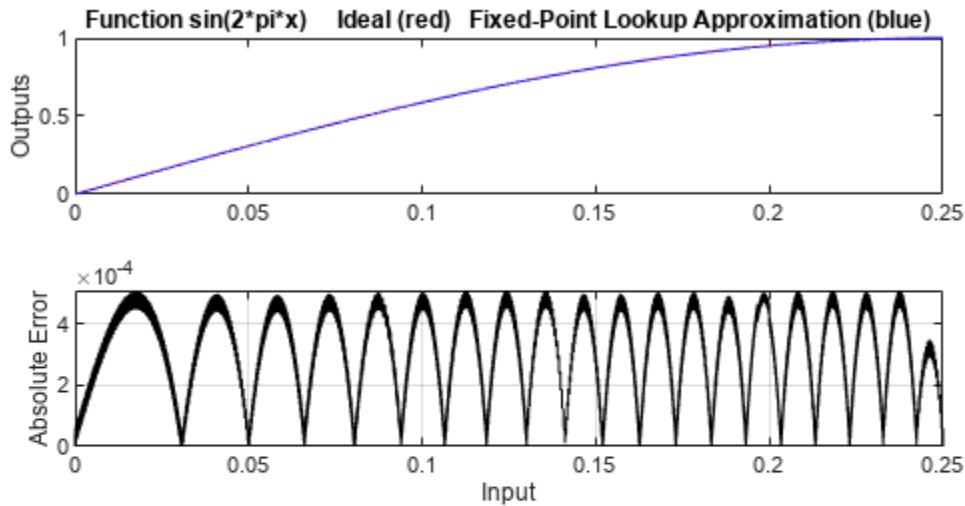


Table uses 21 unevenly spaced data points.
 The input is Unsigned 16 Bit with 16 bits right of binary point
 The output is Signed 16 Bit with 14 bits right of binary point
 Maximum Absolute Error 0.00051139 $\log_2(\text{MAE}) = -10.9333$ MAE/yBit = 8.3785
 The least significant 4 bits of the output can be inaccurate.
 The most significant nonsign bit of the output is used.
 The remaining 11 nonsign bits of the output are used and always accurate.
 The sign bit of the output is not used.
 The rounding mode is to Floor

Restricting the Spacing

In the previous two examples, the function `fixpt_look1_func_approx` creates lookup tables with unrestricted spacing between the points. You can restrict the spacing to improve the computational efficiency of the lookup table using the spacing parameter. Both power of two, 'pow2', and even spacing, 'even', increase the computational speed of the lookup table and use less command read-only memory (ROM). However, specifying either of the spacing restrictions along with `errmax` usually requires more data points in the lookup table than does unrestricted spacing to achieve the same degree of accuracy. See *Effects of Spacing on Speed, Error, and Memory Usage* for more information.

Use `errmax` with Even Spacing

Create a lookup table that has evenly spaced breakpoints and a specified worst-case error.

```
spacing = 'even';
[xdata,ydata,errworst] = fixpt_look1_func_approx(funcstr, ...
xmin,xmax,xdt,xscale,ydt,yscale,rndmeth,errmax,[],spacing);
length(xdata)
```

```
ans = 20
```

```
errworst
```

```
errworst = 9.2109e-04
```

Plot the results:

```
figure(3)
fixpt_look1_func_plot(xdata,ydata,funcstr,xmin,xmax,xdt, ...
xscale,ydt,yscale,rndmeth);
```

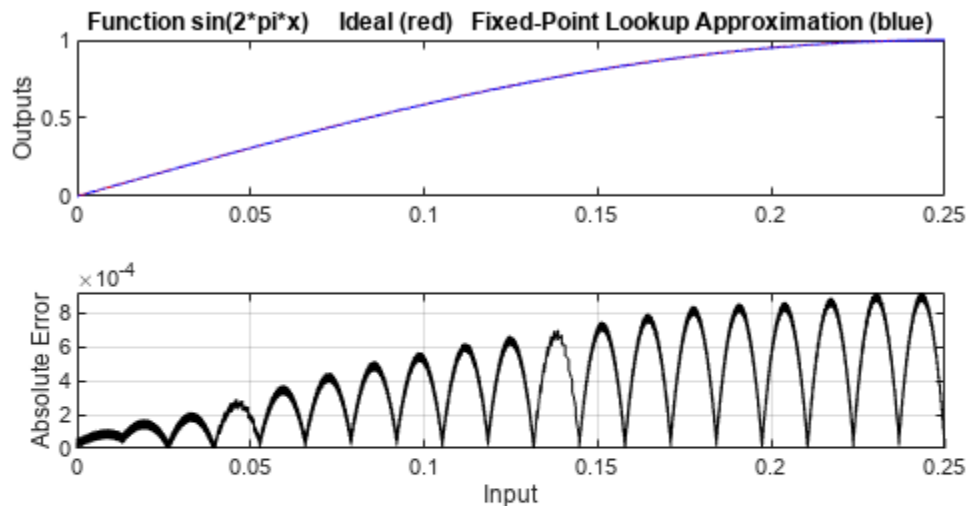


Table uses 20 evenly spaced data points.
The input is Unsigned 16 Bit with 16 bits right of binary point
The output is Signed 16 Bit with 14 bits right of binary point
Maximum Absolute Error 0.00092109 $\log_2(\text{MAE}) = -10.0844$ MAE/yBit = 15.0912
The least significant 4 bits of the output can be inaccurate.
The most significant nonsign bit of the output is not used.
The remaining 10 nonsign bits of the output are used and always accurate.
The sign bit of the output is not used.
The rounding mode is to Floor

Use nptsmax with Even Spacing

Create a lookup table that has evenly spaced breakpoints and minimizes the worst-case error for a specified maximum number of points.

```
spacing = 'even';
[xdata, ydata, errworst] = fixpt_look1_func_approx(funcstr, ...
xmin, xmax, xdt, xscale, ydt, yscale, rndmeth, [], nptsmax, spacing);
length(xdata)

ans = 21

errworst

errworst = 8.3793e-04
```

The result requires 21 evenly spaced points to achieve a maximum absolute error of $2^{-10.2209}$.

Plot the results:

```
figure(4)
fixpt_look1_func_plot(xdata,ydata,funcstr,xmin,xmax,xdt, ...
xscale,ydt,yscale,rndmeth);
```

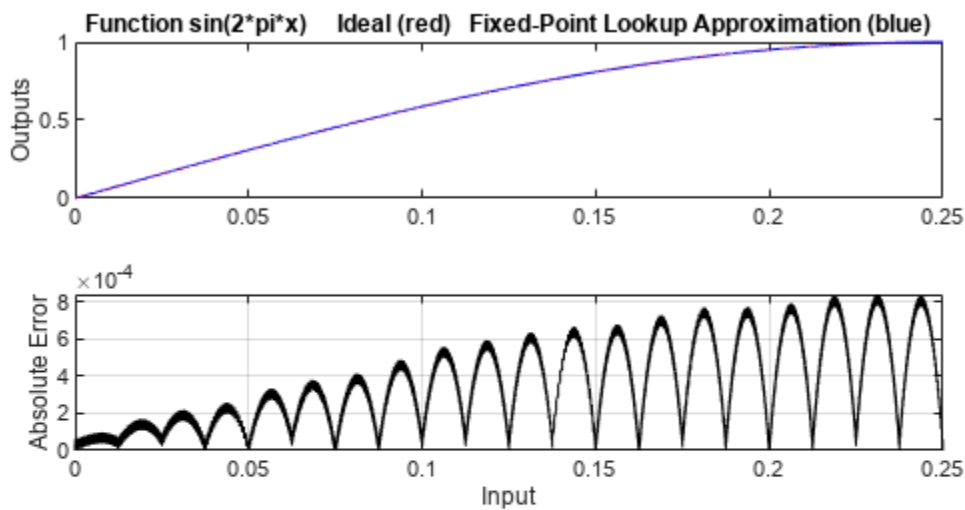


Table uses 21 evenly spaced data points.

The input is Unsigned 16 Bit with 16 bits right of binary point

The output is Signed 16 Bit with 14 bits right of binary point

Maximum Absolute Error 0.00083793 $\log_2(\text{MAE}) = -10.2209$ MAE/yBit = 13.7287

The least significant 4 bits of the output can be inaccurate.

The most significant nonsign bit of the output is not used.

The remaining 10 nonsign bits of the output are used and always accurate.

The sign bit of the output is not used.

The rounding mode is to Floor

Use `errmax` with Power of Two Spacing

Create a lookup table that has power of two spacing and a specified worst-case error.

```
spacing = 'pow2';
[xdata, ydata, errworst] = ...
fixpt_look1_func_approx(funcstr,xmin, ...
xmax,xdt,xscale,ydt,yscale,rndmeth,errmax,[],spacing);
length(xdata)
```

```
ans = 33
```

33 points are required to achieve the worst-case error specified by `errmax`. To verify that these points are evenly spaced, type

```
widths = diff(xdata)
```

```
widths = 32x1
```

```
0.0078
0.0078
0.0078
0.0078
0.0078
0.0078
0.0078
0.0078
0.0078
0.0078
```

```
0.0078
:
```

This generates a vector whose entries are the differences between the consecutive points in `xdata`. Every entry of `widths` is 2^{-7} .

The maximum error is:

```
errworst
```

```
errworst = 3.7209e-04
```

This is less than the value of `errmax`.

Plot the results:

```
figure(5)
fixpt_look1_func_plot(xdata,ydata,funcstr,xmin,xmax,xdt, ...
xscale,ydt,yscale,rndmeth);
```

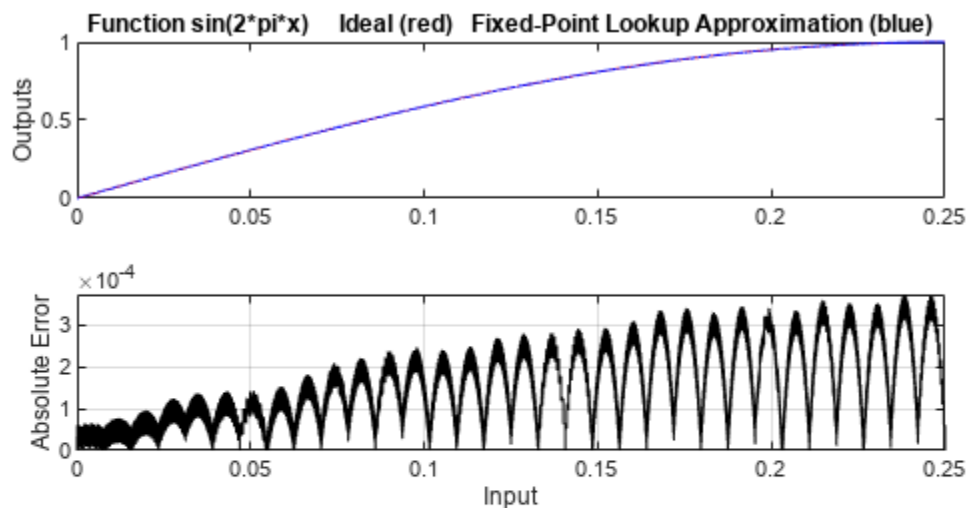


Table uses 33 power of 2 spaced data points.
The input is Unsigned 16 Bit with 16 bits right of binary point
The output is Signed 16 Bit with 14 bits right of binary point
Maximum Absolute Error 0.00037209 $\log_2(\text{MAE}) = -11.3921$ MAE/yBit = 6.0964
The least significant 3 bits of the output can be inaccurate.
The most significant nonsign bit of the output is used.
The remaining 12 nonsign bits of the output are used and always accurate.
The sign bit of the output is not used.
The rounding mode is to Floor

Use `nptsmax` with Power of Two Spacing

Create a lookup table that has power of two spacing and minimizes the worst-case error for a specified number of points.

```
spacing = 'pow2';
[xdata, ydata, errworst] = ...
```

```
fixpt_look1_func_approx(funcstr,xmin, ...
xmax,xdt,xscale,ydt,yscale,rndmeth,[],nptsmax,spacing);
length(xdata)
```

```
ans = 17
```

```
errworst
```

```
errworst = 0.0013
```

The result requires 17 points to achieve a maximum absolute error of $2^{-9.6267}$.

Plot the results:

```
figure(6)
```

```
fixpt_look1_func_plot(xdata,ydata,funcstr,xmin,xmax,xdt, ...
xscale,ydt,yscale,rndmeth);
```

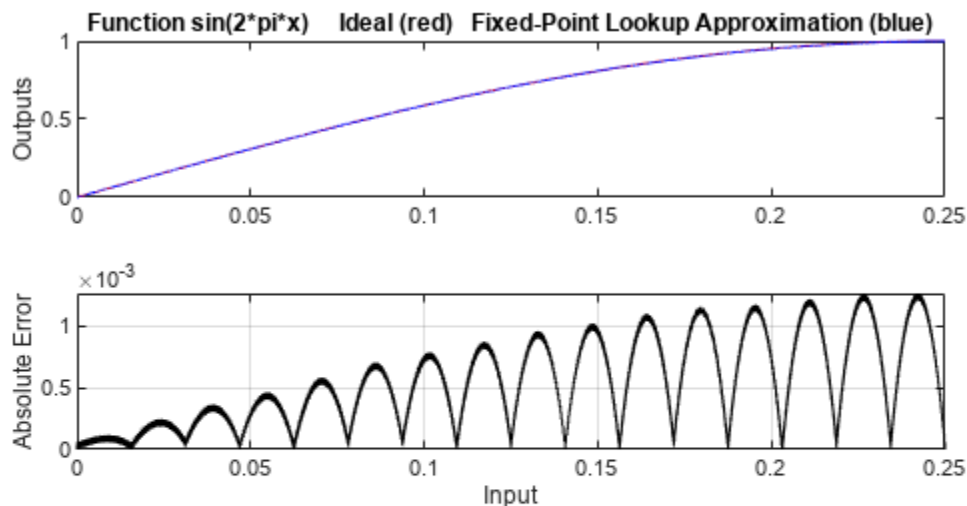


Table uses 17 power of 2 spaced data points.

The input is Unsigned 16 Bit with 16 bits right of binary point

The output is Signed 16 Bit with 14 bits right of binary point

Maximum Absolute Error 0.0012649 $\log_2(\text{MAE}) = -9.6267$ MAE/yBit = 20.7245

The least significant 5 bits of the output can be inaccurate.

The most significant nonsign bit of the output is used.

The remaining 10 nonsign bits of the output are used and always accurate.

The sign bit of the output is not used.

The rounding mode is to Floor

Specifying Both errmax and nptsmax

If you include both the errmax and the nptsmax parameters, the function

fixpt_look1_func_approx tries to find a lookup table with at most nptsmax data points, whose worst-case error is at most errmax. If it can find a lookup table meeting both conditions, it uses the following order of priority for spacing:

- 1 Power of two

- 2 Even
- 3 Unrestricted

If the function cannot find any lookup table satisfying both conditions, it ignores `nptsmax` and returns a lookup table with unrestricted spacing, whose worst-case error is at most `errmax`. In this case, the function behaves the same as if the `nptsmax` parameter were omitted.

The following examples illustrate the results of using different values for `nptsmax` when you enter

```
[xdata ydata errworst] = fixpt_look1_func_approx(funcstr, ...
xmin,xmax,xdt,xscale,ydt,yscale,rndmeth,errmax,nptsmax);
```

The results for three different settings for `nptsmax` are as follows:

- `nptsmax = 33`; — The function creates the lookup table with 33 points having power of two spacing, as in Example 3.
- `nptsmax = 21`; — Because the `errmax` and `nptsmax` conditions cannot be met with power of two spacing, the function creates the lookup table with 20 points having even spacing, as in Example 5.
- `nptsmax = 16`; — Because the `errmax` and `nptsmax` conditions cannot be met with either power of two or even spacing, the function creates the lookup table with 16 points having unrestricted spacing, as in Example 1.

Comparison of Example Results

The following table summarizes the results for the examples. When you specify `errmax`, even spacing requires more data points than unrestricted, and power of two spacing requires more points than even spacing.

| Example | Options | Spacing | Worst-Case Error | Number of Points in Table |
|---------|---------------------------|----------------|------------------|---------------------------|
| 1 | <code>errmax=2^-10</code> | 'unrestricted' | 2^{-10} | 16 |
| 2 | <code>nptsmax=21</code> | 'unrestricted' | $2^{-10.933}$ | 21 |
| 3 | <code>errmax=2^-10</code> | 'even' | $2^{-10.0844}$ | 20 |
| 4 | <code>nptsmax=21</code> | 'even' | $2^{-10.2209}$ | 21 |
| 5 | <code>errmax=2^-10</code> | 'pow2' | $2^{-11.3921}$ | 33 |
| 6 | <code>nptsmax=21</code> | 'pow2' | $2^{-9.627}$ | 17 |

Use Lookup Table Approximation Functions

The following steps summarize how to use the lookup table approximation functions.

- 1** Define:
 - a** The ideal function to approximate
 - b** The range, `xmin` to `xmax`, over which to find X and Y data
 - c** The fixed-point implementation: data type, scaling, and rounding method
 - d** The maximum acceptable error, the maximum number of points, and the spacing
- 2** Run the `fixpt_look1_func_approx` function to generate X and Y data.
- 3** Use the `fixpt_look1_func_plot` function to plot the function and error between the ideal and approximated functions using the selected X and Y data, and to calculate the error and the number of points used.
- 4** Vary input criteria, such as `errmax`, `nptsmax`, and `spacing`, to produce sets of X and Y data that generate functions with varying worst-case error, number of points required, and spacing.
- 5** Compare results of the number of points required and maximum absolute error from various runs to choose the best set of X and Y data.

Optimize Lookup Tables for Memory-Efficiency

The **Lookup Table Optimizer** optimizes the spacing of breakpoints and the data types of lookup table data to reduce the memory used by a lookup table. Using the **Lookup Table Optimizer** and its command-line equivalent, you can:

- Optimize an existing Lookup Table block.
- Generate a lookup table from a Simulink block, including a Math Function block or a subsystem.
- Generate a lookup table from a function or function handle.

Optimize an Existing Lookup Table Using the Lookup Table Optimizer

To optimize an existing lookup table, open the model containing the Lookup Table block.

```
openExample('simulink_automotive/ModelingAFaultTolerantFuelControlSystemExample',...
    'supportingfile','sldemo_fuelsys');
open_system('sldemo_fuelsys/fuel_rate_control/airflow_calc');
```

This example shows how to optimize the Pumping Constant Lookup Table block.


- 1 To open the Lookup Table Optimizer, select the Pumping Constant Lookup Table block. The context-sensitive **Lookup Table** tab appears in the Simulink toolstrip. In the **Lookup Table** tab, select **Lookup Table Optimizer**.
- 2 Select the type of block you want to optimize. To optimize a Simulink block or subsystem, including an existing Lookup Table block or a Math Function block, select **Simulink block or subsystem**. To generate a lookup table approximation for a function handle, select **MATLAB Function Handle**.

In this example, select **Simulink block or subsystem** to optimize the Pumping Constant lookup table. Click **Next**.

- 3 Under **Block Information**, enter the path to the Pumping Constant Lookup Table block. Select the block in the model, then click **Get Current Block** in the Lookup Table Optimizer to fill in the block path automatically.
- 4 Click **Collect Current Values from Model** to update the model diagram and allow the Lookup Table Optimizer to automatically gather information needed for the optimization process including current output data type, and input number, data type, and value range. You can manually edit all of these fields to specify ranges and data types other than those currently specified on the block.
 - Specify the **Desired Output Data Type** of the generated lookup table as a `numericType` or `Simulink.NumericType` object.
 - Specify the data type of each input to the block as a `numericType` or `Simulink.NumericType` object.
 - Specify the minimum and maximum values of each input of the generated lookup table as scalars in the table.

For this example, use the current values specified on the model. Click **Next**.

- 5 Specify constraints to use in the optimization. Set the **Output Error Tolerance** that is acceptable for your design.
 - Absolute tolerance is defined as the absolute value of the difference between the original output value and the output value of the optimized lookup table.

- Relative tolerance measures the error relative to the value at that point, specified as a non-negative.
- 6 Specify the allowed word lengths as a vector based on types that are efficient for your intended hardware target. For example, if you want to allow the optimizer to consider only 8-, 16-, and 32-bit types, specify [8 16 32] in the **Allowed Word Lengths (Vector)** field.
 - 7 To specify additional properties for the optimized lookup table, click **LUT Specification**. For more information on each of the properties, see `FunctionApproximation.Options`. In this example, use the default values for these properties.
 - 8 Specify options for the optimization, such as the maximum time or maximum memory usage for the generated lookup table by clicking the  button.
 - 9 After you set the constraints, click **Optimize**.

Optionally, you can choose to stop the optimization solver before the optimization solver is complete by clicking **Stop**. The optimizer will choose the best solution found at the time the **Stop** button is selected and display it in the app.

When the optimization is complete, the optimizer reports the memory of the optimized lookup table. You can edit the constraints and run the optimization again to achieve further memory reduction.

Using the default settings, the Lookup Table Optimizer reduces the memory used by the Pumping Constant Lookup Table block from 1516 bytes to 505 bytes (66.69%).

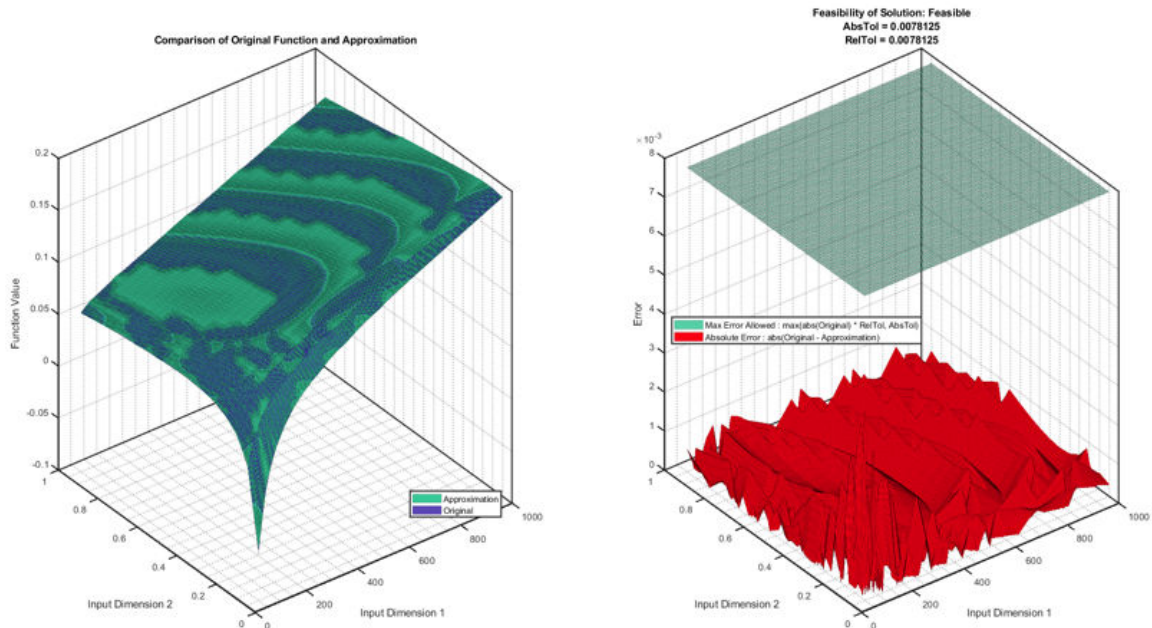
Memory Reduced by 66.69%

| LUT Attributes | Old Memory* | Old Data Type | New Memory* | New Data Type |
|----------------|-------------|-----------------------|-------------|----------------------|
| Table Data | 1368 | numerictype('single') | 493 | numerictype(1,8,9) |
| Breakpoint 1 | 72 | numerictype('single') | 4 | numerictype(0,16,6) |
| Breakpoint 2 | 76 | numerictype('single') | 8 | numerictype(0,32,32) |
| Total | 1516 | | 505 | |

* In bytes

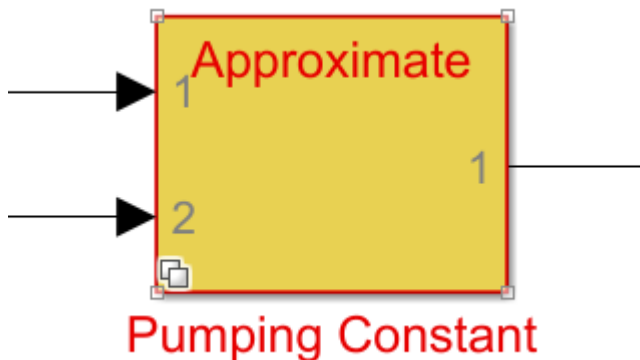
Click **Next**.

- 10 Click **Show Comparison Plot** to view a plot of the original block output compared to the output of the new optimized lookup table.



Click **Replace Original Function** to generate a new lookup table using the optimized settings found by the app, and replace the original block.

The new block is a masked variant subsystem in which the active variant is the optimized lookup table block. The inactive variant is the original block.



Edit the Optimization Settings and Generate a New Approximate

You can iteratively change the approximation block by editing the settings used during the optimization to generate a new lookup table.

- 1 Double-click the Pumping Constant block. To edit the optimization settings, in the Block Parameters dialog, click **Redesign approximate**.
- 2 In the Lookup Table Optimizer, click **Next** to proceed to the **Create** page of the app. In this example, edit the absolute and relative tolerances to a slightly larger value so that you can further reduce the size of the lookup table.

- Set the **Absolute** tolerance to 0.01, or 1%.
- Set the **Relative** tolerance to 0.01, or 1%.

3 Click **Optimize** to optimize the lookup table with the new options.

Using these tolerance values, the new lookup table uses only 304 bytes of memory.

4 Click **Next**. On the **Results** page, click the **Replace Original Function** button to replace the first iteration of the approximation block with this newest iteration.

5 In the model, double-click the Pumping Constant block to open the Block Parameters. The Block Parameters displays the settings used for the approximation

To make the original block or subsystem the active variant, next to **Select desired function version**, select **Original**.

To delete the lookup table approximation from the model, in the Block Parameters, click **Revert to original**.

See Also

Apps

Lookup Table Optimizer

Classes

`FunctionApproximation.Problem` | `FunctionApproximation.Options` |

`FunctionApproximation.LUTSolution` |

`FunctionApproximation.LUTMemoryUsageCalculator`

More About

- “Optimize Lookup Tables for Memory-Efficiency Programmatically” on page 41-19
- “Generate an Optimized Lookup Table as a MATLAB Function” on page 41-36

Optimize Lookup Tables for Memory-Efficiency Programmatically

The following examples show how to generate memory-efficient lookup tables programmatically. Using the command-line equivalent of the **Lookup Table Optimizer**, you can:

- Optimize an existing Lookup Table block.
- Generate a lookup table from a Math Function block.
- Generate a lookup table from a function or function handle.
- Generate a lookup table from a Subsystem block.

Approximate a Function Using a Lookup Table

This example shows how to generate a memory-efficient lookup table that approximates the `sin` function. Define the approximation problem by creating a `FunctionApproximation.Problem` object.

```
P = FunctionApproximation.Problem('sin')
```

```
P =
```

```
1x1 FunctionApproximation.Problem with properties:
    FunctionToApproximate: @(x)sin(x)
      NumberOfInputs: 1
        InputTypes: "numerictype(0,16,13)"
    InputLowerBounds: 0
    InputUpperBounds: 6.2832
      OutputType: "numerictype(1,16,14)"
      Options: [1x1 FunctionApproximation.Options]
```

The `FunctionToApproximate` and `NumberOfInputs` properties of the `Problem` object are inferred from the definition of the object, and cannot be edited after creation. All other properties are writable.

Edit the `FunctionApproximation.Options` object to specify additional constraints to use in the optimization process. For example, constrain the breakpoints of the generated lookup table to even spacing.

```
P.Options.BreakpointSpecification = 'EvenSpacing'
```

```
P =
```

```
1x1 FunctionApproximation.Problem with properties:
    FunctionToApproximate: @(x)sin(x)
      NumberOfInputs: 1
        InputTypes: "numerictype(0,16,13)"
    InputLowerBounds: 0
    InputUpperBounds: 6.2832
```

```
OutputType: "numerictype(1,16,14)"
Options: [1x1 FunctionApproximation.Options]
```

Specify additional constraints, such as the absolute and relative tolerances of the output, and word length constraints.

```
P.Options.AbsTol = 2^-10;
P.Options.RelTol = 2^-6;
P.Options.WordLengths = [8,16];
```

Use the solve method to solve the optimization problem. MATLAB™ displays the iterations of the optimization process. The solve method returns a FunctionApproximation.LUTSolution object.

```
S = solve(P)
```

Searching for fixed-point solutions.

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|---------------|----------|------------|-----------------|--------------|----------------|
| 0 | 48 | 0 | 2 | 8 | 16 | EvenPow |
| 1 | 32 | 0 | 2 | 8 | 8 | EvenPow |
| 2 | 672 | 0 | 41 | 8 | 16 | EvenPow |
| 3 | 560 | 0 | 34 | 8 | 16 | EvenPow |
| 4 | 1104 | 0 | 68 | 8 | 16 | EvenPow |
| 5 | 1648 | 1 | 102 | 8 | 16 | EvenPow |
| 6 | 480 | 0 | 29 | 8 | 16 | EvenPow |
| 7 | 832 | 0 | 51 | 8 | 16 | EvenPow |
| 8 | 320 | 0 | 19 | 8 | 16 | EvenPow |
| 9 | 64 | 0 | 2 | 16 | 16 | EvenPow |
| 10 | 48 | 0 | 2 | 16 | 8 | EvenPow |
| 11 | 640 | 1 | 38 | 16 | 16 | EvenPow |
| 12 | 624 | 0 | 37 | 16 | 16 | EvenPow |
| 13 | 496 | 0 | 29 | 16 | 16 | EvenPow |
| 14 | 480 | 0 | 28 | 16 | 16 | EvenPow |
| 15 | 560 | 0 | 33 | 16 | 16 | EvenPow |
| 16 | 592 | 0 | 35 | 16 | 16 | EvenPow |
| 17 | 608 | 0 | 36 | 16 | 16 | EvenPow |
| 18 | 352 | 1 | 20 | 16 | 16 | EvenPow |
| 19 | 336 | 0 | 19 | 16 | 16 | EvenPow |
| 20 | 208 | 0 | 11 | 16 | 16 | EvenPow |
| 21 | 272 | 0 | 15 | 16 | 16 | EvenPow |
| 22 | 304 | 0 | 17 | 16 | 16 | EvenPow |
| 23 | 320 | 0 | 18 | 16 | 16 | EvenPow |
| 24 | 48 | 0 | 2 | 8 | 16 | EvenPow |
| 25 | 224 | 0 | 13 | 8 | 16 | EvenPow |
| 26 | 64 | 0 | 2 | 16 | 16 | EvenPow |
| 27 | 240 | 0 | 13 | 16 | 16 | EvenPow |
| 28 | 1648 | 1 | 102 | 8 | 16 | EvenPow |

Best Solution

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|---------------|----------|------------|-----------------|--------------|----------------|
| 18 | 352 | 1 | 20 | 16 | 16 | EvenPow |

S =

1x1 FunctionApproximation.LUTSolution with properties:


```
ID: 18
Feasible: "true"
```

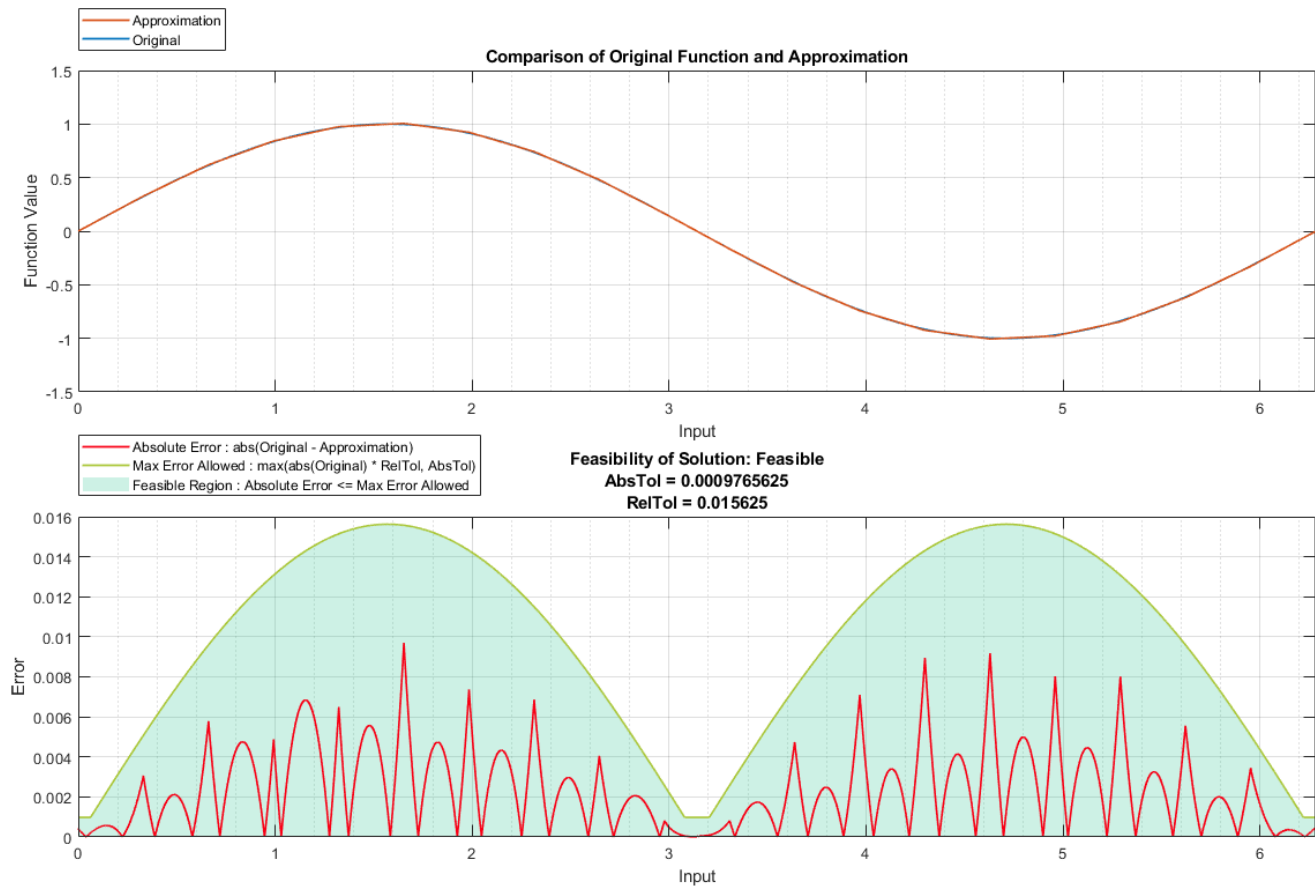
Compare the numerical behavior of the original function with the numerical behavior of the generated lookup table stored in the solution, S.

```
err = compare(S)
```

```
err =
```

```
struct with fields:
```

```
Breakpoints: [51473x1 double]
Original: [51473x1 double]
Approximate: [51473x1 double]
```



You can access the lookup table data stored in the LUTSolution object.

```
t = S.TableData
```

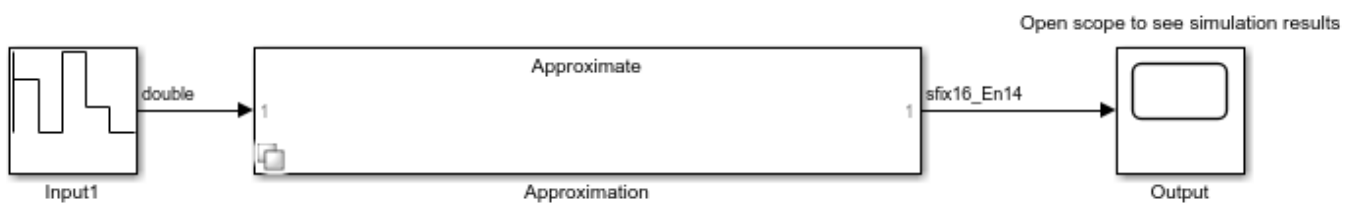
```
t =
```

struct with fields:

```
BreakpointValues: {[0 0.3307 0.6614 0.9921 1.3228 1.6534 1.9841 ... ]}
BreakpointDataTypes: [1x1 embedded.numerictype]
TableValues: [4.2725e-04 0.3278 0.6200 0.8420 0.9759 1.0063 ... ]
TableDataType: [1x1 embedded.numerictype]
IsEvenSpacing: 1
Interpolation: Linear
```

To access the generated Lookup Table block, use the `approximate` method.

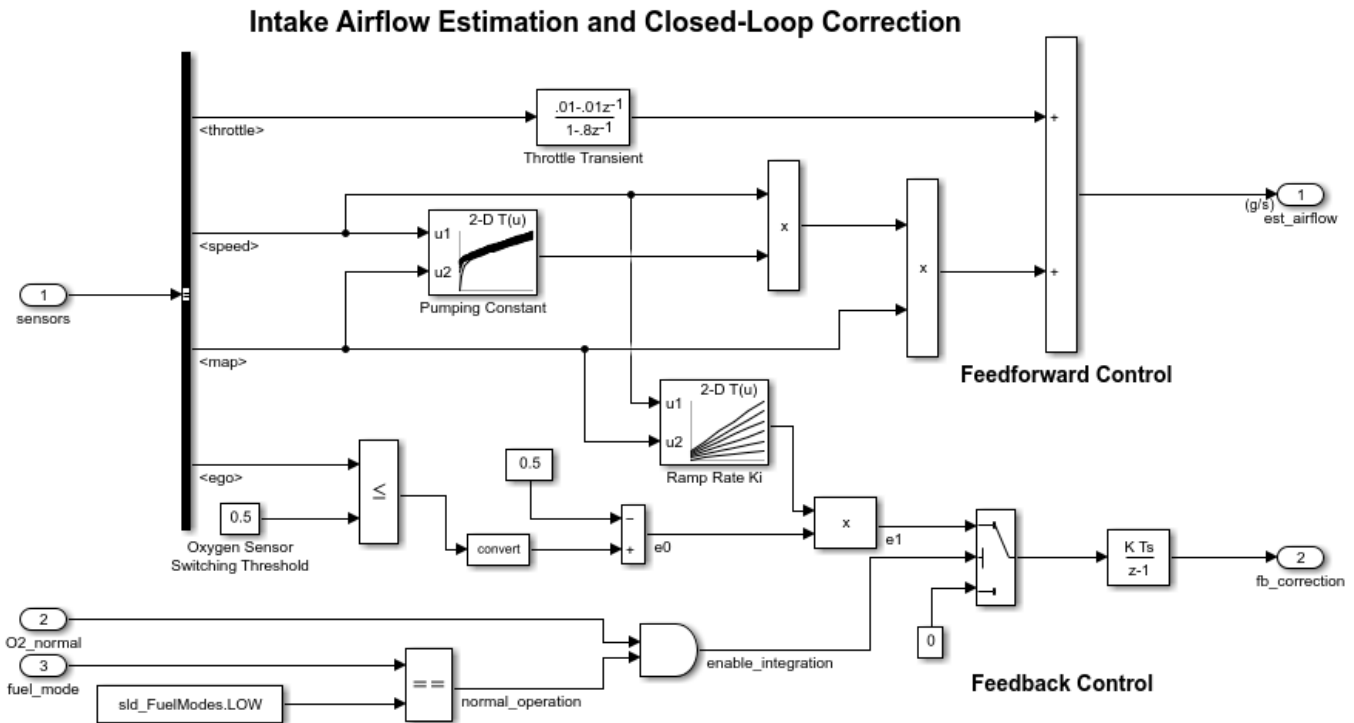
```
approximate(S)
```



Optimize an Existing Lookup Table

This example shows how to optimize an existing Lookup Table block for memory efficiency. Open the model containing the Lookup Table block that you want to optimize.

```
load_system('sldemo_fuelsys'); open_system('sldemo_fuelsys'); save_system('sldemo_fuelsys','my_s');
open_system('my_sldemo_fuelsys/fuel_rate_control/airflow_calc');
```



Create a `FunctionApproximation.Problem` object to define the optimization problem and constraints.

```
P = FunctionApproximation.Problem('my_slldemo_fuelsys/fuel_rate_control/airflow_calc/Pumping Constant')
```

```
P =
```

```
1x1 FunctionApproximation.Problem with properties:
```

```
FunctionToApproximate: 'my_slldemo_fuelsys/fuel_rate_control/airflow_calc/Pumping Constant'
NumberOfInputs: 2
InputTypes: ["numerictype('single')" ... ]
InputLowerBounds: [50 0.0500]
InputUpperBounds: [1000 0.9500]
OutputType: "numerictype('single')"
Options: [1x1 FunctionApproximation.Options]
```

Specify additional constraints by modifying the `Options` object associated with the `Problem` object, `P`.

```
P.Options.BreakpointSpecification = "EvenSpacing"
```

```
P =
```

```
1x1 FunctionApproximation.Problem with properties:
```

```
FunctionToApproximate: 'my_slldemo_fuelsys/fuel_rate_control/airflow_calc/Pumping Constant'
NumberOfInputs: 2
```

```

InputTypes: ["numerictype('single')" ... ]
InputLowerBounds: [50 0.0500]
InputUpperBounds: [1000 0.9500]
OutputType: "numerictype('single')"
Options: [1x1 FunctionApproximation.Options]

```

Solve the optimization problem.

```
S = solve(P)
```

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLs | TableData WL | BreakpointSpec: |
|----|---------------|----------|------------|-----------------|--------------|-----------------|
| 0 | 12128 | 1 | [18 19] | [32 32] | 32 | Expli |
| 1 | 11840 | 1 | [18 19] | [16 32] | 32 | Expli |
| 2 | 80 | 0 | [2 2] | [16 8] | 8 | Eve |
| 3 | 336 | 0 | [6 6] | [16 8] | 8 | Eve |
| 4 | 288 | 0 | [5 6] | [16 8] | 8 | Eve |
| 5 | 288 | 0 | [6 5] | [16 8] | 8 | Eve |
| 6 | 248 | 0 | [5 5] | [16 8] | 8 | Eve |
| 7 | 1016 | 0 | [11 11] | [16 8] | 8 | Eve |
| 8 | 928 | 0 | [10 11] | [16 8] | 8 | Eve |
| 9 | 928 | 0 | [11 10] | [16 8] | 8 | Eve |
| 10 | 848 | 0 | [10 10] | [16 8] | 8 | Eve |
| 11 | 3408 | 0 | [21 20] | [16 8] | 8 | Eve |
| 12 | 3248 | 0 | [20 20] | [16 8] | 8 | Eve |
| 13 | 7008 | 0 | [29 30] | [16 8] | 8 | Eve |
| 14 | 9024 | 0 | [33 34] | [16 8] | 8 | Eve |
| 15 | 9568 | 0 | [35 34] | [16 8] | 8 | Eve |
| 16 | 9840 | 0 | [36 34] | [16 8] | 8 | Eve |
| 17 | 11592 | 0 | [37 39] | [16 8] | 8 | Eve |
| 18 | 840 | 0 | [9 11] | [16 8] | 8 | Eve |
| 19 | 768 | 0 | [9 10] | [16 8] | 8 | Eve |
| 20 | 3088 | 0 | [19 20] | [16 8] | 8 | Eve |
| 21 | 2928 | 0 | [18 20] | [16 8] | 8 | Eve |
| 22 | 11280 | 0 | [36 39] | [16 8] | 8 | Eve |
| 23 | 1848 | 0 | [15 15] | [16 8] | 8 | Eve |
| 24 | 1728 | 0 | [14 15] | [16 8] | 8 | Eve |
| 25 | 1728 | 0 | [15 14] | [16 8] | 8 | Eve |
| 26 | 1616 | 0 | [14 14] | [16 8] | 8 | Eve |
| 27 | 6768 | 0 | [28 30] | [16 8] | 8 | Eve |
| 28 | 6080 | 0 | [29 26] | [16 8] | 8 | Eve |
| 29 | 5872 | 0 | [28 26] | [16 8] | 8 | Eve |
| 30 | 2784 | 0 | [19 18] | [16 8] | 8 | Eve |
| 31 | 2640 | 0 | [18 18] | [16 8] | 8 | Eve |
| 32 | 128 | 0 | [2 2] | [16 32] | 8 | Eve |
| 33 | 384 | 0 | [6 6] | [16 32] | 8 | Eve |
| 34 | 336 | 0 | [5 6] | [16 32] | 8 | Eve |
| 35 | 336 | 0 | [6 5] | [16 32] | 8 | Eve |
| 36 | 296 | 0 | [5 5] | [16 32] | 8 | Eve |
| 37 | 1064 | 0 | [11 11] | [16 32] | 8 | Eve |
| 38 | 976 | 0 | [10 11] | [16 32] | 8 | Eve |
| 39 | 976 | 0 | [11 10] | [16 32] | 8 | Eve |
| 40 | 896 | 0 | [10 10] | [16 32] | 8 | Eve |
| 41 | 3624 | 0 | [21 21] | [16 32] | 8 | Eve |
| 42 | 3456 | 0 | [20 21] | [16 32] | 8 | Eve |
| 43 | 3456 | 0 | [21 20] | [16 32] | 8 | Eve |
| 44 | 3296 | 0 | [20 20] | [16 32] | 8 | Eve |
| 45 | 6824 | 1 | [29 29] | [16 32] | 8 | Eve |

| | | | | | | |
|-----|------|---|---------|---------|----|----|
| 46 | 5096 | 0 | [25 25] | [16 32] | 8 | Ev |
| 47 | 5928 | 0 | [27 27] | [16 32] | 8 | Ev |
| 48 | 6368 | 0 | [28 28] | [16 32] | 8 | Ev |
| 49 | 888 | 0 | [9 11] | [16 32] | 8 | Ev |
| 50 | 816 | 0 | [9 10] | [16 32] | 8 | Ev |
| 51 | 3288 | 0 | [19 21] | [16 32] | 8 | Ev |
| 52 | 3120 | 0 | [18 21] | [16 32] | 8 | Ev |
| 53 | 3136 | 0 | [19 20] | [16 32] | 8 | Ev |
| 54 | 2976 | 0 | [18 20] | [16 32] | 8 | Ev |
| 55 | 4704 | 0 | [24 24] | [16 32] | 8 | Ev |
| 56 | 5504 | 0 | [26 26] | [16 32] | 8 | Ev |
| 57 | 1896 | 0 | [15 15] | [16 32] | 8 | Ev |
| 58 | 1776 | 0 | [14 15] | [16 32] | 8 | Ev |
| 59 | 1776 | 0 | [15 14] | [16 32] | 8 | Ev |
| 60 | 1664 | 0 | [14 14] | [16 32] | 8 | Ev |
| 61 | 6592 | 0 | [28 29] | [16 32] | 8 | Ev |
| 62 | 6592 | 0 | [29 28] | [16 32] | 8 | Ev |
| 63 | 2984 | 0 | [19 19] | [16 32] | 8 | Ev |
| 64 | 2832 | 0 | [18 19] | [16 32] | 8 | Ev |
| 65 | 3576 | 0 | [15 29] | [16 32] | 8 | Ev |
| 66 | 5200 | 0 | [22 29] | [16 32] | 8 | Ev |
| 67 | 5896 | 0 | [25 29] | [16 32] | 8 | Ev |
| 68 | 6360 | 0 | [27 29] | [16 32] | 8 | Ev |
| 69 | 3576 | 0 | [29 15] | [16 32] | 8 | Ev |
| 70 | 5200 | 1 | [29 22] | [16 32] | 8 | Ev |
| 71 | 4272 | 1 | [29 18] | [16 32] | 8 | Ev |
| 72 | 3808 | 0 | [29 16] | [16 32] | 8 | Ev |
| 73 | 4040 | 1 | [29 17] | [16 32] | 8 | Ev |
| 74 | 112 | 0 | [2 2] | [16 8] | 16 | Ev |
| 75 | 624 | 0 | [6 6] | [16 8] | 16 | Ev |
| 76 | 528 | 0 | [5 6] | [16 8] | 16 | Ev |
| 77 | 528 | 0 | [6 5] | [16 8] | 16 | Ev |
| 78 | 448 | 0 | [5 5] | [16 8] | 16 | Ev |
| 79 | 1984 | 0 | [11 11] | [16 8] | 16 | Ev |
| 80 | 1808 | 0 | [10 11] | [16 8] | 16 | Ev |
| 81 | 1808 | 0 | [11 10] | [16 8] | 16 | Ev |
| 82 | 1648 | 0 | [10 10] | [16 8] | 16 | Ev |
| 83 | 2752 | 0 | [13 13] | [16 8] | 16 | Ev |
| 84 | 3184 | 0 | [14 14] | [16 8] | 16 | Ev |
| 85 | 3648 | 0 | [15 15] | [16 8] | 16 | Ev |
| 86 | 1632 | 0 | [9 11] | [16 8] | 16 | Ev |
| 87 | 1488 | 0 | [9 10] | [16 8] | 16 | Ev |
| 88 | 3408 | 0 | [14 15] | [16 8] | 16 | Ev |
| 89 | 3408 | 0 | [15 14] | [16 8] | 16 | Ev |
| 90 | 160 | 0 | [2 2] | [16 32] | 16 | Ev |
| 91 | 672 | 0 | [6 6] | [16 32] | 16 | Ev |
| 92 | 576 | 0 | [5 6] | [16 32] | 16 | Ev |
| 93 | 576 | 0 | [6 5] | [16 32] | 16 | Ev |
| 94 | 496 | 0 | [5 5] | [16 32] | 16 | Ev |
| 95 | 2032 | 0 | [11 11] | [16 32] | 16 | Ev |
| 96 | 1856 | 0 | [10 11] | [16 32] | 16 | Ev |
| 97 | 1856 | 0 | [11 10] | [16 32] | 16 | Ev |
| 98 | 1696 | 0 | [10 10] | [16 32] | 16 | Ev |
| 99 | 2800 | 0 | [13 13] | [16 32] | 16 | Ev |
| 100 | 3232 | 0 | [14 14] | [16 32] | 16 | Ev |
| 101 | 3696 | 0 | [15 15] | [16 32] | 16 | Ev |
| 102 | 1680 | 0 | [9 11] | [16 32] | 16 | Ev |
| 103 | 1536 | 0 | [9 10] | [16 32] | 16 | Ev |

| | | | | | | |
|-----|------|---|---------|---------|----|---------|
| 104 | 3456 | 0 | [14 15] | [16 32] | 16 | EvenPow |
| 105 | 3456 | 0 | [15 14] | [16 32] | 16 | EvenPow |
| 106 | 80 | 0 | [2 2] | [16 8] | 8 | EvenPow |
| 107 | 176 | 0 | [4 4] | [16 8] | 8 | EvenPow |
| 108 | 560 | 0 | [8 8] | [16 8] | 8 | EvenPow |
| 109 | 1848 | 0 | [15 15] | [16 8] | 8 | EvenPow |
| 110 | 128 | 0 | [2 2] | [16 32] | 8 | EvenPow |
| 111 | 224 | 0 | [4 4] | [16 32] | 8 | EvenPow |
| 112 | 608 | 0 | [8 8] | [16 32] | 8 | EvenPow |
| 113 | 1896 | 0 | [15 15] | [16 32] | 8 | EvenPow |
| 114 | 112 | 0 | [2 2] | [16 8] | 16 | EvenPow |
| 115 | 304 | 0 | [4 4] | [16 8] | 16 | EvenPow |
| 116 | 1072 | 0 | [8 8] | [16 8] | 16 | EvenPow |
| 117 | 3648 | 0 | [15 15] | [16 8] | 16 | EvenPow |
| 118 | 160 | 0 | [2 2] | [16 32] | 16 | EvenPow |
| 119 | 352 | 0 | [4 4] | [16 32] | 16 | EvenPow |
| 120 | 1120 | 0 | [8 8] | [16 32] | 16 | EvenPow |
| 121 | 3696 | 0 | [15 15] | [16 32] | 16 | EvenPow |
| 122 | 176 | 0 | [2 2] | [16 8] | 32 | EvenPow |
| 123 | 560 | 0 | [4 4] | [16 8] | 32 | EvenPow |
| 124 | 2096 | 0 | [8 8] | [16 8] | 32 | EvenPow |
| 125 | 224 | 0 | [2 2] | [16 32] | 32 | EvenPow |
| 126 | 1248 | 0 | [6 6] | [16 32] | 32 | EvenPow |
| 127 | 1056 | 0 | [5 6] | [16 32] | 32 | EvenPow |
| 128 | 1056 | 0 | [6 5] | [16 32] | 32 | EvenPow |
| 129 | 896 | 0 | [5 5] | [16 32] | 32 | EvenPow |
| 130 | 3968 | 0 | [11 11] | [16 32] | 32 | EvenPow |
| 131 | 3616 | 0 | [10 11] | [16 32] | 32 | EvenPow |
| 132 | 3616 | 0 | [11 10] | [16 32] | 32 | EvenPow |
| 133 | 3296 | 0 | [10 10] | [16 32] | 32 | EvenPow |
| 134 | 3264 | 0 | [9 11] | [16 32] | 32 | EvenPow |
| 135 | 2976 | 0 | [9 10] | [16 32] | 32 | EvenPow |
| 136 | 224 | 0 | [2 2] | [16 32] | 32 | EvenPow |
| 137 | 608 | 0 | [4 4] | [16 32] | 32 | EvenPow |
| 138 | 2144 | 0 | [8 8] | [16 32] | 32 | EvenPow |
| 139 | 3624 | 0 | [21 21] | [16 32] | 8 | EvenPow |
| 140 | 3456 | 0 | [20 21] | [16 32] | 8 | EvenPow |
| 141 | 3288 | 0 | [19 21] | [16 32] | 8 | EvenPow |
| 142 | 3120 | 0 | [18 21] | [16 32] | 8 | EvenPow |
| 143 | 2984 | 0 | [19 19] | [16 32] | 8 | EvenPow |
| 144 | 2832 | 0 | [18 19] | [16 32] | 8 | EvenPow |
| 145 | 128 | 0 | [2 2] | [16 32] | 8 | EvenPow |
| 146 | 352 | 0 | [4 8] | [16 32] | 8 | EvenPow |
| 147 | 224 | 0 | [4 4] | [16 32] | 8 | EvenPow |
| 148 | 608 | 0 | [8 8] | [16 32] | 8 | EvenPow |
| 149 | 1896 | 0 | [15 15] | [16 32] | 8 | EvenPow |
| 150 | 192 | 0 | [2 2] | [16 16] | 32 | EvenPow |
| 151 | 2112 | 0 | [8 8] | [16 16] | 32 | EvenPow |
| 152 | 1088 | 0 | [4 8] | [16 16] | 32 | EvenPow |
| 153 | 1088 | 0 | [8 4] | [16 16] | 32 | EvenPow |
| 154 | 576 | 0 | [4 4] | [16 16] | 32 | EvenPow |
| 155 | 224 | 0 | [2 2] | [16 32] | 32 | EvenPow |
| 156 | 2144 | 0 | [8 8] | [16 32] | 32 | EvenPow |
| 157 | 1120 | 0 | [4 8] | [16 32] | 32 | EvenPow |
| 158 | 1120 | 0 | [8 4] | [16 32] | 32 | EvenPow |
| 159 | 608 | 0 | [4 4] | [16 32] | 32 | EvenPow |

Best Solution

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|---------------|----------|------------|-----------------|--------------|----------------|
| 73 | 4040 | 1 | [29 17] | [16 32] | 8 | Ev |

S =

1x1 FunctionApproximation.LUTSolution with properties:

ID: 73
Feasible: "true"

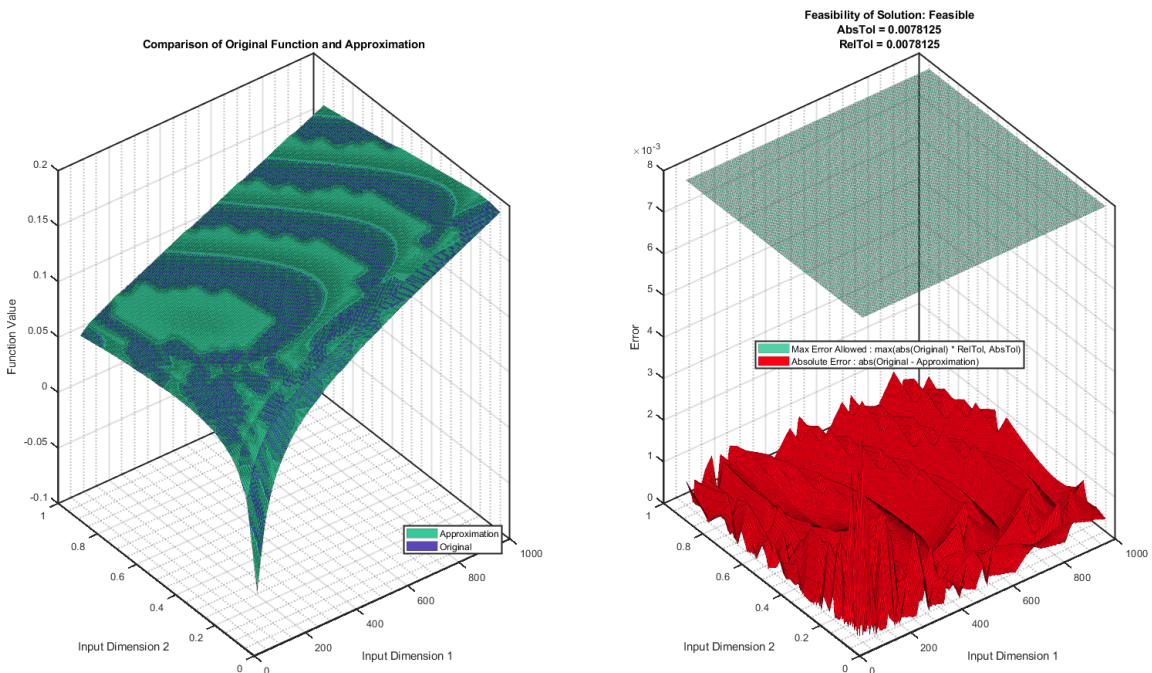
Compare the numerical behavior of the original lookup table, with the optimized lookup table.

compare(S)

ans =

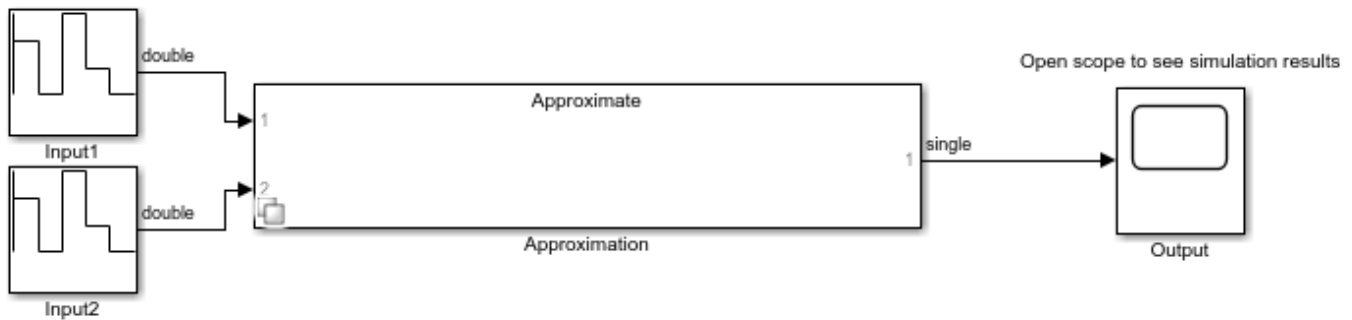
1x2 struct array with fields:

Breakpoints
Original
Approximate



Generate the new Lookup Table block using the approximate method.

S.approximate



```
 %#ok<*NOPTS>
```

Visualize Pareto Front for Memory Optimization Versus Absolute Tolerance

When you want to optimize for both memory and absolute tolerance, it is helpful to visualize the tradeoffs between the two. This example creates a lookup table approximation of the function $1 - \exp(-x)$ with varying levels of absolute tolerance and creates a plot of each solution found. In the final plot you can view the tradeoffs between memory efficiency and numeric fidelity.

```
nTol = 32; % Initialize variables
solutions = cell(1,nTol);
objectiveValues = cell(1,nTol);
constraintValues = cell(1,nTol);
memoryUnits = 'bytes';

% Options for absolute tolerance
absTol = 2.^linspace(-12,-4,nTol);

% Relative tolerance is set to 0
relTol = 0;

% Initialize options
options = FunctionApproximation.Options( ...
    'RelTol', relTol, ...
    'BreakpointSpecification', 'EvenSpacing', ...
    'Display', false, ...
    'WordLengths', 16);

% Setup the approximation problem
problem = FunctionApproximation.Problem( ...
    @(x) 1 - exp(-x), ...
    'InputTypes', numerictype(0,16), ...
    'OutputType', numerictype(1,16,14), ...
    'InputLowerBounds', 0, ...
    'InputUpperBounds', 5, ...
    'Options', options);

% Execute to find solutions with different tolerances
for iTol = 1:nTol
    problem.Options.AbsTol = absTol(iTol);
    solution = solve(problem);
```



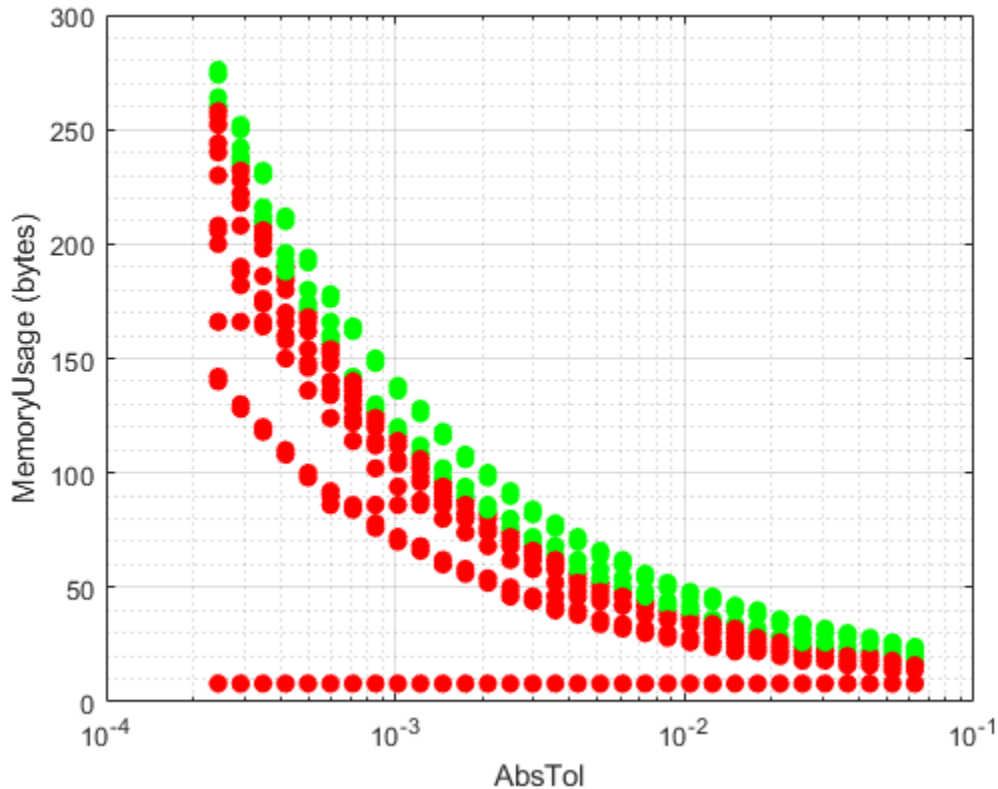
```
objectiveValues{iTol} = arrayfun(@(x) x.totalMemoryUsage(memoryUnits), solution.AllSolutions);
constraintValues{iTol} = arrayfun(@(x) x.Feasible, solution.AllSolutions);
solutions{iTol} = solution;
end

% Plot results

h = figure();
hold on;

for iTol = 1:nTol
    for iObjective = 1:numel(objectiveValues{iTol})
        if constraintValues{iTol}(iObjective)
            markerColor = 'g';
        else
            markerColor = 'r';
        end
        plot(absTol(iTol),objectiveValues{iTol}(iObjective), ...
            'Marker', '.', 'LineStyle', 'none', ...
            'MarkerSize', 24, ...
            'MarkerEdgeColor', markerColor)
    end
end

xlabel('AbsTol')
ylabel(['MemoryUsage (' ,memoryUnits, ')'])
h.Children.XScale = 'log';
h.Children.YMinorGrid = 'on';
grid on
box on
hold off;
```



Solutions that are infeasible, meaning they do not meet the required absolute tolerance are marked in red. Solutions that are feasible are marked in green. As the absolute tolerance increases, the approximation finds solutions which use less memory. When the absolute tolerance is lower, indicating higher numerical fidelity, the required memory also increases.

Compare Approximations Using On Curve and Off Curve Table Values

This example compares the lookup table approximations generated for the `tanh` function when the `OnCurveTableValues` property of the `FunctionApproximation.Options` object is set to `true` and `false`. The `OnCurveTableValues` property specifies whether the table values of the optimized lookup table approximation must be equal to the quantized output of the original function being approximated. In some cases, by setting this value to `false`, the generated lookup table approximation can maintain the same error tolerances while reducing the memory used by the lookup table.

Create a Lookup Table Approximation Using On Curve Table Values

Use the `FunctionApproximation.Problem` object to define a function to approximate with a lookup table. By default, the `OnCurveTableValues` property of the associated `Options` object is set to `false`. Set this property to `true` to constrain table values to the quantized output of the function being approximated.

```
P1 = FunctionApproximation.Problem('tanh');
P1.Options.OnCurveTableValues = 1
```

```
P1 =
  1x1 FunctionApproximation.Problem with properties:

    FunctionToApproximate: @(x)tanh(x)
      NumberOfInputs: 1
        InputTypes: "numerictype(1,16,12)"
      InputLowerBounds: -8
      InputUpperBounds: 8
      OutputType: "numerictype(1,16,15)"
      Options: [1x1 FunctionApproximation.Options]
```

Generate the lookup table approximation.

```
S1 = solve(P1)
```

Searching for fixed-point solutions.

| ID | Memory (bits) | Feasible | Table Size | Breakpoints | WLs | TableData | WL | BreakpointSpec: |
|----|---------------|----------|------------|-------------|-----|-----------|----|-----------------|
| 0 | 64 | 0 | 2 | | 16 | | 16 | Ev |
| 1 | 1248 | 1 | 76 | | 16 | | 16 | Ev |
| 2 | 1232 | 1 | 75 | | 16 | | 16 | Ev |
| 3 | 944 | 1 | 57 | | 16 | | 16 | Ev |
| 4 | 928 | 0 | 56 | | 16 | | 16 | Ev |
| 5 | 656 | 0 | 39 | | 16 | | 16 | Ev |
| 6 | 640 | 0 | 38 | | 16 | | 16 | Ev |
| 7 | 784 | 0 | 47 | | 16 | | 16 | Ev |
| 8 | 864 | 0 | 52 | | 16 | | 16 | Ev |
| 9 | 896 | 0 | 54 | | 16 | | 16 | Ev |
| 10 | 912 | 0 | 55 | | 16 | | 16 | Ev |
| 11 | 496 | 0 | 29 | | 16 | | 16 | Ev |
| 12 | 720 | 0 | 43 | | 16 | | 16 | Ev |
| 13 | 832 | 0 | 50 | | 16 | | 16 | Ev |
| 14 | 880 | 0 | 53 | | 16 | | 16 | Ev |
| 15 | 448 | 1 | 14 | | 16 | | 16 | Expli |

Best Solution

| ID | Memory (bits) | Feasible | Table Size | Breakpoints | WLs | TableData | WL | BreakpointSpec: |
|----|---------------|----------|------------|-------------|-----|-----------|----|-----------------|
| 15 | 448 | 1 | 14 | | 16 | | 16 | Expli |

```
S1 =
  1x1 FunctionApproximation.LUTSolution with properties:
```

```
    ID: 15
    Feasible: "true"
```

Create a Lookup Table Approximation Using Any Table Values

Create another `FunctionApproximation.Problem` object. Set the `OnCurveTableValues` property of this object to false to allow the optimization to optimize the table values as well as the breakpoints.

```
P2 = FunctionApproximation.Problem('tanh');
P2.Options.OnCurveTableValues = 0
```

```
P2 =
  1x1 FunctionApproximation.Problem with properties:
```

```
FunctionToApproximate: @(x)tanh(x)
  NumberOfInputs: 1
    InputTypes: "numerictype(1,16,12)"
  InputLowerBounds: -8
  InputUpperBounds: 8
    OutputType: "numerictype(1,16,15)"
    Options: [1x1 FunctionApproximation.Options]
```

Generate the lookup table approximation.

```
S2 = solve(P2)
```

Searching for fixed-point solutions.

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|---------------|----------|------------|-----------------|--------------|----------------|
| 0 | 64 | 0 | 2 | 16 | 16 | Even |
| 1 | 1248 | 1 | 76 | 16 | 16 | Even |
| 2 | 1232 | 1 | 75 | 16 | 16 | Even |
| 3 | 944 | 1 | 57 | 16 | 16 | Even |
| 4 | 928 | 1 | 56 | 16 | 16 | Even |
| 5 | 656 | 0 | 39 | 16 | 16 | Even |
| 6 | 640 | 0 | 38 | 16 | 16 | Even |
| 7 | 784 | 1 | 47 | 16 | 16 | Even |
| 8 | 704 | 1 | 42 | 16 | 16 | Even |
| 9 | 672 | 1 | 40 | 16 | 16 | Even |
| 10 | 368 | 0 | 21 | 16 | 16 | Even |
| 11 | 512 | 0 | 30 | 16 | 16 | Even |
| 12 | 592 | 0 | 35 | 16 | 16 | Even |
| 13 | 624 | 0 | 37 | 16 | 16 | Even |
| 14 | 384 | 1 | 12 | 16 | 16 | Explicit |
| 15 | 384 | 0 | 12 | 16 | 16 | Explicit |
| 16 | 384 | 1 | 12 | 16 | 16 | Explicit |

Best Solution

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|---------------|----------|------------|-----------------|--------------|----------------|
| 14 | 384 | 1 | 12 | 16 | 16 | Explicit |

S2 =

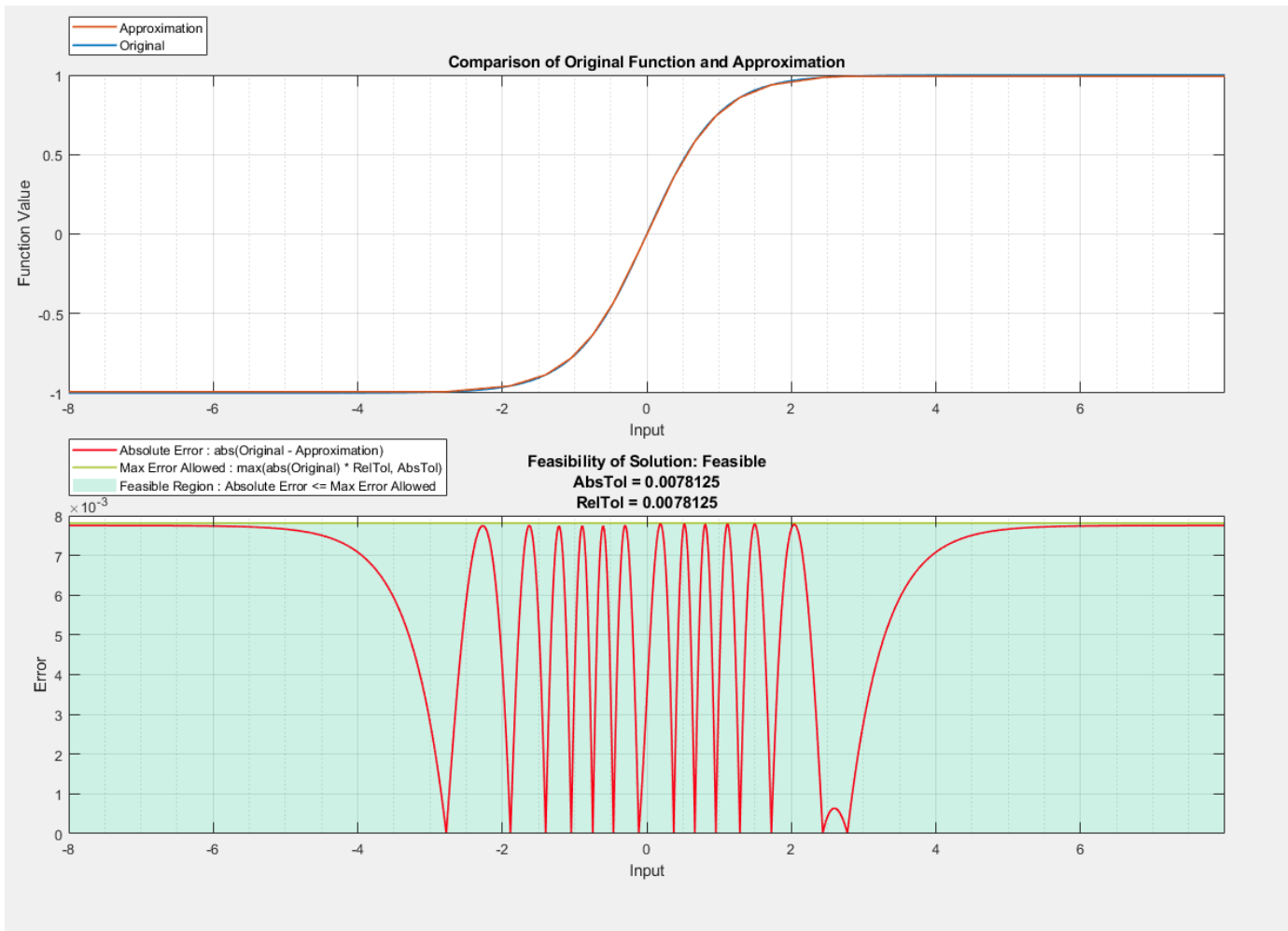
1x1 FunctionApproximation.LUTSolution with properties:

```
    ID: 14
  Feasible: "true"
```

View Results

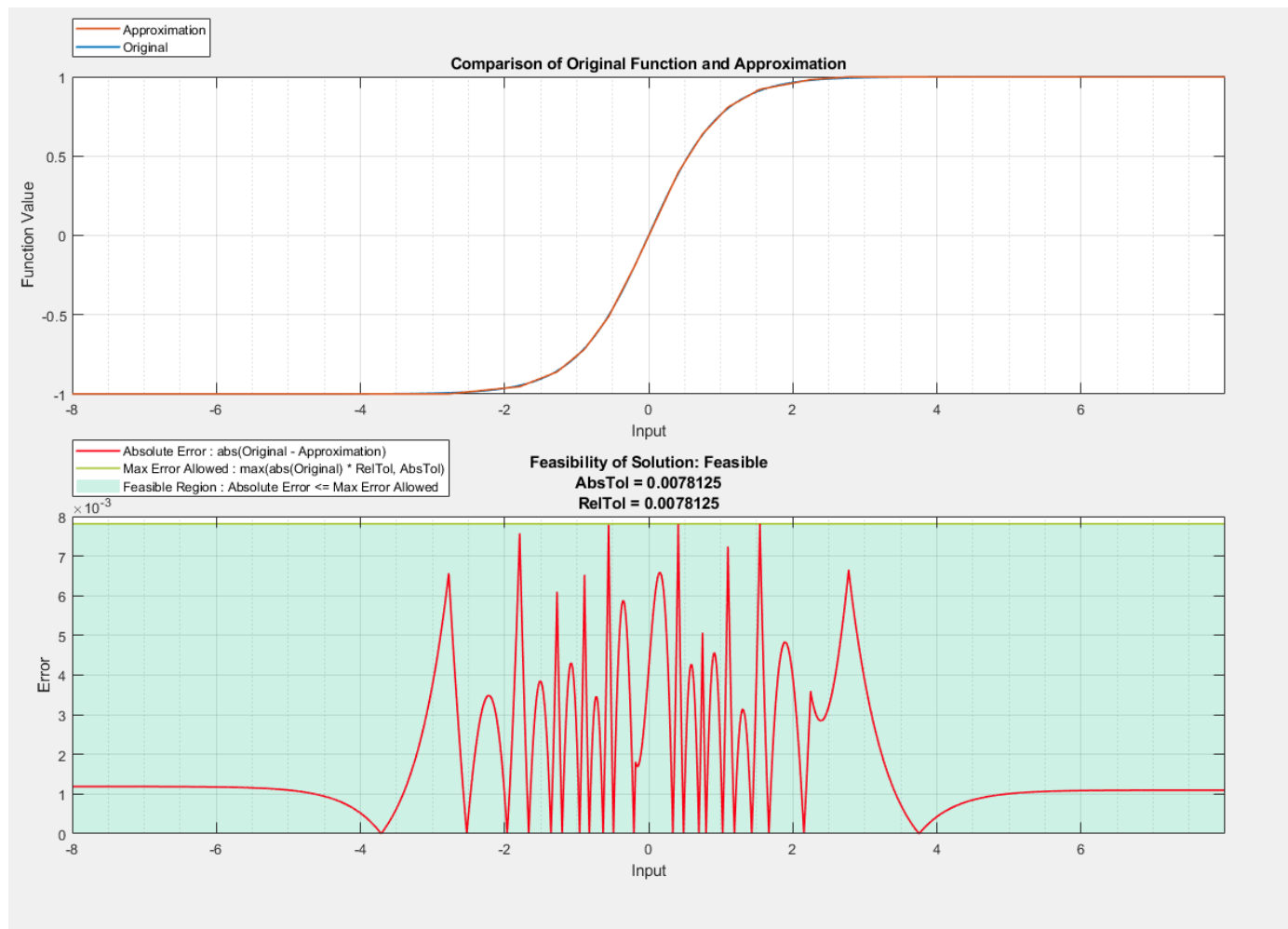
Compare the best solutions for each lookup table approximation.

```
compare(S1)
```



```
ans = struct with fields:
  Breakpoints: [65536x1 double]
  Original: [65536x1 double]
  Approximate: [65536x1 double]
```

```
compare(S2)
```



```
ans = struct with fields:
  Breakpoints: [65536x1 double]
  Original: [65536x1 double]
  Approximate: [65536x1 double]
```

The maximum error between the original function and the two lookup table approximations are approximately equal, however the memory used by the lookup table that was not constrained to using only on curve table values is significantly lower.

```
percent_reduction = S2.totalMemoryUsage/S1.totalMemoryUsage
```

```
percent_reduction = 0.8571
```

See Also

Apps

Lookup Table Optimizer

Classes

FunctionApproximation.Problem | FunctionApproximation.Options |
FunctionApproximation.LUTSolution |
FunctionApproximation.LUTMemoryUsageCalculator

More About

- “Optimize Lookup Tables for Memory-Efficiency” on page 41-15

Generate an Optimized Lookup Table as a MATLAB Function

This example shows how to approximate $y = 1/(1+\exp(-x))$ as a MATLAB function lookup table using the Lookup Table Optimizer.

- 1 To open the Lookup Table Optimizer, on the Simulink **Apps** tab, in the **Code Generation** gallery, click **Lookup Table Optimizer**.
- 2 In the **Objective** pane of the app, select the source as **MATLAB Function Handle**. Click **Next**.
- 3 In the **Setup** pane, provide the function handle $@(x)(1/(1+\exp(-x)))$.

The attributes populate in the table below. You can manually edit the fields to specify ranges and data types other than those populated. For this example, set **Minimum** to 0 and **Maximum** to 0.25.

Click **Next**.

- 4 In the **Create** pane, specify the **Output Error Tolerance** that is acceptable for your design.

To specify additional properties for the optimized lookup table, click **LUT Specification**. Change the **Solution Type** to **MATLAB**.

After you are satisfied with the constraints and additional options, click **Optimize**. When the optimization is complete, the Lookup Table Optimizer reports the memory of the optimized lookup table. You can edit the constraints and run the optimization again to achieve further memory reduction.

Constraints

| | | |
|-------------------------------|--|--|
| Output Error Tolerance | | Allowed Word Lengths (Vector) |
| Absolute | <input type="text" value="0.0078125"/> | <input type="text" value="[8 16 32]"/> |
| Relative | <input type="text" value="0.0078125"/> | |

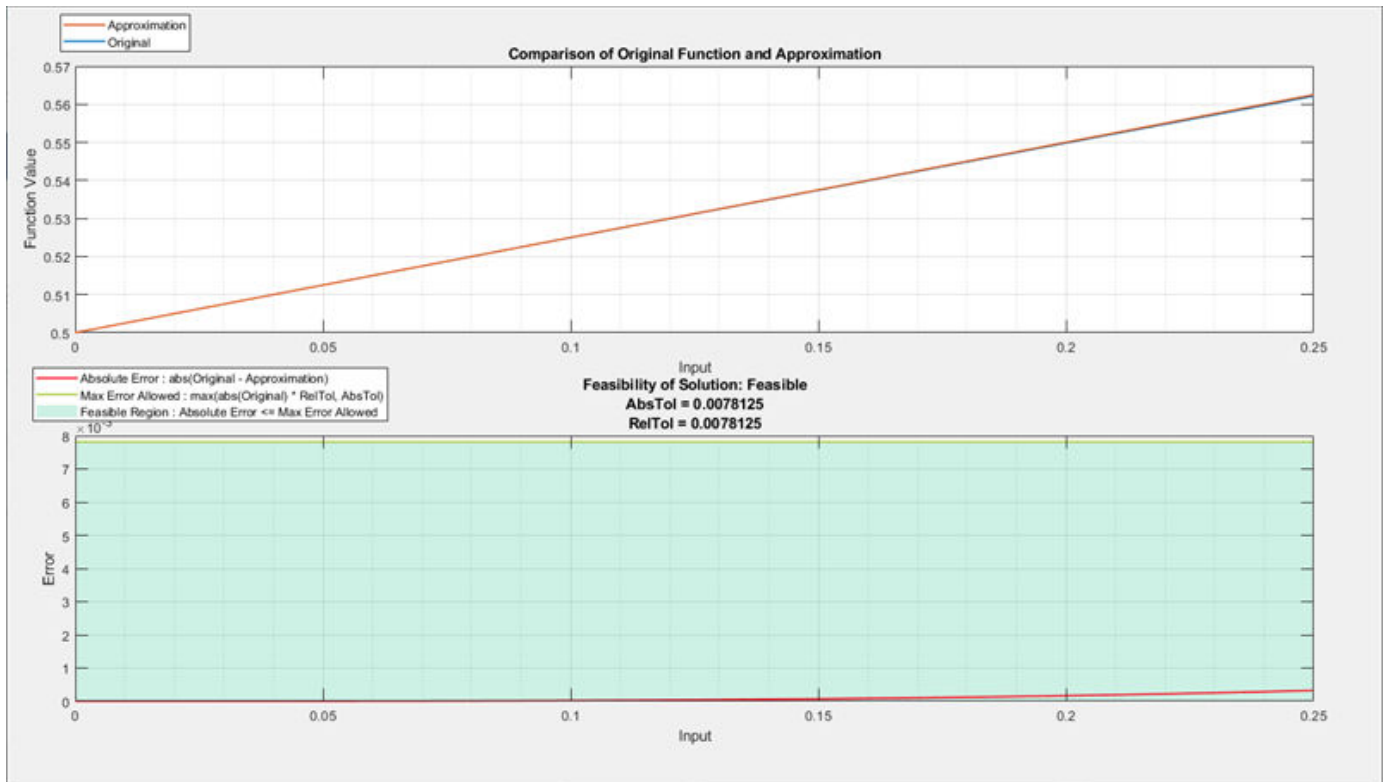
Optimization

| LUT Attributes | New Memory* | New Data Type |
|----------------|-------------|--------------------|
| Table Data | 2 | numerictype(0,8,8) |
| Breakpoint 1 | 2 | numerictype(0,8,9) |
| Total | 4 | |

* In bytes

Click **Next** to view the **Results** pane.

- 5 Click **Show Comparison Plot** to view a plot of the original function output compared to the output of the new optimized lookup table.



- 6 Click **Show Optimized LUT** to view the lookup table function launched in the MATLAB Command Window.

See Also

Lookup Table Optimizer

Related Examples

- "Optimize Lookup Tables for Memory-Efficiency" on page 41-15
- "Generate an Optimized Lookup Table as a MATLAB Function Programmatically" on page 41-38

Generate an Optimized Lookup Table as a MATLAB Function Programmatically

This example shows how to generate an optimized lookup table as a MATLAB® function to approximate hyperbolic tangent. The MATLAB function lookup table approximation can then be used to replace the hyperbolic tangent function and generate C code.

Use the `FunctionApproximation.Options` object to specify a MATLAB function as the solution type. Use the default values for accuracy and word length constraints.

```
options = FunctionApproximation.Options();
options.ApproximateSolutionType = 'MATLAB';
```

Specify the function to approximate and the input ranges and data types in the `FunctionApproximation.Problem` object.

```
functionToApproximate = 'tanh';
```

```
problem = FunctionApproximation.Problem(functionToApproximate, 'Options', options);
problem.InputLowerBounds = 0;
problem.InputUpperBounds = 0.25;
```

Use the `solve` method to solve the optimization problem and create a lookup table solution.

```
solution = solve(problem)
```

```
Searching for fixed-point solutions.
```

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|---------------|----------|------------|-----------------|--------------|----------------|
| 0 | 64 | 1 | 2 | 16 | 16 | EvenPow |
| 1 | 64 | 1 | 2 | 16 | 16 | EvenPow |

```
Best Solution
```

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|---------------|----------|------------|-----------------|--------------|----------------|
| 1 | 64 | 1 | 2 | 16 | 16 | EvenPow |

```
solution =
```

```
  1x1 FunctionApproximation.LUTSolution with properties:
```

```
      ID: 1
  Feasible: "true"
```

To obtain the generated lookup table as a MATLAB function, use the `approximate` method. Use optional name-value arguments to specify the name and path for the lookup table function.

```
filename = 'tanhApproximate';
filepath = cd;
approximate(solution, 'Name', filename, 'Path', filepath);
```

```

1 function output = tanhApproximate(inputValues1)
2 %#codegen
3 %This code is an approximation for @(x)tanh(x)
4 %
5 %Inputs for approximation:
6 % InputLowerBounds: 0
7 % InputTypes: numerictype(1,16,12)
8 % InputUpperBounds: 0.25
9 % OutputType: numerictype(1,16,15)
10 % Options (non-default)
11 % AbsTol: 0.0078125
12 % ApproximateSolutionType: MATLAB
13 % BreakpointSpecification: EvenPow2Spacing
14 % RelTol: 0.0078125
15 %
16 %Steps to evaluate function:
17 %Check input type
18 %Create constants for table value and breakpoint values
19 %Initialize variables and set their type
20 %Perform pre-lookup to get the index and fraction corresponding to input
21 %Calculate the output value
22
23
24 inputValues1 = cast(inputValues1,'like',fi('numerictype',numerictype(1,16,12),'Value','[]'));
25
26 breakpointValues1 = coder.const(fi([0 0.25],numerictype(1,16,12)));
27
28 bpSpaceExponent1 = 2;
29
30 tableValues = coder.const(fi('numerictype',numerictype(1,16,15),'Value','[0 0.244903564453125]'));
31
32 idxType = coder.const(uint32([]));
33 fracType = coder.const(fi([],numerictype(0,32,32)));
34
35 f = fimath('RoundingMethod','Floor',...
36     'OverflowAction','Saturate',...
37     'ProductMode','SpecifyPrecision',...
38     'ProductWordLength',16,...
39     'ProductFractionLength',15,...
40     'SumMode','SpecifyPrecision',...
41     'SumWordLength',16,...
42     'SumFractionLength',15,...
43     'CastBeforeSum',true);
44 output = zeros(size(inputValues1),'like',fi([],numerictype(1,16,15),f));
45

```

If you have MATLAB Coder™ installed, you can use the `codegen` command to generate C code from the approximate lookup table function.

```
inputArgs = linspace(1,10,10);
codegen tanhApproximate.m -args {inputArgs}
```

Code generation successful.

See Also

FunctionApproximation.Problem | FunctionApproximation.Options | solve | approximate

Related Examples

- “Generate an Optimized Lookup Table as a MATLAB Function” on page 41-36

Convert Neural Network Algorithms to Fixed-Point Using fxpopt and Generate HDL Code

This example shows how to convert a neural network regression model in Simulink® to fixed point using the fxpopt function and Lookup Table Optimizer.

Overview

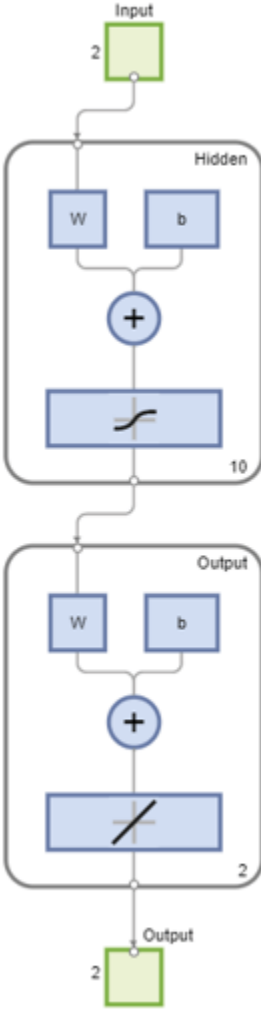
Fixed-Point Designer™ provides workflows via the Fixed Point Tool that can convert a design from floating-point data types to fixed-point data types. The fxpopt function optimizes data types in a model based on specified system behavioral constraints. The Lookup Table Optimizer generates memory-efficient lookup table replacements for unbounded functions such as exp and log2. Using these tools, this example showcases how to convert a trained floating-point neural network regression model to use embedded-efficient fixed-point data types.

Data and Neural Network Training

The engine_dataset contains data representing the relationship between the fuel rate and speed of the engine, and its torque and gas emissions.


Use the Neural Net Fitting app (nftool) from Deep Learning Toolbox™ to train a neural network to estimate torque and gas emissions of an engine given the fuel rate and speed. Use the following commands to train the neural network.

```
load engine_dataset;  
x = engineInputs;  
t = engineTargets;  
net = fitnet(10);  
net = train(net,x,t);  
view(net)
```



Network Diagram

Training Results

Training finished: Met validation criterion 

Training Progress

| Unit | Initial Value | Stopped Value | Target Value |
|-------------------|---------------|---------------|--------------|
| Epoch | 0 | 58 | 1000 |
| Elapsed Time | - | 00:00:08 | - |
| Performance | 7.1e+05 | 1.74e+03 | 0 |
| Gradient | 1.35e+06 | 725 | 1e-07 |
| Mu | 0.001 | 10 | 1e+10 |
| Validation Checks | 0 | 6 | 6 |

Training Algorithms

Data Division: Random dividerand

Training: Levenberg-Marquardt trainlm

Performance: Mean Squared Error mse

Calculations: MEX

Training Plots

Performance Training State

Error Histogram Regression

Fit

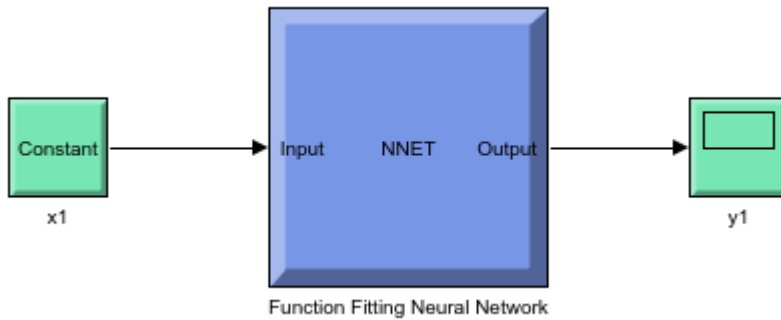
Close the view of the network.

```
nnet.guis.closeAllViews();
```

Model Preparation for Fixed-Point Conversion

Once the network is trained, use the `gensim` function from the Deep Learning Toolbox to generate a Simulink model.

```
[sysName, netName] = gensim(net, 'Name', 'mTrainedNN');
```

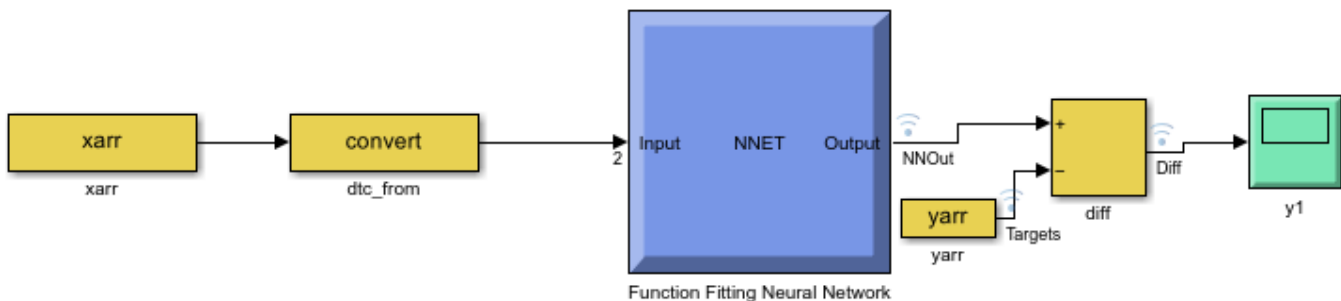


The model generated by the `gensim` function contains the neural network with trained weights and biases. To prepare this generated model for fixed-point conversion, follow the preparation steps in the best practices guidelines.

After applying these principles, the trained neural network is further modified to enable signal logging at the output of the network, add input stimuli and verification blocks.

Open and inspect the model. The model is already configured for HDL compatibility by using the `hdlsetup` function.

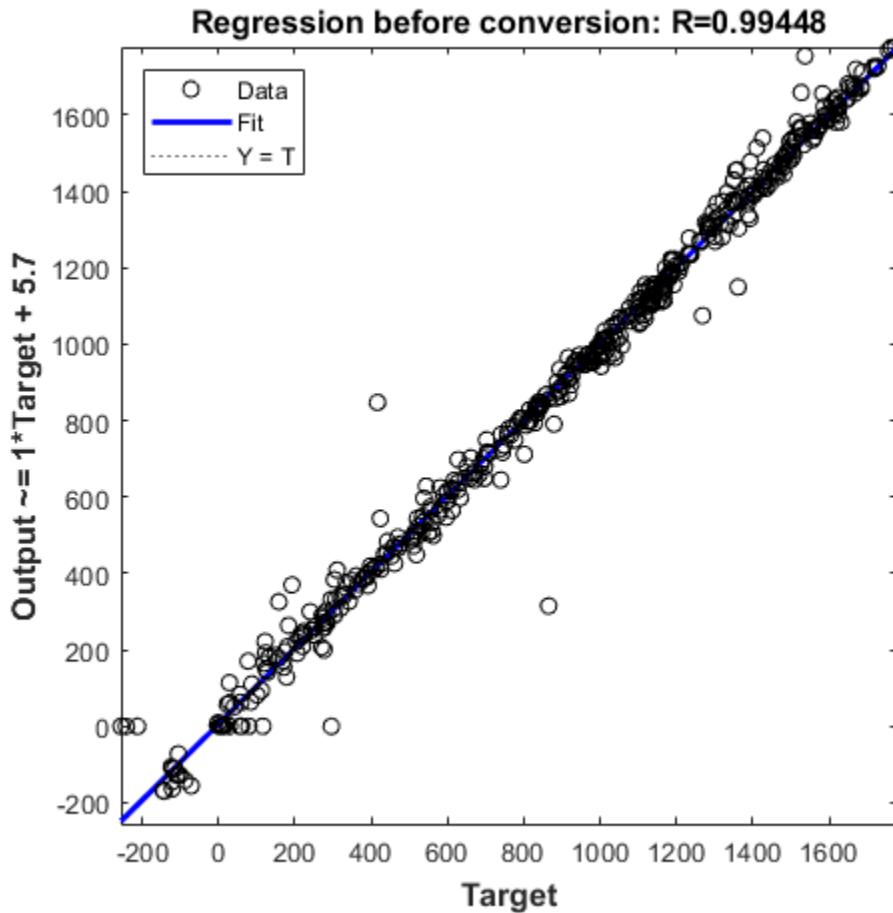
```
model = 'ex_fxpdemo_neuralnet_regression_hdlsetup';
system_under_design = [model '/Function Fitting Neural Network'];
baseline_output = [model '/yarr'];
open_system(model);
```



Simulate the model to observe model performance when using double-precision floating-point data types.

```
loggingInfo = get_param(model, 'DataLoggingOverride');
sim_out = sim(model, 'SaveFormat', 'Dataset');

plotRegression(sim_out, baseline_output, system_under_design, 'Regression before conversion');
```

Define System Behavioral Constraints for Fixed Point Conversion

```
opts = fxpOptimizationOptions();
opts.addTolerance(system_under_design, 1, 'RelTol', 0.05);
opts.addTolerance(system_under_design, 1, 'AbsTol', 50)
opts.AllowableWordLengths = 8:32;
```

Optimize Data Types

Use the `fxpopt` function to optimize the data types in the system under design and explore the solution. The software analyzes the range of objects in `system_under_design` and wordlength and tolerance constraints specified in `opts` to apply heterogeneous data types to the model while minimizing total bit width.

```
solution = fxpopt(model, system_under_design, opts);
best_solution = solution.explore;
```

```
+ Starting data type optimization...
+ Checking for unsupported constructs.
- The paths printed in the Command Window have constructs that do not support fixed-point
'ex_fxpdemo_neuralnet_regression_hdlsetup/Function Fitting Neural Network/Layer 1/tansig/tanl
```

```
+ Preprocessing
+ Modeling the optimization problem
  - Constructing decision variables
+ Running the optimization solver
Exporting logged dataset prior to deleting run...done.
- Evaluating new solution: cost 515, does not meet the behavioral constraints.
- Evaluating new solution: cost 577, does not meet the behavioral constraints.
- Evaluating new solution: cost 639, does not meet the behavioral constraints.
- Evaluating new solution: cost 701, does not meet the behavioral constraints.
- Evaluating new solution: cost 763, does not meet the behavioral constraints.
- Evaluating new solution: cost 825, does not meet the behavioral constraints.
- Evaluating new solution: cost 887, does not meet the behavioral constraints.
- Evaluating new solution: cost 949, meets the behavioral constraints.
- Updated best found solution, cost: 949
- Evaluating new solution: cost 945, meets the behavioral constraints.
- Updated best found solution, cost: 945
- Evaluating new solution: cost 944, meets the behavioral constraints.
- Updated best found solution, cost: 944
- Evaluating new solution: cost 943, meets the behavioral constraints.
- Updated best found solution, cost: 943
- Evaluating new solution: cost 942, meets the behavioral constraints.
- Updated best found solution, cost: 942
- Evaluating new solution: cost 941, meets the behavioral constraints.
- Updated best found solution, cost: 941
- Evaluating new solution: cost 940, meets the behavioral constraints.
- Updated best found solution, cost: 940
- Evaluating new solution: cost 939, meets the behavioral constraints.
- Updated best found solution, cost: 939
- Evaluating new solution: cost 938, meets the behavioral constraints.
- Updated best found solution, cost: 938
- Evaluating new solution: cost 937, meets the behavioral constraints.
- Updated best found solution, cost: 937
- Evaluating new solution: cost 936, meets the behavioral constraints.
- Updated best found solution, cost: 936
- Evaluating new solution: cost 926, meets the behavioral constraints.
- Updated best found solution, cost: 926
- Evaluating new solution: cost 925, meets the behavioral constraints.
- Updated best found solution, cost: 925
- Evaluating new solution: cost 924, meets the behavioral constraints.
- Updated best found solution, cost: 924
- Evaluating new solution: cost 923, meets the behavioral constraints.
- Updated best found solution, cost: 923
- Evaluating new solution: cost 922, meets the behavioral constraints.
- Updated best found solution, cost: 922
- Evaluating new solution: cost 917, meets the behavioral constraints.
- Updated best found solution, cost: 917
- Evaluating new solution: cost 916, meets the behavioral constraints.
- Updated best found solution, cost: 916
- Evaluating new solution: cost 914, meets the behavioral constraints.
- Updated best found solution, cost: 914
- Evaluating new solution: cost 909, meets the behavioral constraints.
- Updated best found solution, cost: 909
- Evaluating new solution: cost 908, meets the behavioral constraints.
- Updated best found solution, cost: 908
- Evaluating new solution: cost 906, meets the behavioral constraints.
- Updated best found solution, cost: 906
- Evaluating new solution: cost 898, meets the behavioral constraints.
- Updated best found solution, cost: 898
```

```

- Evaluating new solution: cost 897, meets the behavioral constraints.
- Updated best found solution, cost: 897
- Evaluating new solution: cost 893, does not meet the behavioral constraints.
- Evaluating new solution: cost 896, meets the behavioral constraints.
- Updated best found solution, cost: 896
- Evaluating new solution: cost 895, meets the behavioral constraints.
- Updated best found solution, cost: 895
- Evaluating new solution: cost 894, meets the behavioral constraints.
- Updated best found solution, cost: 894
- Evaluating new solution: cost 893, meets the behavioral constraints.
- Updated best found solution, cost: 893
- Evaluating new solution: cost 892, meets the behavioral constraints.
- Updated best found solution, cost: 892
- Evaluating new solution: cost 891, meets the behavioral constraints.
- Updated best found solution, cost: 891
- Evaluating new solution: cost 890, meets the behavioral constraints.
- Updated best found solution, cost: 890
- Evaluating new solution: cost 889, meets the behavioral constraints.
- Updated best found solution, cost: 889
- Evaluating new solution: cost 888, meets the behavioral constraints.
- Updated best found solution, cost: 888
- Evaluating new solution: cost 878, meets the behavioral constraints.
- Updated best found solution, cost: 878
- Evaluating new solution: cost 877, meets the behavioral constraints.
- Updated best found solution, cost: 877
- Evaluating new solution: cost 876, meets the behavioral constraints.
- Updated best found solution, cost: 876
- Evaluating new solution: cost 875, meets the behavioral constraints.
- Updated best found solution, cost: 875
- Evaluating new solution: cost 874, meets the behavioral constraints.
- Updated best found solution, cost: 874
- Evaluating new solution: cost 869, meets the behavioral constraints.
- Updated best found solution, cost: 869
- Evaluating new solution: cost 868, does not meet the behavioral constraints.
- Evaluating new solution: cost 867, meets the behavioral constraints.
- Updated best found solution, cost: 867
- Evaluating new solution: cost 862, does not meet the behavioral constraints.
- Evaluating new solution: cost 866, does not meet the behavioral constraints.
- Evaluating new solution: cost 865, does not meet the behavioral constraints.
+ Optimization has finished.
  - Neighborhood search complete.
  - Maximum number of iterations completed.
+ Fixed-point implementation that satisfies the behavioral constraints found. The best found
  - Total cost: 867
  - Maximum absolute difference: 49.714162
  - Use the explore method of the result to explore the implementation.

```

Verify model accuracy after conversion by simulating the model.

```

set_param(model, 'DataLoggingOverride', loggingInfo);
Simulink.sdi.markSignalForStreaming([model '/yarr'], 1, 'on');
Simulink.sdi.markSignalForStreaming([model '/diff'], 1, 'on');
sim_out = sim(model, 'SaveFormat', 'Dataset');

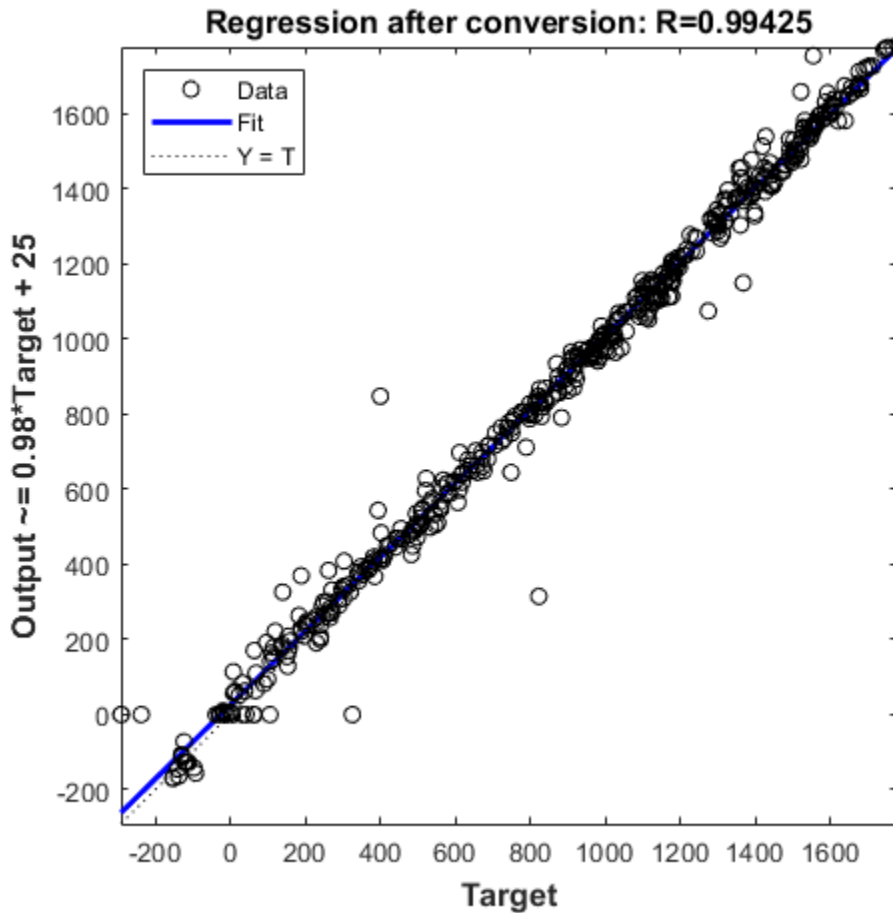
```

Plot the regression accuracy of the fixed-point model.

```

plotRegression(sim_out, baseline_output, system_under_design, 'Regression after conversion');

```



Replace Activation Function with an Optimized Lookup Table

The Tanh Activation function in Layer 1 can be replaced with either a lookup table or a CORDIC implementation for more efficient fixed-point code generation. In this example, we will be using the Lookup Table Optimizer to get a lookup table as a replacement for `tanh`. We will be using `EvenPow2Spacing` for faster execution speed.

```
block_path = [system_under_design '/Layer 1/tansig'];
p = FunctionApproximation.Problem(block_path);
p.Options.WordLengths = 8:32;
p.Options.BreakpointSpecification = 'EvenPow2Spacing';
solution = p.solve;
solution.replaceWithApproximate;
```

Searching for fixed-point solutions.

| ID | Memory (bits) | Feasible | Table Size | Breakpoints | WLs | TableData WL | BreakpointSpec |
|----|---------------|----------|------------|-------------|-----|--------------|----------------|
| 0 | 44 | 0 | 2 | | 14 | 8 | EvenPow |
| 1 | 4124 | 1 | 512 | | 14 | 8 | EvenPow |
| 2 | 4114 | 1 | 512 | | 9 | 8 | EvenPow |
| 3 | 2076 | 0 | 256 | | 14 | 8 | EvenPow |
| 4 | 2064 | 0 | 256 | | 8 | 8 | EvenPow |

| | | | | | | |
|----|------|---|-----|----|----|---------|
| 5 | 46 | 0 | 2 | 14 | 9 | EvenPow |
| 6 | 2332 | 0 | 256 | 14 | 9 | EvenPow |
| 7 | 2320 | 0 | 256 | 8 | 9 | EvenPow |
| 8 | 48 | 0 | 2 | 14 | 10 | EvenPow |
| 9 | 2588 | 0 | 256 | 14 | 10 | EvenPow |
| 10 | 2576 | 0 | 256 | 8 | 10 | EvenPow |
| 11 | 50 | 0 | 2 | 14 | 11 | EvenPow |
| 12 | 2844 | 0 | 256 | 14 | 11 | EvenPow |
| 13 | 2832 | 0 | 256 | 8 | 11 | EvenPow |
| 14 | 52 | 0 | 2 | 14 | 12 | EvenPow |
| 15 | 3100 | 0 | 256 | 14 | 12 | EvenPow |
| 16 | 3088 | 0 | 256 | 8 | 12 | EvenPow |
| 17 | 54 | 0 | 2 | 14 | 13 | EvenPow |
| 18 | 3356 | 0 | 256 | 14 | 13 | EvenPow |
| 19 | 3344 | 0 | 256 | 8 | 13 | EvenPow |

Best Solution

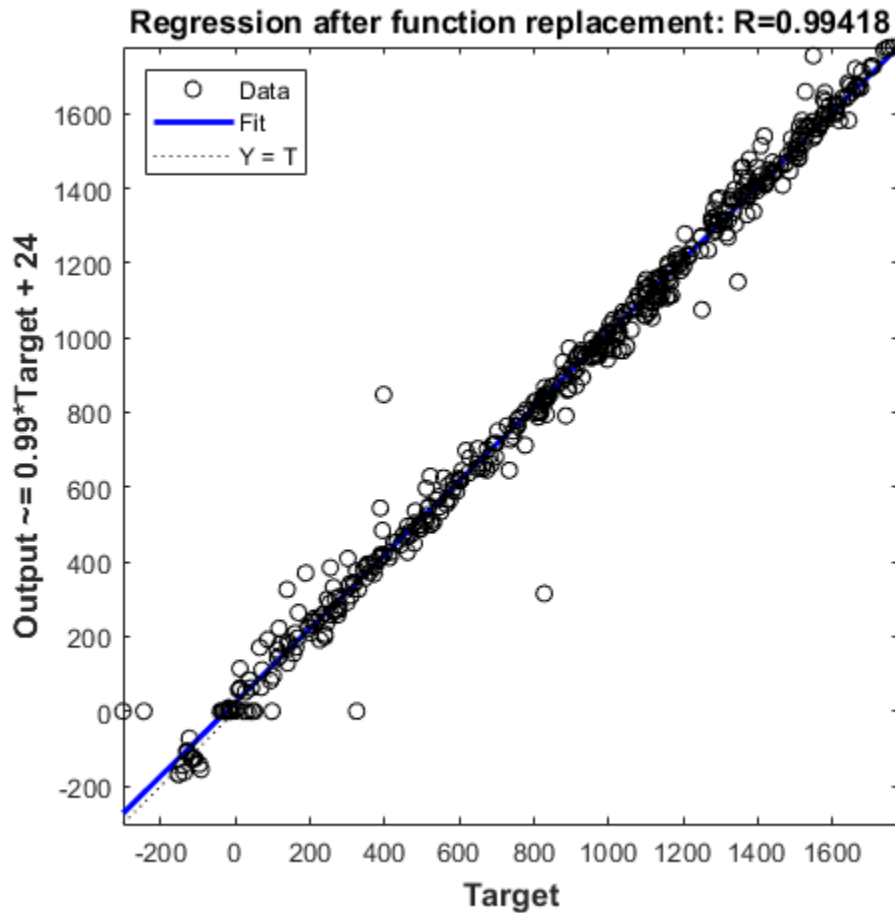
| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec: |
|----|---------------|----------|------------|-----------------|--------------|-----------------|
| 2 | 4114 | 1 | 512 | 9 | 8 | EvenPow |

Verify model accuracy after function replacement

```
sim_out = sim(model, 'SaveFormat', 'Dataset');
```

Plot regression accuracy after function replacement.

```
plotRegression(sim_out, baseline_output, system_under_design, 'Regression after function replacem
```



Generate HDL Code and Test Bench

Generating HDL code requires an HDL Coder™ license.

Choose the model for which to generate HDL code and a test bench.

```
systemname = 'ex_fxpdemo_neuralnet_regression/Function Fitting Neural Network';
```

Use a temporary directory for the generated files.

```
workingdir = tempname;
```

You can run the following command to check for HDL code generation compatibility.

```
checkhdl(systemname, 'TargetDirectory', workingdir);
```

Run the following command to generate HDL code.

```
makehdl(systemname, 'TargetDirectory', workingdir);
```

Run the following command to generate the test bench.

```
makehdltb(systemname, 'TargetDirectory', workingdir);
```

See Also

[fxpopt](#) | [FunctionApproximation.Options](#) | “Best Practices for Fixed-Point Conversion Workflow” on page 42-5

Convert Neural Network Algorithms to Fixed Point and Generate C Code

This example shows how to convert a neural network regression model in Simulink to fixed point using the Fixed-Point Tool and Lookup Table Optimizer and generate C code using Simulink Coder.

Overview

Fixed-Point Designer provides workflows via the Fixed Point Tool that can convert a design from floating-point data types to fixed-point data types. The Lookup Table Optimizer generates memory-efficient lookup table replacements for unbounded functions such as `exp` and `log2`. Using these tools, this example showcases how to convert a trained floating-point neural network regression model to use embedded-efficient fixed-point data types.

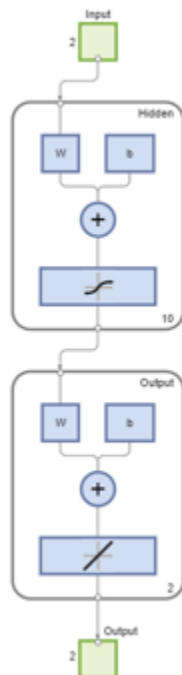
Data and Neural Network Training

The `engine_dataset` contains data representing the relationship between the fuel rate, speed of the engine, and its torque and gas emissions.

```
load engine_dataset;
```


Use the function fitting tool `nftool` from Deep Learning Toolbox™ to train a neural network to estimate torque and gas emissions of an engines given the fuel rate and speed. Use the following commands to train the neural network.

```
x = engineInputs;
t = engineTargets;
net = fitnet(10);
net = train(net,x,t);
view(net)
```



Network Diagram

Training Results

Training finished: Met validation criterion 

Training Progress

| Unit | Initial Value | Stopped Value | Target Value |
|-------------------|---------------|---------------|--------------|
| Epoch | 0 | 58 | 1000 |
| Elapsed Time | - | 00:00:15 | - |
| Performance | 7.1e+05 | 1.74e+03 | 0 |
| Gradient | 1.35e+06 | 725 | 1e-07 |
| Mu | 0.001 | 10 | 1e+10 |
| Validation Checks | 0 | 6 | 6 |

Training Algorithms

Data Division: Random dividerand

Training: Levenberg-Marquardt trainlm

Performance: Mean Squared Error mse

Calculations: MEX

Training Plots

Performance Training State

Error Histogram Regression

Fit

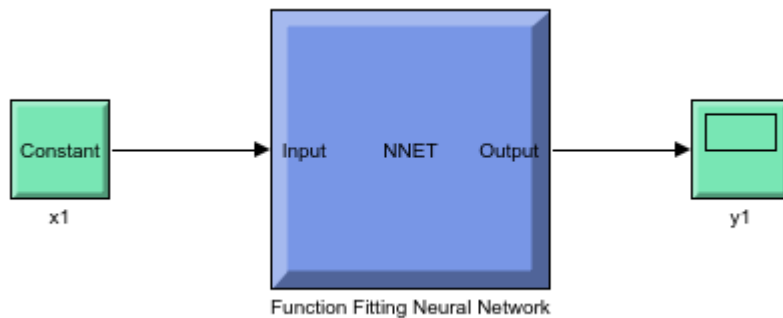
Close the view of the network.

```
nnet.guis.closeAllViews();
```

Model Preparation for Fixed-Point Conversion

Once the network is trained, use the `gensim` function from the Deep Learning Toolbox™ to generate a simulink model.

```
sys_name = gensim(net, 'Name', 'mTrainedNN');
```

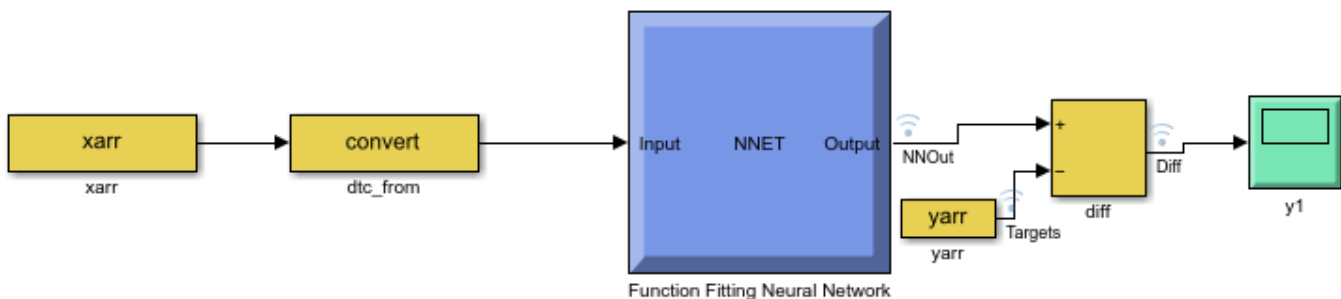


The model generated by the `gensim` function contains the neural network with trained weights and biases. To prepare this generated model for fixed-point conversion, follow the preparation steps in the best practices guidelines. <https://www.mathworks.com/help/fixedpoint/ug/best-practices-for-using-the-fixed-point-tool-to-propose-data-types-for-your-simulink-model.html>

After applying these principles, the trained neural network is further modified to enable signal logging at the output of the network, add input stimuli and verification blocks.

Open and inspect the model.

```
model = 'ex_fxpdemo_neuralnet_regression';
system_under_design = [model '/Function Fitting Neural Network'];
baseline_output = [model '/yarr'];
open_system(model);
```



To open the Fixed-Point Tool, right click on the Function Fitting Neural Network subsystem and select **Fixed-Point Tool**. Alternatively, use the command-line interface of the Fixed-Point Tool. Fixed Point Tool and the command-line interface provide workflow steps for model preparation for fixed point conversion, range and overflow instrumentation of objects via simulation and range analysis, homogeneous wordlength exploration for fixed point data typing and additional overflow diagnostics.

```
converter = DataTypeWorkflow.Converter(system_under_design);
```

Run Simulation to Collect Ranges

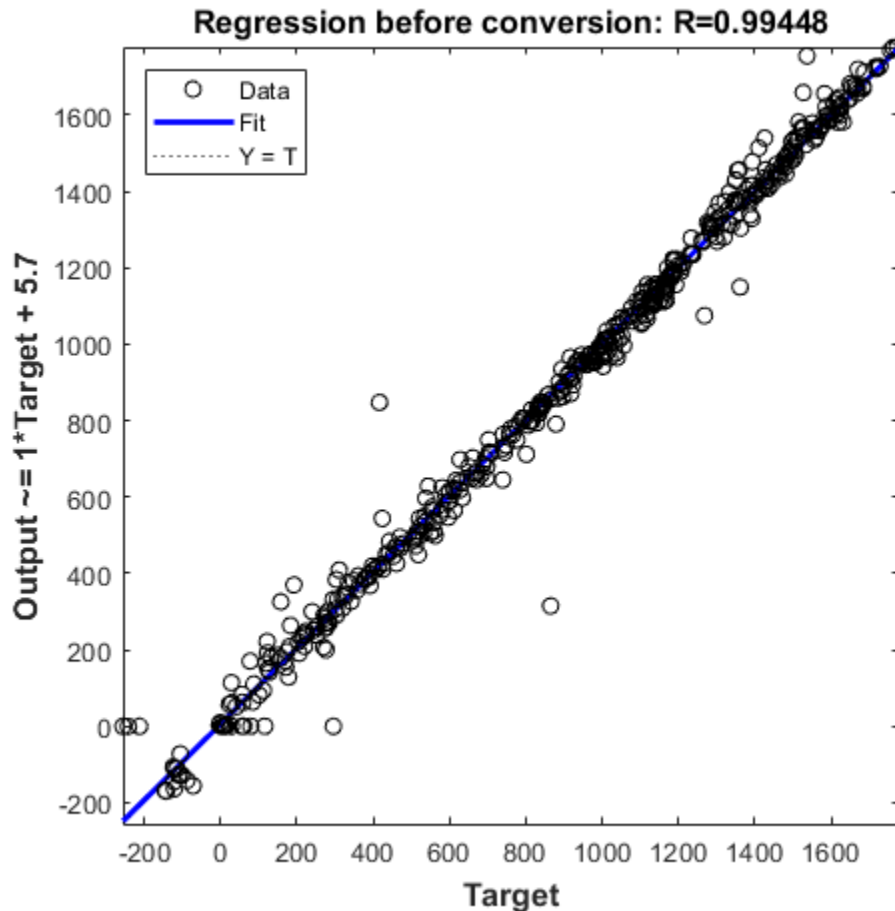
Simulate the model with instrumentation to collect ranges. This is achieved by clicking the **Collect Ranges** button in the tool or the following commands.

```
converter.applySettingsFromShortcut('Range collection using double override');
```

```
% Save simulation run name generated as collect_ranges. This run name is used in
% later steps to propose fixed point data types.
collect_ranges = converter.CurrentRunName;
sim_out = converter.simulateSystem();
```

Plot the regression accuracy before the conversion.

```
plotRegression(sim_out, baseline_output, system_under_design, 'Regression before conversion');
```



Propose Fixed-Point Data Types

Range information obtained from simulation can be used by the Fixed-Point Tool to propose fixed-point data types for blocks in the system under design. In this example, to ensure that the tools propose signed data types for all blocks in the subsystem, disable the ProposeSignedness option in the ProposalSettings object.

```
ps = DataTypeWorkflow.ProposalSettings;
ps.ProposeSignedness = false;
converter.proposeDataTypes(collect_ranges, ps);
```

Apply Proposed Data Types

By default, the Fixed-Point Tool applies all of the proposed data types. Use the `applyDataTypes` method to apply the data types. If you want to only apply a subset of the proposals, in the Fixed-Point Tool use the **Accept** check box to specify the proposals that you want to apply.

```
converter.applyDataTypes(collect_ranges);
```

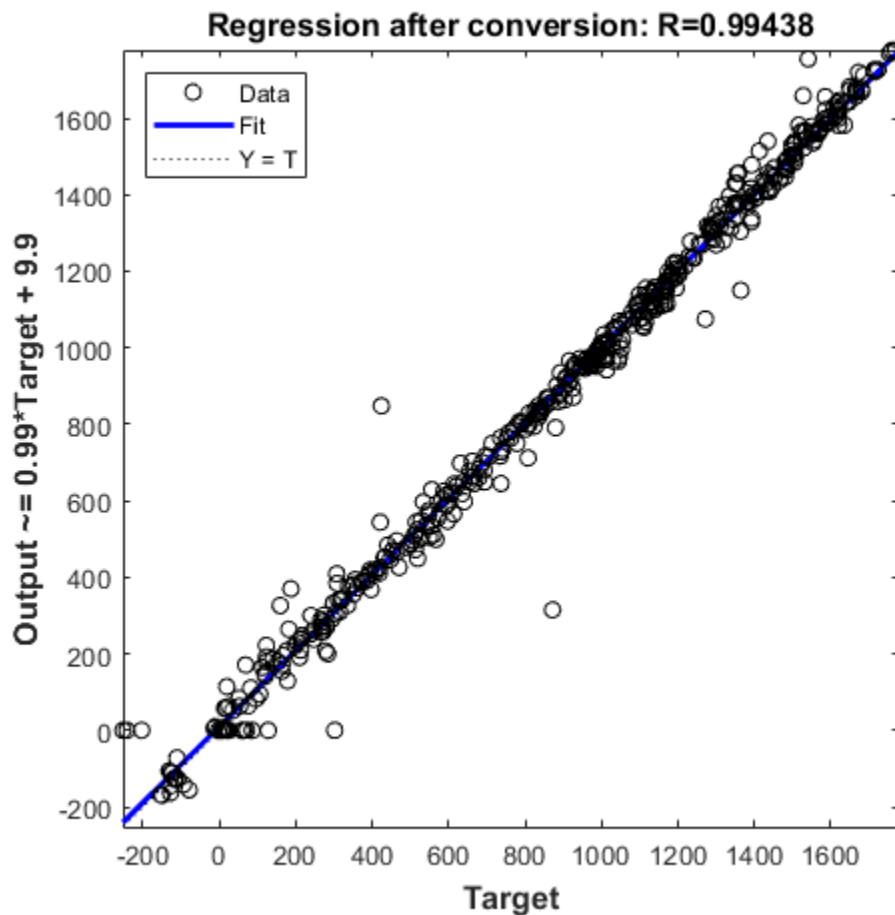
Verify Data Types

Proposed types should handle all possible inputs correctly. Set the model to simulate using the newly applied types, simulate the model, and observe that the neural network regression accuracy is retained post fixed-point conversion.

```
converter.applySettingsFromShortcut('Range collection with specified data types');
sim_out = converter.simulateSystem();
```

Plot the regression accuracy of the fixed-point model.

```
plotRegression(sim_out, baseline_output, system_under_design, 'Regression after conversion');
```



Replace Activation Function With an Optimized Lookup Table

The Tanh Activation function in Layer 1 can be replaced with either a lookup table or a CORDIC implementation for more efficient fixed-point code generation. In this example, we will be using the Lookup Table Optimizer to get a lookup table as a replacement for `tanh`. We will be using `EvenPow2Spacing` for faster execution speed. For more information, see <https://www.mathworks.com/help/fixedpoint/ref/functionapproximation.options-class.html>.

```
block_path = [system_under_design '/Layer 1/tansig'];
p = FunctionApproximation.Problem(block_path);
p.Options.WordLengths = 16;
p.Options.BreakpointSpecification = 'EvenPow2Spacing';
solution = p.solve;
solution.replaceWithApproximate;
```

| ID | Memory (bits) | Feasible | Table Size | Breakpoints | WLs | TableData | WL | BreakpointSpec |
|----|---------------|----------|------------|-------------|-----|-----------|----|----------------|
| 0 | 64 | 0 | 2 | | 16 | | 16 | EvenPow |
| 1 | 8224 | 1 | 512 | | 16 | | 16 | EvenPow |
| 2 | 4128 | 0 | 256 | | 16 | | 16 | EvenPow |

Best Solution

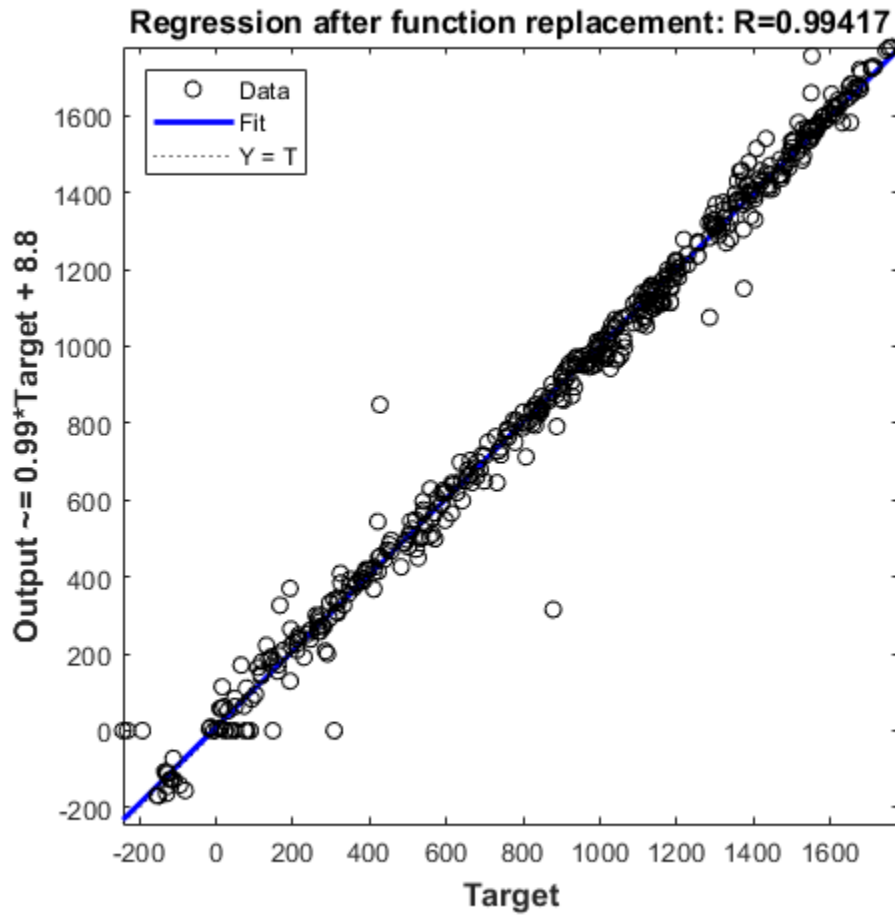
| ID | Memory (bits) | Feasible | Table Size | Breakpoints | WLs | TableData | WL | BreakpointSpec |
|----|---------------|----------|------------|-------------|-----|-----------|----|----------------|
| 1 | 8224 | 1 | 512 | | 16 | | 16 | EvenPow |

Verify model accuracy after function approximation replacement.

```
converter.applySettingsFromShortcut(converter.ShortcutsForSelectedSystem{2});
sim_out = converter.simulateSystem;
```

Plot regression accuracy.

```
plotRegression(sim_out, baseline_output, system_under_design, 'Regression after function replacer
```



Generate C Code

To generate C code using Simulink Coder™, right-click on the Function Fitting Neural Network subsystem, select **C/C++ Code > Build Subsystem**, then click the **Build** button when prompted for tunable parameters. You can also generate code by using the following command.

```
slbuild('fxpdemo_neuralnet_regression_toconvert/Function Fitting Neural Network')
```

Effects of Spacing on Speed, Error, and Memory Usage

In this section...

“Criteria for Comparing Types of Breakpoint Spacing” on page 41-59
 “Model That Illustrates Effects of Breakpoint Spacing” on page 41-59
 “Data ROM Required for Each Lookup Table” on page 41-59
 “Determination of Out-of-Range Inputs” on page 41-60
 “How the Lookup Tables Determine Input Location” on page 41-60
 “Interpolation for Each Lookup Table” on page 41-62
 “Summary of the Effects of Breakpoint Spacing” on page 41-63

Criteria for Comparing Types of Breakpoint Spacing

The sections that follow compare implementations of lookup tables that use breakpoints whose spacing is uneven, even, and power of two. The comparison focuses on:

- Execution speed of commands
- Rounding error during interpolation
- The amount of read-only memory (ROM) for data
- The amount of ROM for commands

This comparison is valid only when the breakpoints are not tunable. If the breakpoints are tunable in the generated code, all three cases generate the same code. For a summary of the effects of breakpoint spacing on execution speed, error, and memory usage, see “Summary of the Effects of Breakpoint Spacing” on page 41-63.

Model That Illustrates Effects of Breakpoint Spacing

This comparison uses the model `fxpdemo_approx_sin`. Three fixed-point lookup tables appear in this model. All three tables approximate the function $\sin(2\pi u)$ over the first quadrant and achieve a worst-case error of less than 2^{-8} . However, they have different restrictions on their breakpoint spacing.

You can use the model `fxpdemo_approx`, which `fxpdemo_approx_sin` opens, to generate Simulink Coder code (Simulink Coder software license required). The sections that follow present several segments of generated code to emphasize key differences.

To open the model, type at the MATLAB prompt:

```
fxpdemo_approx_sin
```

Data ROM Required for Each Lookup Table

This section looks at the data ROM required by each of the three spacing options.

Uneven Case

Uneven spacing requires both Y data points and breakpoints:

```
int16_T yuneven[8];
uint16_T xuneven[8];
```

The total bytes used are 32.

Even Case

Even spacing requires only Y data points:

```
int16_T yeven[10];
```

The total bytes used are 20. The breakpoints are not explicitly required. The code uses the spacing between the breakpoints, and might use the smallest and largest breakpoints. At most, three values related to the breakpoints are necessary.

Power of Two Case

Power of two spacing requires only Y data points:

```
int16_T ypow2[17];
```

The total bytes used are 34. The breakpoints are not explicitly required. The code uses the spacing between the breakpoints, and might use the smallest and largest breakpoints. At most, three values related to the breakpoints are necessary.

Determination of Out-of-Range Inputs

In all three cases, you must guard against the chance that the input is less than the smallest breakpoint or greater than the biggest breakpoint. There can be differences in how occurrences of these possibilities are handled. However, the differences are generally minor and are normally not a key factor in deciding to use one spacing method over another. The subsequent sections assume that out-of-range inputs are impossible or have already been handled.

How the Lookup Tables Determine Input Location

This section describes how the three fixed-point lookup tables determine where the current input is relative to the breakpoints.

Uneven Case

Unevenly spaced breakpoints require a general-purpose algorithm such as a binary search to determine where the input lies in relation to the breakpoints. The following code provides an example:

```
iLeft = 0;
iRght = 7; /* number of breakpoints minus 1 */

while ( ( iRght - iLeft ) > 1 )
{
    i = ( iLeft + iRght ) >> 1;

    if ( uAngle < xuneven[i] )
    {
        iRght = i;
    }
}
```



```

else
{
    iLeft = i;
}
}

```

The while loop executes up to $\log_2(N)$ times, where N is the number of breakpoints.

Even Case

Evenly spaced breakpoints require only one step to determine where the input lies in relation to the breakpoints:

```
iLeft = uAngle / 455U;
```

The divisor 455U represents the spacing between breakpoints. In general, the dividend would be $(uAngle - \text{SmallestBreakPoint})$. In this example, the smallest breakpoint is zero, so the code optimizes out the subtraction.

Power of Two Case

Power of two spaced breakpoints require only one step to determine where the input lies in relation to the breakpoints:

```
iLeft = uAngle >> 8;
```

The number of shifts is 8 because the breakpoints have spacing 2^8 . The smallest breakpoint is zero, so $uAngle$ replaces the general case of $(uAngle - \text{SmallestBreakPoint})$.

Comparison

To determine where the input lies with respect to the breakpoints, the unevenly spaced case requires much more code than the other two cases. This code requires additional command ROM. If many lookup tables share the binary search algorithm as a function, you can reduce this ROM penalty. Even if the code is shared, the number of clock cycles required to determine the location of the input is much higher for the unevenly spaced cases than the other two cases. If the code is shared, function call overhead decreases the speed of execution a little more.

In the evenly spaced case and the power of two spaced case, you can determine the location of the input with a single line of code. The evenly spaced case uses a general integer division. The power of two case uses a shift instead of general division because the divisor is an exact power of two. Without knowing the specific processor, you cannot be certain that a shift is better than division.

Many processors can implement division with a single assembly language instruction, so the code will be small. However, this instruction often takes many clock cycles to complete. Many processors do not provide a division instruction. Division on these processors occurs through repeated subtractions. This process is slow and requires a lot of machine code, but this code can be shared.

Most processors provide a way to do logical and arithmetic shifts left and right. A key difference is whether the processor can do N shifts in one instruction (barrel shift) or requires N instructions that shift one bit at a time. The barrel shift requires less code. Whether the barrel shift also increases speed depends on the hardware that supports the operation.

The compiler can also complicate the comparison. In the previous example, the command `uAngle >> 8` essentially takes the upper 8 bits in a 16-bit word. The compiler can detect this situation and

replace the bit shifts with an instruction that takes the bits directly. If the number of shifts is some other value, such as 7, this optimization would not occur.

Interpolation for Each Lookup Table

In theory, you can calculate the interpolation with the following code:

```
y = ( yData[iRight] - yData[iLeft] ) * ( u - xData[iLeft] ) ...  
    / ( xData[iRight] - xData[iLeft] ) + yData[iLeft]
```

The term $(xData[iRight] - xData[iLeft])$ is the spacing between neighboring breakpoints. If this value is constant, due to even spacing, some simplification is possible. If spacing is not just even but also a power of two, significant simplifications are possible for fixed-point implementations.

Uneven Case

For the uneven case, one possible implementation of the ideal interpolation in fixed point is:

```
xNum = uAngle          - xuneven[iLeft];  
xDen  = xuneven[iRight] - xuneven[iLeft];  
yDiff = yuneven[iRight] - yuneven[iLeft];  
  
MUL_S32_S16_U16( bigProd, yDiff, xNum );  
  
    DIV_NZP_S16_S32_U16_FLOOR( yDiff, bigProd, xDen );  
  
    yUneven = yuneven[iLeft] + yDiff;
```

The multiplication and division routines are not shown here. These routines can be complex and depend on the target processor. For example, these routines look different for a 16-bit processor than for a 32-bit processor.

Even Case

Evenly spaced breakpoints implement interpolation using slightly different calculations than the uneven case. The key difference is that the calculations do not directly use the breakpoints. When the breakpoints are not required in ROM, you can save a lot of memory.

```
xNum = uAngle - ( iLeft * 455U );  
  
    yDiff = yeven[iLeft+1] - yeven[iLeft];  
  
    MUL_S32_S16_U16( bigProd, yDiff, xNum );  
  
    DIV_NZP_S16_S32_U16_FLOOR( yDiff, bigProd, 455U );  
  
    yEven = yeven[iLeft] + yDiff;
```

Power of Two Case

Power of two spaced breakpoints implement interpolation using very different calculations than the other two cases. As in the even case, breakpoints are not used in the generated code and therefore not required in ROM.

```
lambda = uAngle & 0x00FFU;
```

```

yPow2 = ypow2[iLeft]+1] - ypow2[iLeft];

MUL_S16_U16_S16_SR8(yPow2, lambda, yPow2);

yPow2 += ypow2[iLeft];

```

This implementation has significant advantages over the uneven and even implementations:

- A bitwise AND combined with a shift right at the end of the multiplication replaces a subtraction and a division.
- The term $(u - xData[iLeft]) / (xData[iRight] - xData[iLeft])$ results in no loss of precision, because the spacing is a power of two.

In contrast, the uneven and even cases usually introduce rounding error in this calculation.

Summary of the Effects of Breakpoint Spacing

The following table summarizes the effects of breakpoint spacing on execution speed, error, and memory usage.

| Parameter | Even Power of 2 Spaced Data | Evenly Spaced Data | Unevenly Spaced Data |
|-----------------|---|--|---|
| Execution speed | The execution speed is the fastest. The position search and interpolation are the same as for evenly spaced data. However, to increase the speed more, a bit shift replaces the position search, and a bit mask replaces the interpolation. | The execution speed is faster than that for unevenly spaced data, because the position search is faster and the interpolation requires a simple division. | The execution speed is the slowest of the different spacings because the position search is slower, and the interpolation requires more operations. |
| Error | The error can be larger than that for unevenly spaced data because approximating a function with nonuniform curvature requires more points to achieve the same accuracy. | The error can be larger than that for unevenly spaced data because approximating a function with nonuniform curvature requires more points to achieve the same accuracy. | The error can be smaller because approximating a function with nonuniform curvature requires fewer points to achieve the same accuracy. |
| ROM usage | Uses less command ROM, but more data ROM. | Uses less command ROM, but more data ROM. | Uses more command ROM, but less data ROM. |
| RAM usage | Not significant. | Not significant. | Not significant. |

The number of Y data points follows the expected pattern. For the same worst-case error, unrestricted spacing (uneven) requires the fewest data points, and power-of-two-spaced breakpoints require the most. However, the implementation for the evenly spaced and the power of two cases does not need the breakpoints in the generated code. This reduces their data ROM requirements by half. As a result, the evenly spaced case actually uses less data ROM than the unevenly spaced case. Also, the power of two case requires only slightly more ROM than the uneven case. Changing the worst-case error can change these rankings. Nonetheless, when you compare data ROM usage, you should always take into account the fact that the evenly spaced and power of two spaced cases do not require their breakpoints in ROM.

The effort of determining where the current input is relative to the breakpoints strongly favors the evenly spaced and power of two spaced cases. With uneven spacing, you use a binary search method that loops up to $\log_2(N)$ times. With even and power of two spacing, you can determine the location with the execution of one line of C code. But you cannot decide the relative advantages of power of two versus evenly spaced without detailed knowledge of the hardware and the C compiler.

The effort of calculating the interpolation favors the power of two case, which uses a bitwise AND operation and a shift to replace a subtraction and a division. The advantage of this behavior depends on the specific hardware, but you can expect an advantage in code size, speed, and also in accuracy. The evenly spaced case calculates the interpolation with a minor improvement in efficiency over the unevenly spaced case.

Approximate Functions with a Direct Lookup Table

Using the Lookup Table Optimizer, you can generate a direct lookup table approximating a Simulink block, or a function. Direct lookup tables are efficient to implement on hardware because they do not require any calculations.

Generate a Two-Dimensional Direct Lookup Table Approximation

Create a `FunctionApproximation.Problem` object specifying the function for which to generate the approximate. To generate a direct lookup table, set the interpolation method to `None` in the `FunctionApproximation.Options` object.

```
problem = FunctionApproximation.Problem('atan2');
problem.InputTypes = [numerictype(0,4,2) numerictype(0,8,4)];
problem.OutputType = fixdt(0,8,7);
problem.Options.Interpolation = "None";
problem.Options.AbsTol = 2^-4;
problem.Options.RelTol = 0;
problem.Options.WordLengths = 1:8;
```

Use the `solve` method to generate the optimal lookup table.

```
solution = solve(problem)
```

| ID | Memory (bits) | Feasible | Table Size | Intermediate WLS | TableData WL | |
|----|---------------|----------|------------|------------------|--------------|----------|
| 0 | 32768 | 1 | [16 256] | [4 8] | 8 | 6.250000 |
| 1 | 28672 | 1 | [16 256] | [4 8] | 7 | 6.250000 |
| 2 | 24576 | 1 | [16 256] | [4 8] | 6 | 6.250000 |
| 3 | 16384 | 1 | [16 128] | [4 7] | 8 | 6.250000 |
| 4 | 14336 | 1 | [16 128] | [4 7] | 7 | 6.250000 |
| 5 | 12288 | 1 | [16 128] | [4 7] | 6 | 6.250000 |
| 6 | 10240 | 0 | [16 128] | [4 7] | 5 | 6.250000 |
| 7 | 8192 | 0 | [16 128] | [4 7] | 4 | 6.250000 |

Best Solution

| ID | Memory (bits) | Feasible | Table Size | Intermediate WLS | TableData WL | |
|----|---------------|----------|------------|------------------|--------------|----------|
| 5 | 12288 | 1 | [16 128] | [4 7] | 6 | 6.250000 |

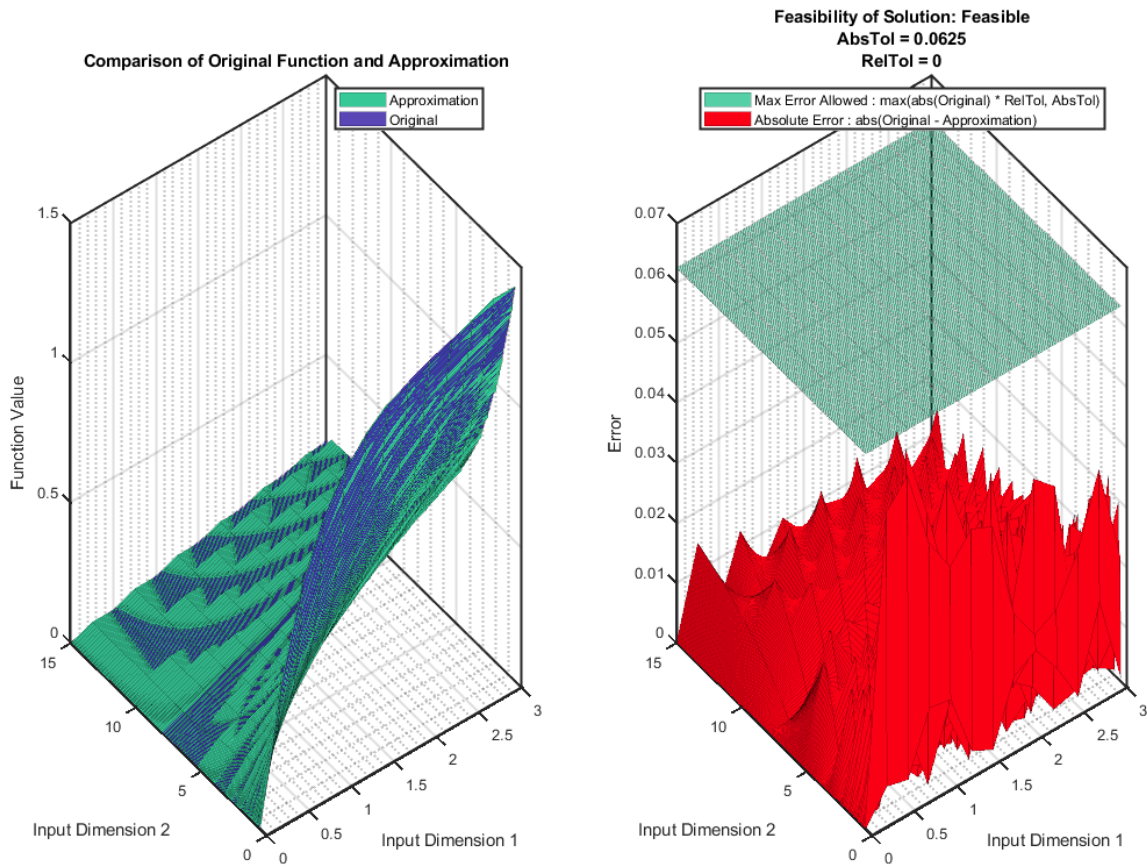
```
solution =
```

```
1x1 FunctionApproximation.LUTSolution with properties:
```

```
    ID: 5
  Feasible: "true"
```

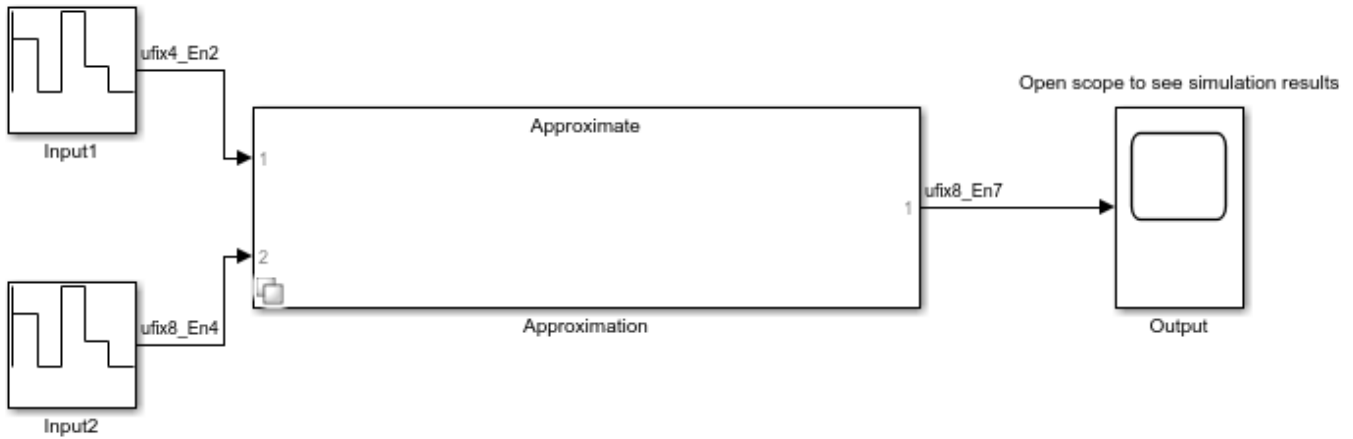
Use the `compare` method to compare the output of the original function and the approximate.

```
compare(solution);
```



Use the approximate method to generate a Simulink™ subsystem containing the generated direct lookup table.

```
approximate(solution)
```

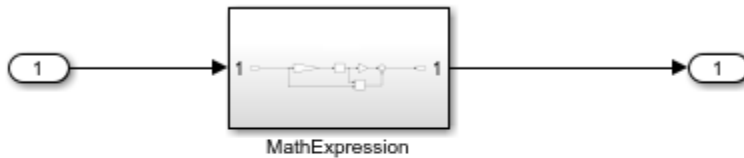


Generate a Direct Lookup Table Approximation for a Subsystem

This example shows how to approximate a Simulink™ subsystem with a direct lookup table.

Open the model containing the subsystem to approximate.

```
functionToApproximate = 'ex_direct_approximation/MathExpression';
open_system('ex_direct_approximation');
```



Copyright 2018 The MathWorks, Inc.

To generate a direct lookup table, set the interpolation method to None.

```
problem = FunctionApproximation.Problem(functionToApproximate);
problem.Options.Interpolation = 'None';
problem.Options.RelTol = 0;
problem.Options.AbsTol = 0.2;
problem.Options.WordLengths = [7 8 9 16];
solution = solve(problem);
```

| ID | Memory (bits) | Feasible | Table Size | Intermediate WLS | TableData WL | |
|----|---------------|----------|------------|------------------|--------------|----------|
| 0 | 2097152 | 1 | 65536 | 16 | 32 | 2.000000 |
| 1 | 896 | 0 | 128 | 7 | 7 | 2.000000 |
| 2 | 1024 | 0 | 128 | 7 | 8 | 2.000000 |
| 3 | 1152 | 0 | 128 | 7 | 9 | 2.000000 |
| 4 | 2048 | 0 | 128 | 7 | 16 | 2.000000 |

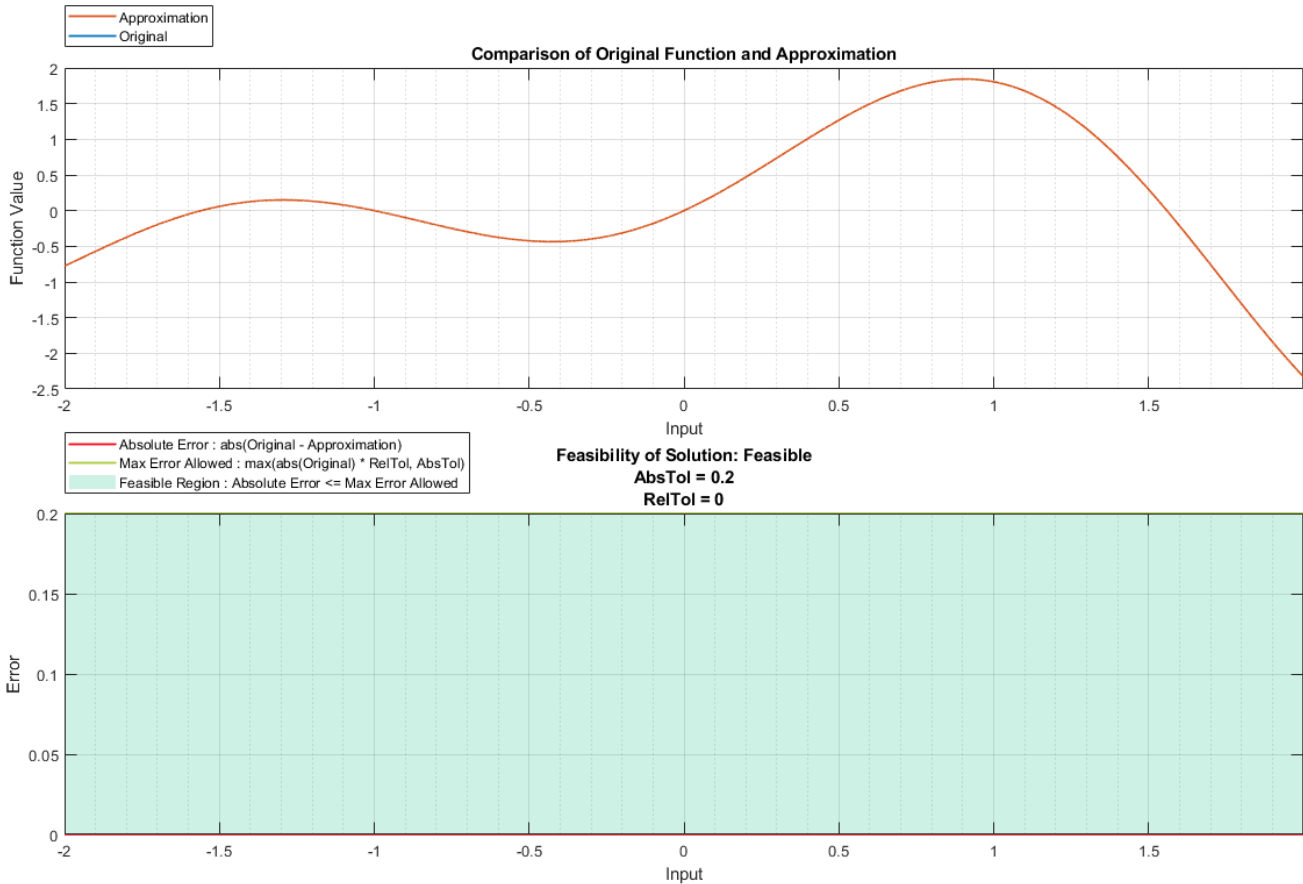
Best Solution

| ID | Memory (bits) | Feasible | Table Size | Intermediate WLS | TableData WL |
|----|---------------|----------|------------|------------------|--------------|
| 0 | 2097152 | 1 | 65536 | 16 | 32 |

2.00000

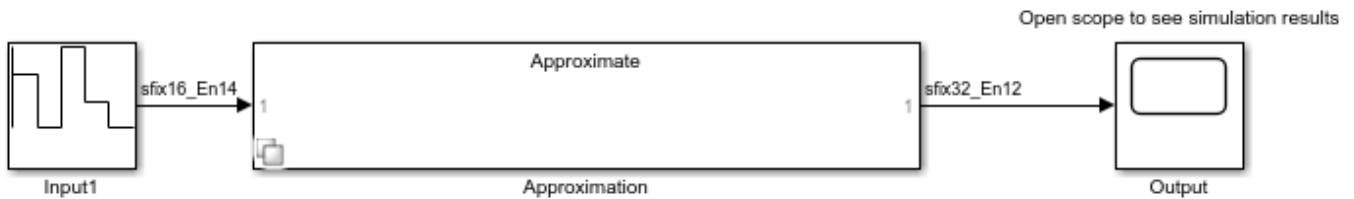
Compare the original subsystem behavior to the lookup table approximation.

```
compare(solution);
```



Generate a new subsystem containing the lookup table approximation.

```
approximate(solution);
```



Replace the original subsystem with the new subsystem containing the lookup table approximation.

```
replaceWithApproximate(solution);
```


You can revert your model back to its original state using the `revertToOriginal` function. This function places the original subsystem back into the model.

```
revertToOriginal(solution);
```

Convert Digit Recognition Neural Network to Fixed Point and Generate C Code

This example shows how to convert a neural network classification model in Simulink™ to fixed point using the Fixed-Point Tool and Lookup Table Optimizer. Following the conversion, you can generate C code using Simulink Coder.

Overview

Using the Fixed-Point Tool, you can convert a design from floating point to fixed point. Use the Lookup Table Optimizer to generate memory-efficient lookup table replacements for unbounded functions such as `exp` and `log2`. Using these tools, this example shows how to convert a trained floating-point neural network classification model to use embedded-efficient fixed-point data types.

Digit Classification and MNIST Dataset

MNIST handwritten digit dataset is a commonly used dataset in the field of neural networks. For an example showing a simple way to create a two-layered neural network using this dataset, see *Artificial Neural Networks for Beginners*.

Data and Neural Network Training

Download the training and test MNIST files according to the directions in *Artificial Neural Networks for Beginners*. Load the data and train the network.

```
%Load Data
tr = csvread('train.csv', 1, 0);           % read train.csv
sub = csvread('test.csv', 1, 0);          % read test.csv

% Prepare Data
n = size(tr, 1);                          % number of samples in the dataset
targets = tr(:,1);                        % 1st column is |label|
targets(targets == 0) = 10;               % use '10' to present '0'
targetsd = dummyvar(targets);             % convert label into a dummy variable
inputs = tr(:,2:end);                     % the rest of columns are predictors

inputs = inputs';                          % transpose input
targets = targets';                        % transpose target
targetsd = targetsd';                     % transpose dummy variable

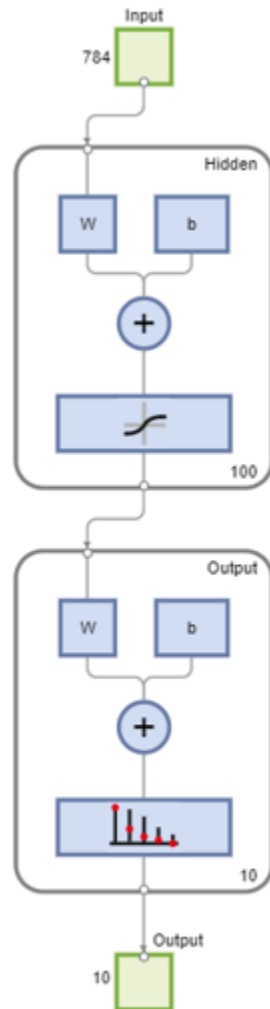
rng(1);                                    % for reproducibility
c = cvpartition(n, 'Holdout', n/3);       % hold out 1/3 of the dataset

Xtrain = inputs(:, training(c));          % 2/3 of the input for training
Ytrain = targetsd(:, training(c));        % 2/3 of the target for training
Xtest = inputs(:, test(c));               % 1/3 of the input for testing
Ytest = targetsd(test(c));                % 1/3 of the target for testing
Ytestd = targetsd(:, test(c));            % 1/3 of the dummy variable for testing

% Train Network
hiddenLayerSize = 100;
net = patternnet(hiddenLayerSize);

[net, tr] = train(net, Xtrain, Ytrain);
view(net);
```

```
outputs = net(Xtest);  
errors = gsubtract(Ytest, outputs);  
performance = perform(net, Ytest, outputs);  
  
figure, plotperform(tr);
```



Network Diagram

Training Results

Training finished: Met validation criterion ✔

Training Progress

| Unit | Initial Value | Stopped Value | Target Value |
|-------------------|---------------|---------------|--------------|
| Epoch | 0 | 95 | 1000 |
| Elapsed Time | - | 00:00:42 | - |
| Performance | 0.66 | 0.00169 | 0 |
| Gradient | 3.38 | 0.00396 | 1e-06 |
| Validation Checks | 0 | 6 | 6 |

Training Algorithms

Data Division: Random dividerand

Training: Scaled Conjugate Gradient trainscg

Performance: Cross Entropy crossentropy

Calculations: MEX

Training Plots

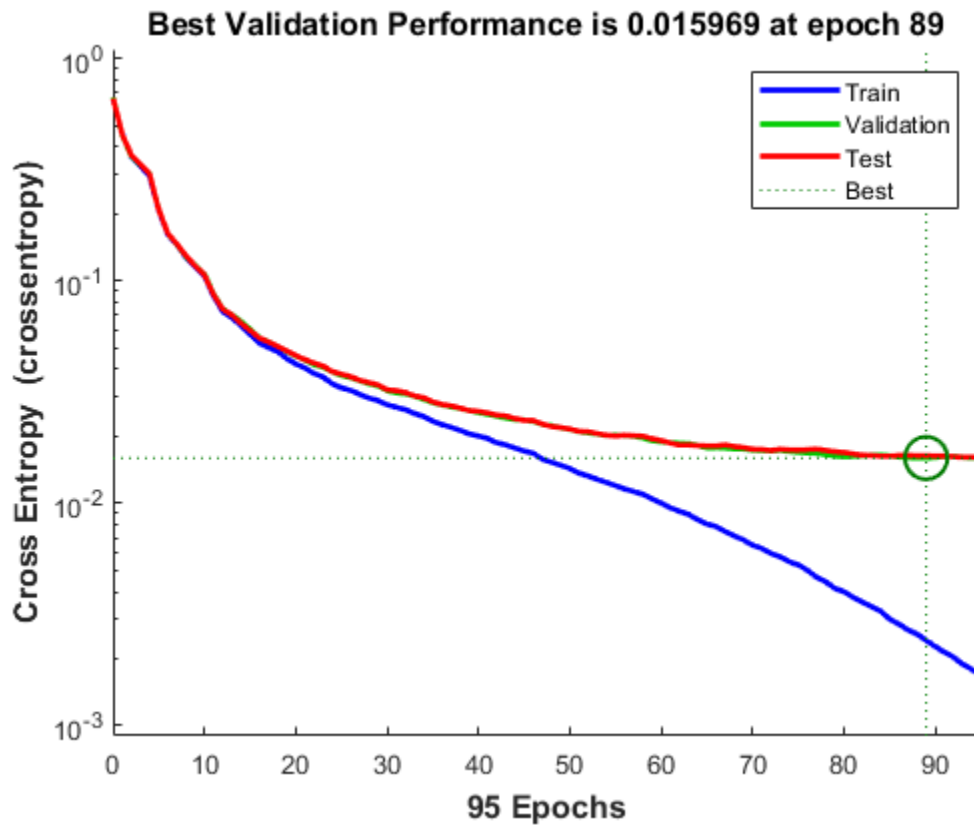
Performance

Training State

Error Histogram

Confusion

Receiver Operating Characteristic



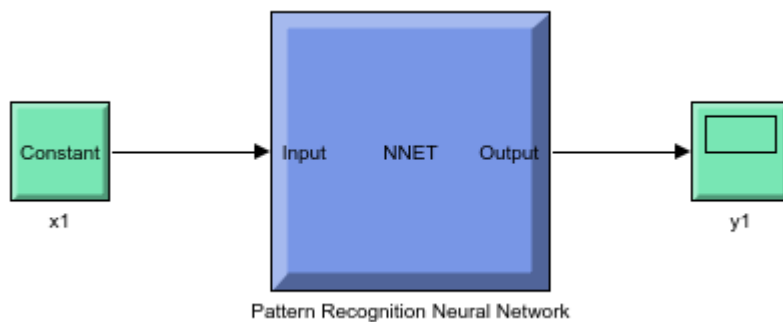
Close the view of the network.

```
nnet.guis.closeAllViews();
```

Model Preparation for Fixed-Point Conversion

After training the network, use the `gensim` function from the Deep Learning Toolbox™ to generate a Simulink model.

```
sys_name = gensim(net, 'Name', 'mTrainedNN');
```

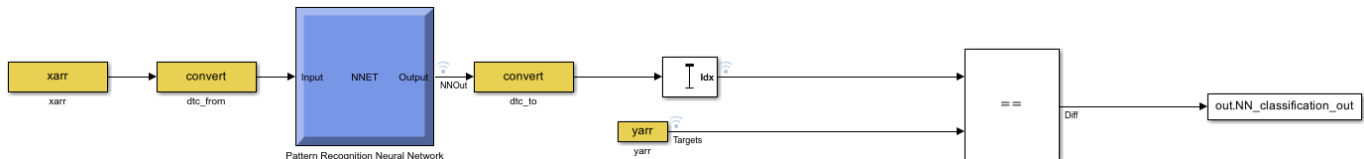


The model generated by the `gensim` function contains the neural network with trained weights and biases. Prepare the trained neural network for conversion to fixed point by enabling signal logging at

the output of the network, and adding input stimuli and verification blocks. The modified model is `fxpdemo_mnist_classification`.

Open and inspect the model.

```
model = 'fxpdemo_mnist_classification';
system_under_design = [model '/Pattern Recognition Neural Network'];
baseline_output = [model '/yarr'];
open_system(model);
```



To open the Fixed-Point Tool, right-click the Function Fitting Neural Network subsystem and select **Fixed-Point Tool**. Alternatively, use the command-line interface of the Fixed-Point Tool. The Fixed-Point Tool and its command-line interface help you prepare your model for conversion, and convert your system to fixed point. You can use the Fixed-Point Tool to collect range and overflow instrumentation of objects in your model via simulation and range analysis. In this example, use the command-line interface of the Fixed-Point Tool to convert the neural network to fixed point.

```
converter = DataTypeWorkflow.Converter(system_under_design);
```

Run Simulation to Collect Ranges

Simulate the model with instrumentation to collect ranges. Enable instrumentation using the 'Range collection using double override' shortcut. Save the simulation run name for use in later steps.

```
converter.applySettingsFromShortcut('Range collection using double override');
collect_ranges = converter.CurrentRunName;
sim_out = converter.simulateSystem();
```

Plot the correct classification rate before the conversion to establish baseline behavior.

```
plotConfusionMatrix(sim_out, baseline_output, system_under_design, 'Classification rate before c
```

Classification rate before conversion Confusion Matrix

| | | | | | | | | | | | | |
|--------------|----|---------------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| Output Class | 1 | 1632 11.7% | 10 0.1% | 3 0.0% | 2 0.0% | 0 0.0% | 0 0.0% | 6 0.0% | 5 0.0% | 2 0.0% | 0 0.0% | 98.3% 1.7% |
| | 2 | 10 0.1% | 1325 9.5% | 9 0.1% | 10 0.1% | 5 0.0% | 5 0.0% | 15 0.1% | 15 0.1% | 7 0.0% | 7 0.0% | 94.1% 5.9% |
| | 3 | 2 0.0% | 25 0.2% | 1368 9.8% | 0 0.0% | 30 0.2% | 3 0.0% | 8 0.1% | 24 0.2% | 10 0.1% | 5 0.0% | 92.7% 7.3% |
| | 4 | 3 0.0% | 11 0.1% | 0 0.0% | 1227 8.8% | 0 0.0% | 8 0.1% | 3 0.0% | 5 0.0% | 33 0.2% | 1 0.0% | 95.0% 5.0% |
| | 5 | 4 0.0% | 4 0.0% | 18 0.1% | 1 0.0% | 1181 8.4% | 8 0.1% | 2 0.0% | 16 0.1% | 14 0.1% | 4 0.0% | 94.3% 5.7% |
| | 6 | 1 0.0% | 0 0.0% | 1 0.0% | 9 0.1% | 13 0.1% | 1344 9.6% | 2 0.0% | 12 0.1% | 1 0.0% | 6 0.0% | 96.8% 3.2% |
| | 7 | 8 0.1% | 11 0.1% | 2 0.0% | 3 0.0% | 3 0.0% | 0 0.0% | 1410 10.1% | 2 0.0% | 15 0.1% | 5 0.0% | 96.6% 3.4% |
| | 8 | 13 0.1% | 9 0.1% | 12 0.1% | 3 0.0% | 14 0.1% | 8 0.1% | 3 0.0% | 1253 8.9% | 15 0.1% | 8 0.1% | 93.6% 6.4% |
| | 9 | 6 0.0% | 0 0.0% | 12 0.1% | 14 0.1% | 4 0.0% | 2 0.0% | 24 0.2% | 7 0.0% | 1301 9.3% | 5 0.0% | 94.6% 5.4% |
| | 10 | 1 0.0% | 6 0.0% | 0 0.0% | 0 0.0% | 2 0.0% | 8 0.1% | 0 0.0% | 3 0.0% | 0 0.0% | 1334 9.5% | 98.5% 1.5% |
| | | | 97.1% 2.9% | 94.6% 5.4% | 96.0% 4.0% | 96.7% 3.3% | 94.3% 5.7% | 97.0% 3.0% | 95.7% 4.3% | 93.4% 6.6% | 93.1% 6.9% | 97.0% 3.0% |
| | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | |
| | | Target Class | | | | | | | | | | |

Propose Fixed-Point Data Types

The Fixed-Point Tool uses range information obtained through simulation to propose fixed-point data types for blocks in the system under design. In this example, to ensure that the tools propose signed data types for all blocks in the subsystem, disable the `ProposeSignedness` option in the `ProposalSettings` object.

```
ps = DataTypeWorkflow.ProposalSettings;
converter.proposeDataTypes(collect_ranges, ps);
```

Apply Proposed Data Types

By default, the Fixed-Point Tool applies all of the proposed data types. Use the `applyDataTypes` method to apply the data types. If you want to only apply a subset of the proposals, in the Fixed-Point Tool use the **Accept** check box to specify the proposals that you want to apply.

```
converter.applyDataTypes(collect_ranges);
```

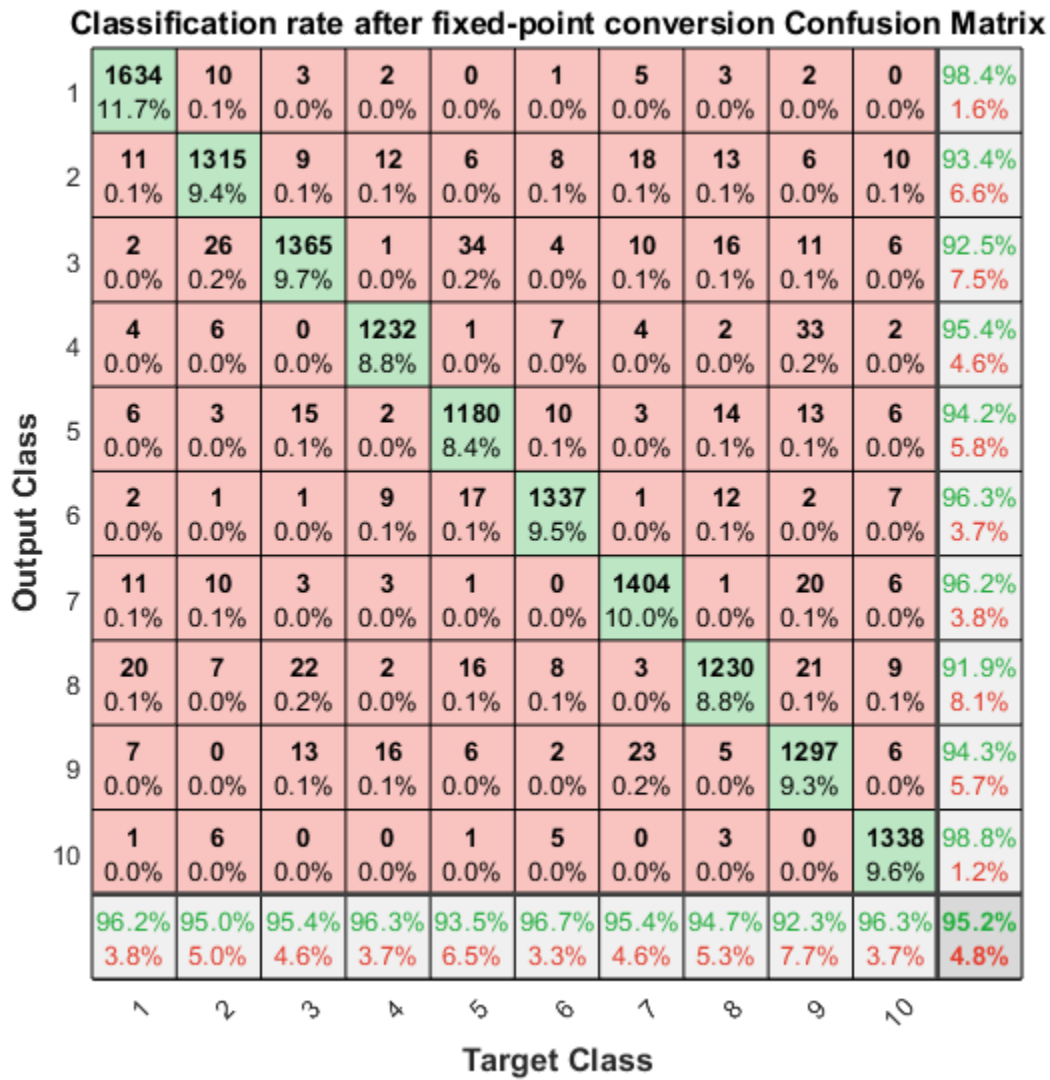
Verify Data Types

Proposed types should handle all possible inputs correctly. Set the model to simulate using the newly applied types, simulate the model, and observe that the neural network regression accuracy remains the same after the conversion.

```
converter.applySettingsFromShortcut('Range collection with specified data types');
sim_out = converter.simulateSystem();
```

Plot the correct classification rate of the fixed-point model.

```
plotConfusionMatrix(sim_out, baseline_output, system_under_design, 'Classification rate after fi
```



Replace Activation Functions With an Optimized Lookup Table

For more efficient code, replace the Tanh Activation function in the first layer with either a lookup table or a CORDIC implementation. In this example, use the Lookup Table Optimizer to get a lookup table to replace `tanh`. In this example, specify `EvenPow2Spacing` for the breakpoint spacing for faster execution speed.

```
block_path = [system_under_design '/Layer 1/tansig'];
p = FunctionApproximation.Problem(block_path);
p.Options.WordLengths = 16;
p.Options.BreakpointSpecification = 'EvenPow2Spacing';
solution = p.solve;
solution.replaceWithApproximate;
```

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|---------------|----------|------------|-----------------|--------------|----------------|
| 0 | 64 | 0 | 2 | 16 | 16 | EvenPow |
| 1 | 8224 | 1 | 512 | 16 | 16 | EvenPow |
| 2 | 4128 | 1 | 256 | 16 | 16 | EvenPow |
| 3 | 2080 | 0 | 128 | 16 | 16 | EvenPow |

Best Solution

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|---------------|----------|------------|-----------------|--------------|----------------|
| 2 | 4128 | 1 | 256 | 16 | 16 | EvenPow |

Follow the same steps to replace the `exp` function in the softmax implementation in the second layer with a lookup table.

```
block_path = [system_under_design '/Layer 2/softmax/Exp'];
p = FunctionApproximation.Problem(block_path);
p.Options.WordLengths = 16;
p.Options.BreakpointSpecification = 'EvenPow2Spacing';
```

To get an optimized lookup table, define finite lower and upper bounds for the inputs.

```
p.InputLowerBounds = -40;
p.InputUpperBounds = 0;
solution = p.solve;
solution.replaceWithApproximate;
```

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|---------------|----------|------------|-----------------|--------------|----------------|
| 0 | 64 | 0 | 2 | 16 | 16 | EvenPow |
| 1 | 2608 | 1 | 161 | 16 | 16 | EvenPow |
| 2 | 1328 | 0 | 81 | 16 | 16 | EvenPow |

Best Solution

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|---------------|----------|------------|-----------------|--------------|----------------|
| 1 | 2608 | 1 | 161 | 16 | 16 | EvenPow |

Verify model accuracy after replacing the functions with the lookup table approximations.

```
converter.applySettingsFromShortcut(converter.ShortcutsForSelectedSystem{2});
sim_out = converter.simulateSystem;
```

```
plotConfusionMatrix(sim_out, baseline_output, system_under_design, 'Classification rate after fun
```

Classification rate after function replacement Confusion Matrix

| | | | | | | | | | | | | |
|--------------|----|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| Output Class | 1 | 1634 11.7% | 10 0.1% | 3 0.0% | 2 0.0% | 0 0.0% | 1 0.0% | 5 0.0% | 3 0.0% | 2 0.0% | 0 0.0% | 98.4% 1.6% |
| | 2 | 11 0.1% | 1314 9.4% | 10 0.1% | 12 0.1% | 6 0.0% | 8 0.1% | 18 0.1% | 13 0.1% | 6 0.0% | 10 0.1% | 93.3% 6.7% |
| | 3 | 2 0.0% | 26 0.2% | 1365 9.7% | 1 0.0% | 35 0.2% | 4 0.0% | 10 0.1% | 15 0.1% | 11 0.1% | 6 0.0% | 92.5% 7.5% |
| | 4 | 4 0.0% | 6 0.0% | 0 0.0% | 1232 8.8% | 1 0.0% | 7 0.0% | 4 0.0% | 2 0.0% | 33 0.2% | 2 0.0% | 95.4% 4.6% |
| | 5 | 6 0.0% | 3 0.0% | 15 0.1% | 2 0.0% | 1181 8.4% | 10 0.1% | 3 0.0% | 13 0.1% | 13 0.1% | 6 0.0% | 94.3% 5.7% |
| | 6 | 2 0.0% | 1 0.0% | 1 0.0% | 9 0.1% | 17 0.1% | 1337 9.5% | 1 0.0% | 12 0.1% | 2 0.0% | 7 0.0% | 96.3% 3.7% |
| | 7 | 11 0.1% | 10 0.1% | 3 0.0% | 3 0.0% | 1 0.0% | 0 0.0% | 1404 10.0% | 1 0.0% | 20 0.1% | 6 0.0% | 96.2% 3.8% |
| | 8 | 20 0.1% | 7 0.0% | 22 0.2% | 2 0.0% | 16 0.1% | 8 0.1% | 3 0.0% | 1230 8.8% | 21 0.1% | 9 0.1% | 91.9% 8.1% |
| | 9 | 7 0.0% | 0 0.0% | 13 0.1% | 16 0.1% | 6 0.0% | 2 0.0% | 23 0.2% | 5 0.0% | 1297 9.3% | 6 0.0% | 94.3% 5.7% |
| | 10 | 1 0.0% | 6 0.0% | 0 0.0% | 0 0.0% | 1 0.0% | 5 0.0% | 0 0.0% | 3 0.0% | 0 0.0% | 1338 9.6% | 98.8% 1.2% |
| | | 96.2% 3.8% | 95.0% 5.0% | 95.3% 4.7% | 96.3% 3.7% | 93.4% 6.6% | 96.7% 3.3% | 95.4% 4.6% | 94.8% 5.2% | 92.3% 7.7% | 96.3% 3.7% | 95.2% 4.8% |
| | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | |
| | | Target Class | | | | | | | | | | |

Generate C code

To generate C code, right-click the Function Fitting Neural Network subsystem, select C/C++ Code > Build Subsystem. Click the **Build** button when prompted for tunable parameters.

References

[1] LeCun, Y., C. Cortes, and C. J. C. Burges. "The MNIST Database of Handwritten Digits." <http://yann.lecun.com/exdb/mnist/>.

See Also

Fixed-Point Tool | Lookup Table Optimizer

Calculate Complex dB Using a Direct Lookup Table

You can calculate complex decibel levels using the following formula.

$$dB = 20 \times \log_{10}(\sqrt{\Re^2 + \Im^2})$$

However, this equation contains expressions, such as the log calculation, that are not efficient to implement on hardware. Using a direct lookup table, you can very closely approximate this expression in a way that is efficient on hardware.

To begin, define the function to approximate with a lookup table.

```
f = @(re,im) 20*log10(sqrt(re.^2 + im.^2));
```

To specify the tolerances that are acceptable, and the desired word lengths to use in the lookup table, use the `FunctionApproximation.Options` object. To generate a direct lookup table, set the `Interpolation` property of the `Options` object to `None`. Use the `ApproximateSolutionType` property to specify whether to return the lookup table as a Simulink™ subsystem or as a MATLAB® function.

```
options = FunctionApproximation.Options('Interpolation', 'None', 'AbsTol', 0.25, 'RelTol', 0, 'W
```

```
% Problem setup
```

```
problem = FunctionApproximation.Problem(f, 'Options', options);
problem.InputTypes = [numerictype(0,5,0) numerictype(0,5,0)];
problem.InputLowerBounds = [1 1];
problem.InputUpperBounds = [Inf Inf]; % upper bound will clip to input types range
problem.OutputType = numerictype(0,10,4);
```

The `solve` function returns the optimal lookup table as a `FunctionApproximation.LUTSolution` object. As the software optimizes the parameters of the lookup table, MATLAB® displays information about each iteration of the optimization, including the total memory used by the lookup table, the word lengths used for data in the lookup table, and the maximum difference in output between the original function and the lookup table approximation. The best solution is defined as the lookup table using the smallest memory that meets the tolerances and other constraints defined in the `Options` object.

```
solution = solve(problem)
```

```
Upper bound for input 2 has been set to the maximum representable value of the type numerictype(0,5,0)
```

```
Upper bound for input 1 has been set to the maximum representable value of the type numerictype(0,5,0)
```

| ID | Memory (bits) | Feasible | Table Size | Intermediate WLS | TableData WL | |
|----|---------------|----------|------------|------------------|--------------|----------|
| 0 | 10240 | 1 | [32 32] | [5 5] | 10 | 2.500000 |
| 1 | 9216 | 1 | [32 32] | [5 5] | 9 | 2.500000 |
| 2 | 8192 | 1 | [32 32] | [5 5] | 8 | 2.500000 |
| 3 | 7168 | 1 | [32 32] | [5 5] | 7 | 2.500000 |

```
Best Solution
```

| ID | Memory (bits) | Feasible | Table Size | Intermediate WLS | TableData WL | |
|----|---------------|----------|------------|------------------|--------------|----------|
| 3 | 7168 | 1 | [32 32] | [5 5] | 7 | 2.500000 |

```
solution =
```

1x1 FunctionApproximation.LUTSolution with properties:

ID: 3
Feasible: "true"

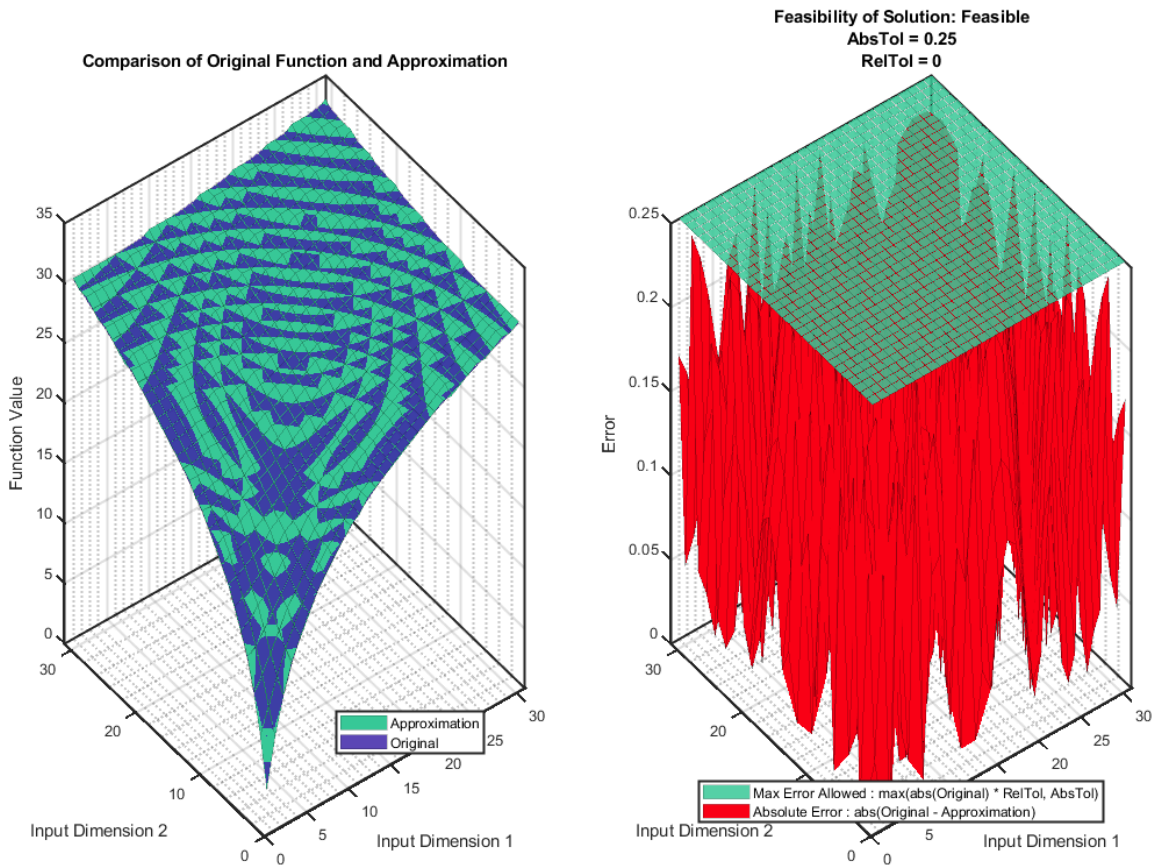
Compare the output of the original function and the lookup table approximation. The left plot shows the output of the original function defined in the Problem object, and the lookup table approximation. The plot on the right shows the difference between the output of the original function and the corresponding output from the generated lookup table approximation. The difference between the two outputs is less than the tolerance specified in the Options object.

```
compareData = compare(solution)
```

```
compareData =
```

1x2 struct array with fields:

```
Breakpoints  
Original  
Approximate
```



Access the `TableData` property of the `solution` to use the lookup table in a MATLAB® application.

```
tableData = solution.TableData
```

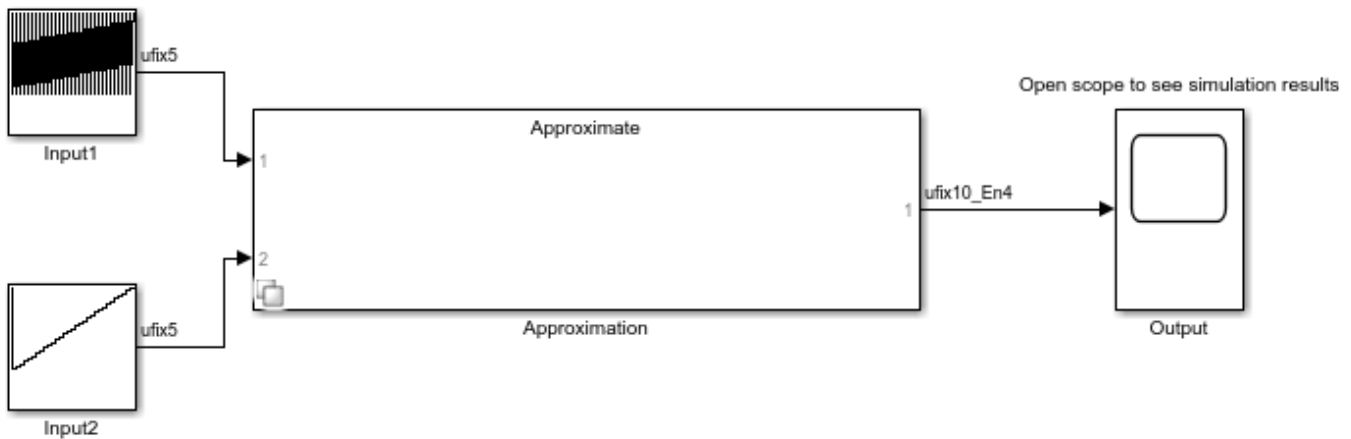
```
tableData =
```

```
struct with fields:
```

```
BreakpointValues: {1x2 cell}
BreakpointDataTypes: [2x1 embedded.numericity]
TableValues: [32x32 double]
TableDataType: [1x1 embedded.numericity]
IsEvenSpacing: 1
Interpolation: None
```

Use the `approximate` function to generate a Simulink™ subsystem containing the lookup table approximation.

```
approximate(solution)
```



You can use the generated subsystem containing the lookup table in an HDL application. To check that the lookup table is compatible with HDL code generation using HDL Coder™, use the `checkhdl` function.

Optimize Lookup Tables for Periodic Functions

This example shows how to use translational and reflectional symmetries in functions in order to optimize lookup tables.

Functions with Symmetry

Functions have symmetry if the functions are unchanged by simple mathematical operations such as translation, rotation, or reflection. When a function has symmetry, the entire range of the function can be generated from a smaller region of the function. After constructing this smaller region, you can translate, reflect, and rotate it to obtain the remainder of the function. This property of functions with symmetry is useful for embedded applications.

Periodic Functions and Discrete Translational Symmetry

Periodic functions are functions with discrete translational symmetry. Discrete translational symmetry means that for the period T , $f(x) = f(x + nT)$ for any integer n . Given any x , you can calculate the value of $f(x)$ by using the identity $f(x) = f(x \bmod T)$. This calculation is equivalent to translating x by an integer multiple of T , such that $0 \leq x - nT < T$.

Periodic Functions and Reflectional Symmetry

Functions can also have reflectional symmetries about lines or points in the Cartesian plane. The most common examples are odd and even functions. Even functions are symmetric about the y-axis, meaning $f(x) = f(-x)$. Odd functions are symmetric under reflection through the origin, which is expressed as $-f(x) = f(-x)$. Functions may also be even or odd with respect to other lines and points in the plane. A function that is even with respect to the line $x = C$ will obey $f(C + \Delta x) = f(C - \Delta x)$. Similarly, a function that is odd with respect to the point $(x, y) = (C, 0)$ will obey $-f(C + \Delta x) = f(C - \Delta x)$.

Lookup Tables for Symmetric Functions

Lookup tables are used to store the results of expensive calculations. These tables approximate a function by finding the appropriate entry in the table for a given input. Often, the input does not exactly match one that is stored. In this case, further approximation, such as interpolation or rounding, is needed. Further approximation causes numerical errors that can be reduced by creating a larger table. Thus, there is a tradeoff between the accuracy of the approximating lookup table and the memory efficiency of the table.

When approximating a symmetric function, you only need to store a region of the full image of the function. You can construct the remainder by applying the symmetry operations to this region. This process creates smaller, more memory-efficient lookup tables without sacrificing numerical performance.

Use a Lookup Table to Approximate a Superposition of Sinusoids

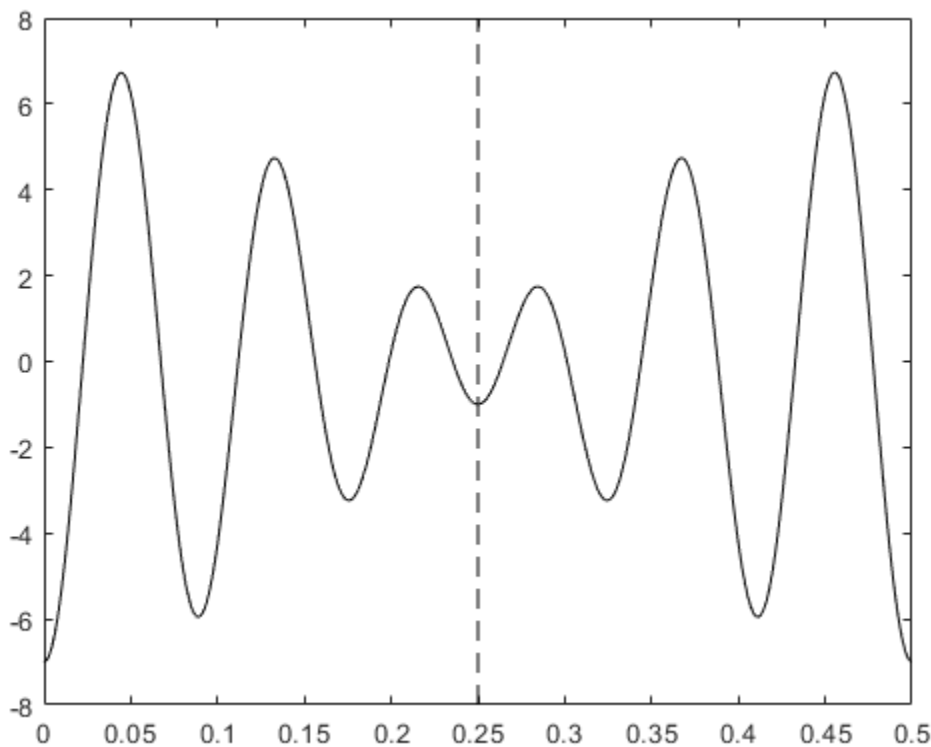
To demonstrate how to use symmetries to reduce the space in a lookup table, we will use a superposition of two sinusoid functions. This function should be viewed as a non-trivial model that illustrates this concept.

When two or more sinusoids are superimposed, they create a distinct interference pattern that is dependent upon the amplitude and relative phase of each sinusoid. If the ratio of the periods of the sinusoids is rational, then the superposition of the sinusoids is periodic as well, with a period equal to the least common multiple of each period. The interference can lead to a symmetric substructure within each period, though not always.

Consider the function $3 \times \sin(20 \times \pi \times (x - \frac{1}{8})) + 4 \times \cos(24 \times \pi \times (x - \frac{1}{8}))$. This function has a period of $\frac{1}{2}$, as well as an even axis of symmetry about the line $x = \frac{1}{4}$. You can exploit both of these values to create an efficient lookup table.

The plot below shows the first period of the function, with the axis of even symmetry shown as a dotted vertical line.

```
h = figure;
x = 0:0.001:0.5;
f = @(x) 3.*sin(20.*pi.*(x - 0.125)) + 4 .* cos(24.*pi.*(x-0.125));
p0 = plot(x, f(x), 'k-');
YLim = [-8, 8];
hold on;
a = h.Children(1);
a.YLim = YLim;
p1 = plot(a, [0.25 0.25], YLim, 'k--');
```



Consider the case where inputs range from 0 to 10, with a precision of approximately three decimal places. To model this, use an input numeric type with 14 bits, 10 of them fractional.

```
inputNt = numeric(0, 14, 10);
```

Use function approximation to generate a memory-efficient lookup table for any test function. Define the approximation problem by creating a `FunctionApproximation.Problem` object. Use the `solve` method to solve the optimization problem.

```
problem = FunctionApproximation.Problem(f);
problem.InputLowerBounds = 0;
problem.InputUpperBounds = 10;
problem.InputTypes = inputNt;
s1 = problem.solve();
fprintf("Lookup table uses %3.2f Kilobytes\n", s1.totalMemoryUsage ./ 2^10);
```

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|---------------|----------|------------|-----------------|--------------|----------------|
| 0 | 92 | 0 | 2 | 14 | 32 | EvenPow |
| 1 | 163900 | 1 | 5121 | 14 | 32 | EvenPow |
| 2 | 109276 | 0 | 3414 | 14 | 32 | EvenPow |
| 3 | 65596 | 0 | 2049 | 14 | 32 | EvenPow |
| 4 | 81980 | 0 | 2561 | 14 | 32 | EvenPow |
| 5 | 156 | 0 | 2 | 14 | 64 | EvenPow |
| 6 | 128 | 0 | 2 | 32 | 32 | EvenPow |
| 7 | 131164 | 0 | 2049 | 14 | 64 | EvenPow |
| 8 | 92 | 0 | 2 | 14 | 32 | EvenPow |
| 9 | 81980 | 0 | 2561 | 14 | 32 | EvenPow |
| 10 | 156 | 0 | 2 | 14 | 64 | EvenPow |
| 11 | 128 | 0 | 2 | 32 | 32 | EvenPow |
| 12 | 138000 | 0 | 3000 | 14 | 32 | Explicit |
| 13 | 150006 | 1 | 3261 | 14 | 32 | Explicit |
| 14 | 150006 | 0 | 3261 | 14 | 32 | Explicit |
| 15 | 150006 | 0 | 3261 | 14 | 32 | Explicit |
| 16 | 150006 | 0 | 3261 | 14 | 32 | Explicit |
| 17 | 150006 | 0 | 3261 | 14 | 32 | Explicit |

Best Solution

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|---------------|----------|------------|-----------------|--------------|----------------|
| 13 | 150006 | 1 | 3261 | 14 | 32 | Explicit |

```
Lookup table uses 146.49 Kilobytes
```

In this case, the lookup table is large.

Use Function Symmetries to Reduce Size of Lookup Table

Optimizing a lookup table over a smaller range of inputs results in a smaller lookup table, as fewer function values are stored in the table. For this example, only the range for x between 0 and $1/2$ needs be stored in the lookup table approximation. Adjusting the input upper and lower bounds and re-solving demonstrates the memory saved.

```
problem.InputLowerBounds = 0;
problem.InputUpperBounds = 0.25;
s2 = problem.solve();
sprintf("Lookup table uses %1.2f Kilobytes", s2.totalMemoryUsage ./ 2^10)
```

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|---------------|----------|------------|-----------------|--------------|----------------|
| 0 | 48 | 0 | 2 | 8 | 16 | EvenPow |
| 1 | 32 | 0 | 2 | 8 | 8 | EvenPow |
| 2 | 2080 | 1 | 129 | 8 | 16 | EvenPow |
| 3 | 1056 | 0 | 65 | 8 | 16 | EvenPow |

| | | | | | | |
|----|------|---|-----|----|----|---------|
| 4 | 80 | 0 | 2 | 8 | 32 | Even |
| 5 | 48 | 0 | 2 | 8 | 16 | EvenPow |
| 6 | 32 | 0 | 2 | 8 | 8 | EvenPow |
| 7 | 1056 | 0 | 65 | 8 | 16 | EvenPow |
| 8 | 80 | 0 | 2 | 8 | 32 | EvenPow |
| 9 | 2544 | 1 | 106 | 8 | 16 | Expli |
| 10 | 2448 | 0 | 102 | 8 | 16 | Expli |
| 11 | 2448 | 0 | 102 | 8 | 16 | Expli |
| 12 | 2712 | 1 | 113 | 8 | 16 | Expli |
| 13 | 144 | 0 | 2 | 8 | 64 | Even |
| 14 | 128 | 0 | 2 | 32 | 32 | Even |
| 15 | 144 | 0 | 2 | 8 | 64 | EvenPow |
| 16 | 128 | 0 | 2 | 32 | 32 | EvenPow |

Best Solution

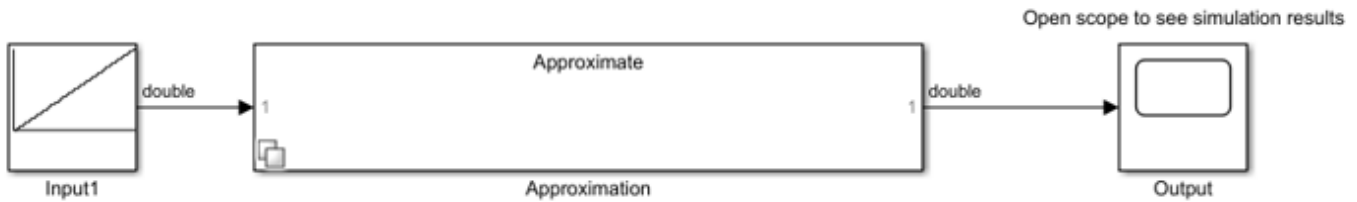
| ID | Memory (bits) | Feasible | Table Size | Breakpoints Wls | TableData WL | BreakpointSpec |
|----|---------------|----------|------------|-----------------|--------------|----------------|
| 2 | 2080 | 1 | 129 | 8 | 16 | Even |

```
ans =
    "Lookup table uses 2.03 Kilobytes"
```

Create Lookup Table Block

Use the approximate method to generate a Lookup Table block from the reduced-size lookup table in the previous example.

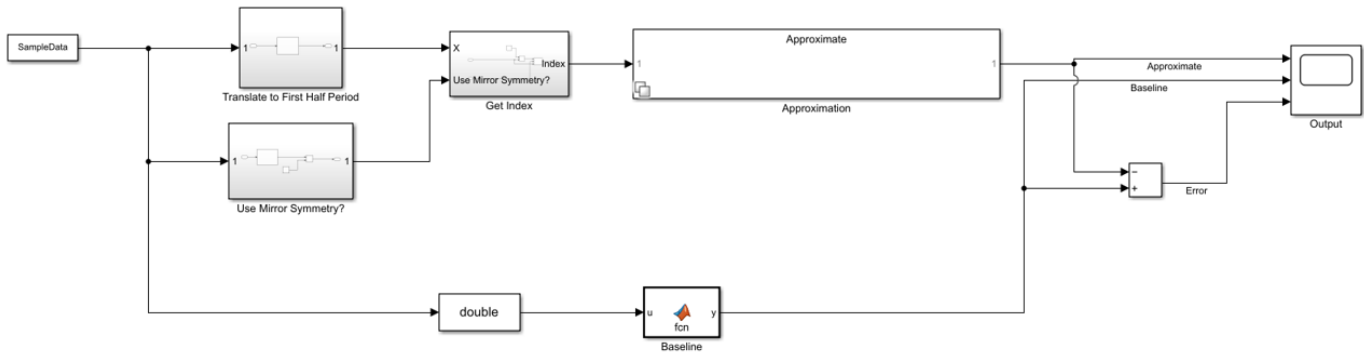
```
s2.approximate();
```



Use the Reduced-Size Lookup Table in a Simulink Model

This model demonstrates how to combine the symmetries of the example function with the lookup table to obtain an accurate approximation. The block diagram illustrates how to use symmetry operations in conjunction with the lookup table found in the previous example. The lookup table stored in the block named **Approximation** contains the function values for inputs between 0 and 0.25.

```
open_system('ModelWithApproximation')
```

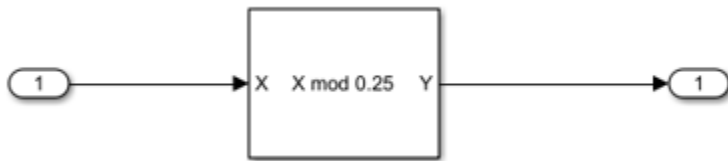


Copyright 2021 The MathWorks, Inc.

Translate to First Half Period

The subsystem named `Translate to First Half Period` demonstrates how to use the `Modulo by Constant` block to reduce the function input to the half open interval $[0, 0.25)$. Note that this block uses the value 0.25 instead of 0.5, because you use both the translational and mirror symmetries of the function for efficiency. Since this block is designed for embedded efficiency, it can reduce this calculation to a cast, as 0.25 is a power of 2.

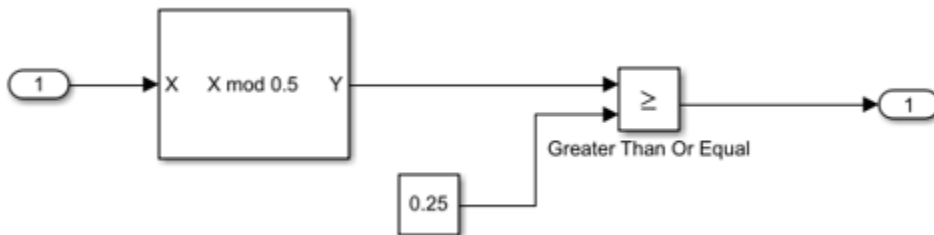
```
open_system('ModelWithApproximation/Translate to First Half Period')
```



Use Mirror Symmetry

To use the mirror symmetry, first check which half of a full period each input lies in. If $x \bmod 0.5 > 0.25$, reflect about the axis of symmetry in the first period. Use the `Modulo by Constant` block to perform this operation.

```
open_system('ModelWithApproximation/Use Mirror Symmetry?')
```

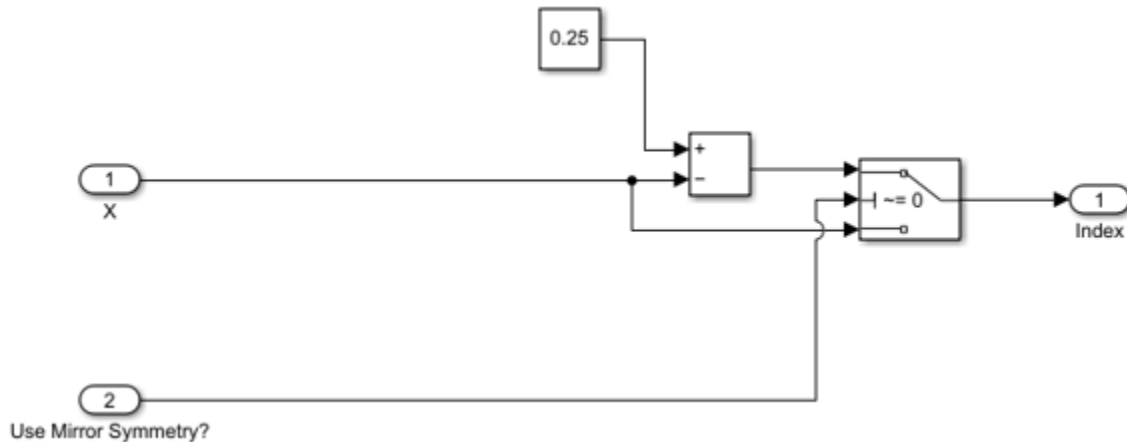


Get the Lookup Table Index

The subsystem named `Get Index` converts $X \bmod 0.25$ into a lookup table index. It takes in a Boolean argument to indicate whether you need to reflect about the axis of symmetry. If so,

$X \bmod 0.25$ is subtracted from 0.5 to get the index into the lookup table. Otherwise, $X \bmod 0.25$ is used to get the index.

```
open_system('ModelWithApproximation/Get Index');
```



Generate a Floating-Point Baseline

Use a MATLAB function block to generate a suitable floating-point baseline which to compare the lookup table approximation. The body of the function is shown here.

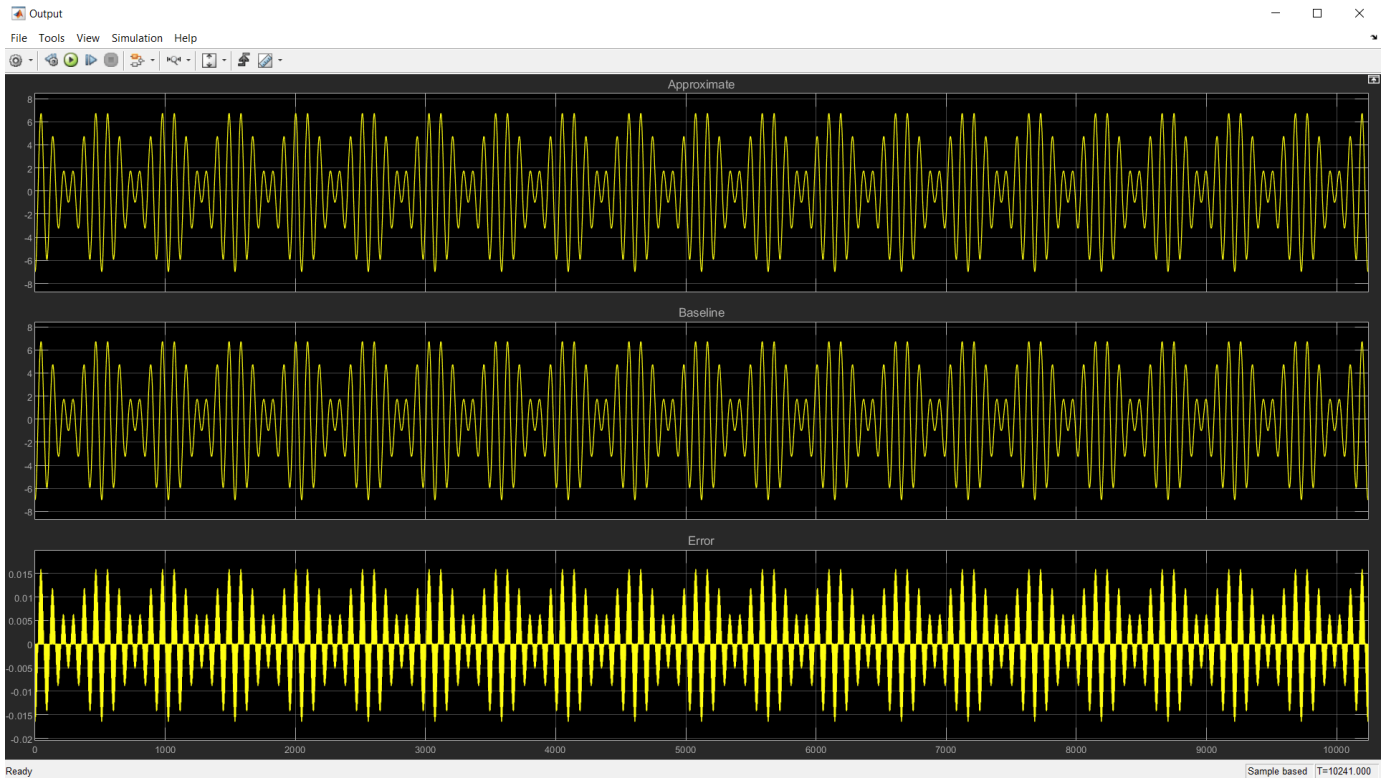
```
1 function y = fcn(u)
2
3 y = 3 .* sin(20.*pi.*(u-0.125)) + 4.*cos(24.*pi.*(u-0.125));
4
```

Simulate the System

Create data from 0 to 10 and simulate the model. This range will make it easy to see the repeating interference pattern exhibited by both the lookup table approximation as well as the baseline.

```
SampleData.signals.values = fi((0:pow2(-10):10)', inputNt);
SampleData.signals.dims = [1,1];
SampleData.time = (0:1:length(SampleData.signals.values)-1)';
sim('ModelWithApproximation');
```

The scope below compares the approximation and baseline. The error plot shows that the approximation using the compressed lookup table agrees well with the double-precision value of the example function.



See Also

`FunctionApproximation.Problem | Modulo by Constant`

Related Examples

- “Optimize Lookup Tables for Memory-Efficiency Programmatically” on page 41-19

Replace Fitted Curve with Optimized Lookup Table

This example shows how to approximate a fitted curve or surface and generate a lookup table.

In this example, you fit a surface to two-dimensional data then use the Lookup Table Optimizer command line interface to approximate the fitted curve with a lookup table.

Fit Surface to Data

Load the data set `sample_dataset`, which contains data that shows the relationship between engine speed, fuel rate, and torque.

```
load sample_dataset
```

Two variables are created in the workspace. The `inputs` variable is two-dimensional data. The `targets` variable is a column vector.

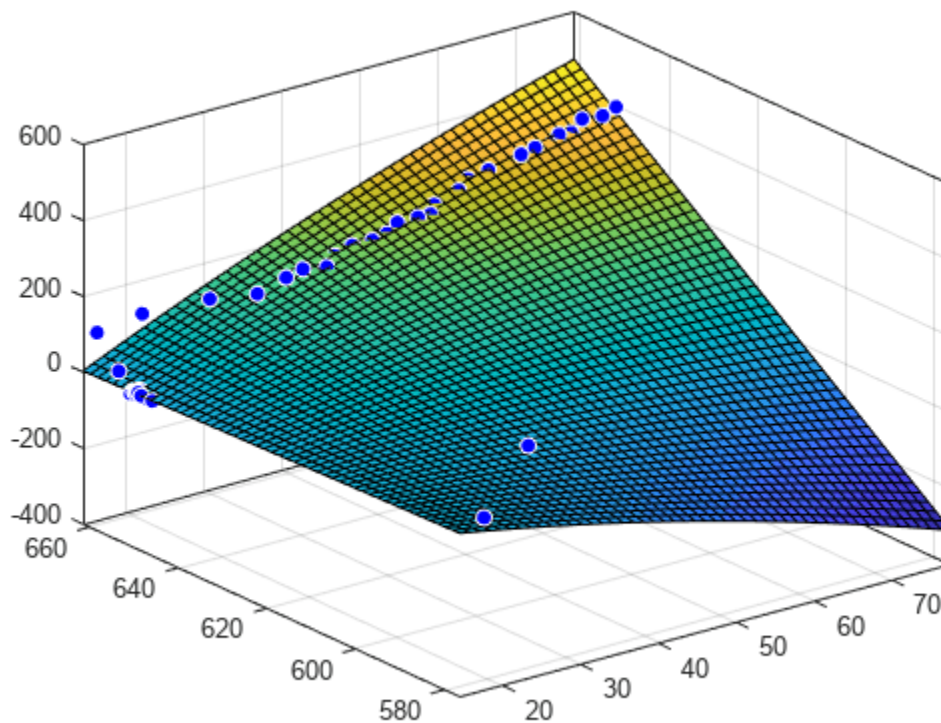
Use the `fit` function to fit a polynomial surface to these two variables. Specify the second-degree polynomial function to fit as `'poly21'`.

```
f = fit(inputs,targets,'poly21')
```

```
Linear model Poly21:  
f(x,y) = p00 + p10*x + p01*y + p20*x^2 + p11*x*y  
Coefficients (with 95% confidence bounds):  
p00 =      1539   (599.4, 2479)  
p10 =     -89.93  (-126.3, -53.61)  
p01 =     -2.532  (-3.983, -1.08)  
p20 =    -0.02819 (-0.04941, -0.006961)  
p11 =      0.1511 (0.09482, 0.2074)
```

Plot the fit and the data.

```
plot(f,inputs,targets)
```



Approximate Surface

Approximate the surface with an optimized lookup table. Create a `FunctionApproximation.Problem` object and specify the fitted surface curve f as the function to approximate. Set the lower and upper limits to be the range of the input values.

```
problem = FunctionApproximation.Problem ('f', 'InputLowerBounds', min(inputs), 'InputUpperBounds', max(inputs))
```

```
problem =
```

```
1x1 FunctionApproximation.Problem with properties:
```

```
FunctionToApproximate: [1x1 sfit]
NumberOfInputs: 2
InputTypes: ["numerictype('double')" "numerictype('double')"]
InputLowerBounds: [14.4000 576.2000]
InputUpperBounds: [77.5000 661.5000]
OutputType: "numerictype('double')"
Options: [1x1 FunctionApproximation.Options]
```

You can edit the `FunctionApproximation.Options` object to specify additional constraints to use in the lookup table optimization process. Specify the word lengths and the maximum time to find a solution.

```
problem.Options.WordLengths = [8,32];
problem.Options.MaxTime = 240;
```

Use the solve method to create an optimized lookup table approximation. The solve method returns a FunctionApproximation.LUTSolution object.

```
solution = solve(problem)
```

ExplicitValues specification is only available for a function with 1 input dimension. Trying only

Searching for fixed-point solutions.

| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|---------------|----------|------------|-----------------|--------------|----------------|
| 0 | 160 | 0 | [2 2] | [8 8] | 32 | Ev |
| 1 | 672 | 0 | [4 5] | [8 8] | 32 | Ev |
| 2 | 512 | 0 | [3 5] | [8 8] | 32 | Ev |
| 3 | 544 | 0 | [4 4] | [8 8] | 32 | Ev |
| 4 | 416 | 0 | [3 4] | [8 8] | 32 | Ev |
| 5 | 1824 | 0 | [7 8] | [8 8] | 32 | Ev |
| 6 | 1568 | 0 | [6 8] | [8 8] | 32 | Ev |
| 7 | 5024 | 0 | [13 12] | [8 8] | 32 | Ev |
| 8 | 4640 | 0 | [12 12] | [8 8] | 32 | Ev |
| 9 | 19168 | 0 | [26 23] | [8 8] | 32 | Ev |
| 10 | 16224 | 0 | [22 23] | [8 8] | 32 | Ev |
| 11 | 31680 | 0 | [43 23] | [8 8] | 32 | Ev |
| 12 | 94240 | 1 | [128 23] | [8 8] | 32 | Ev |
| 13 | 47136 | 0 | [64 23] | [8 8] | 32 | Ev |
| 14 | 2080 | 0 | [8 8] | [8 8] | 32 | Ev |
| 15 | 5792 | 0 | [15 12] | [8 8] | 32 | Ev |
| 16 | 23584 | 0 | [32 23] | [8 8] | 32 | Ev |
| 17 | 4256 | 0 | [11 12] | [8 8] | 32 | Ev |
| 18 | 3872 | 0 | [10 12] | [8 8] | 32 | Ev |
| 19 | 14016 | 0 | [19 23] | [8 8] | 32 | Ev |
| 20 | 49184 | 1 | [128 12] | [8 8] | 32 | Ev |
| 21 | 24608 | 0 | [128 6] | [8 8] | 32 | Ev |
| 22 | 32800 | 0 | [128 8] | [8 8] | 32 | Ev |
| 23 | 256 | 0 | [2 2] | [32 32] | 32 | Ev |
| 24 | 768 | 0 | [4 5] | [32 32] | 32 | Ev |
| 25 | 608 | 0 | [3 5] | [32 32] | 32 | Ev |
| 26 | 640 | 0 | [4 4] | [32 32] | 32 | Ev |
| 27 | 512 | 0 | [3 4] | [32 32] | 32 | Ev |
| 28 | 2144 | 0 | [7 9] | [32 32] | 32 | Ev |
| 29 | 1856 | 0 | [6 9] | [32 32] | 32 | Ev |
| 30 | 1920 | 0 | [7 8] | [32 32] | 32 | Ev |
| 31 | 1664 | 0 | [6 8] | [32 32] | 32 | Ev |
| 32 | 7200 | 0 | [13 17] | [32 32] | 32 | Ev |
| 33 | 6656 | 0 | [12 17] | [32 32] | 32 | Ev |
| 34 | 6784 | 0 | [13 16] | [32 32] | 32 | Ev |
| 35 | 6272 | 0 | [12 16] | [32 32] | 32 | Ev |
| 36 | 26528 | 0 | [25 33] | [32 32] | 32 | Ev |
| 37 | 25472 | 0 | [24 33] | [32 32] | 32 | Ev |
| 38 | 25728 | 0 | [25 32] | [32 32] | 32 | Ev |
| 39 | 24704 | 0 | [24 32] | [32 32] | 32 | Ev |
| 40 | 14240 | 0 | [21 21] | [32 32] | 32 | Ev |
| 41 | 28928 | 0 | [30 30] | [32 32] | 32 | Ev |
| 42 | 37120 | 0 | [34 34] | [32 32] | 32 | Ev |
| 43 | 41600 | 0 | [36 36] | [32 32] | 32 | Ev |
| 44 | 43936 | 0 | [37 37] | [32 32] | 32 | Ev |
| 45 | 46336 | 0 | [38 38] | [32 32] | 32 | Ev |
| 46 | 2432 | 0 | [8 9] | [32 32] | 32 | Ev |
| 47 | 2176 | 0 | [8 8] | [32 32] | 32 | Ev |

```
| 48 |          8288 |          0 | [15 17] | [32 32] |          32 | Ev
| 49 |          7744 |          0 | [14 17] | [32 32] |          32 | Ev
| 50 |          7808 |          0 | [15 16] | [32 32] |          32 | Ev
Searching for floating-point solutions.
```

Time limit for finding a solution has been reached.

Best Solution

```
| ID | Memory (bits) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec
| 20 |          49184 |          1 | [128 12] | [8 8] |          32 | Ev
```

```
solution =
  1x1 FunctionApproximation.LUTSolution with properties:
```

```
    ID: 20
    Feasible: "true"
```

View the lookup table data.

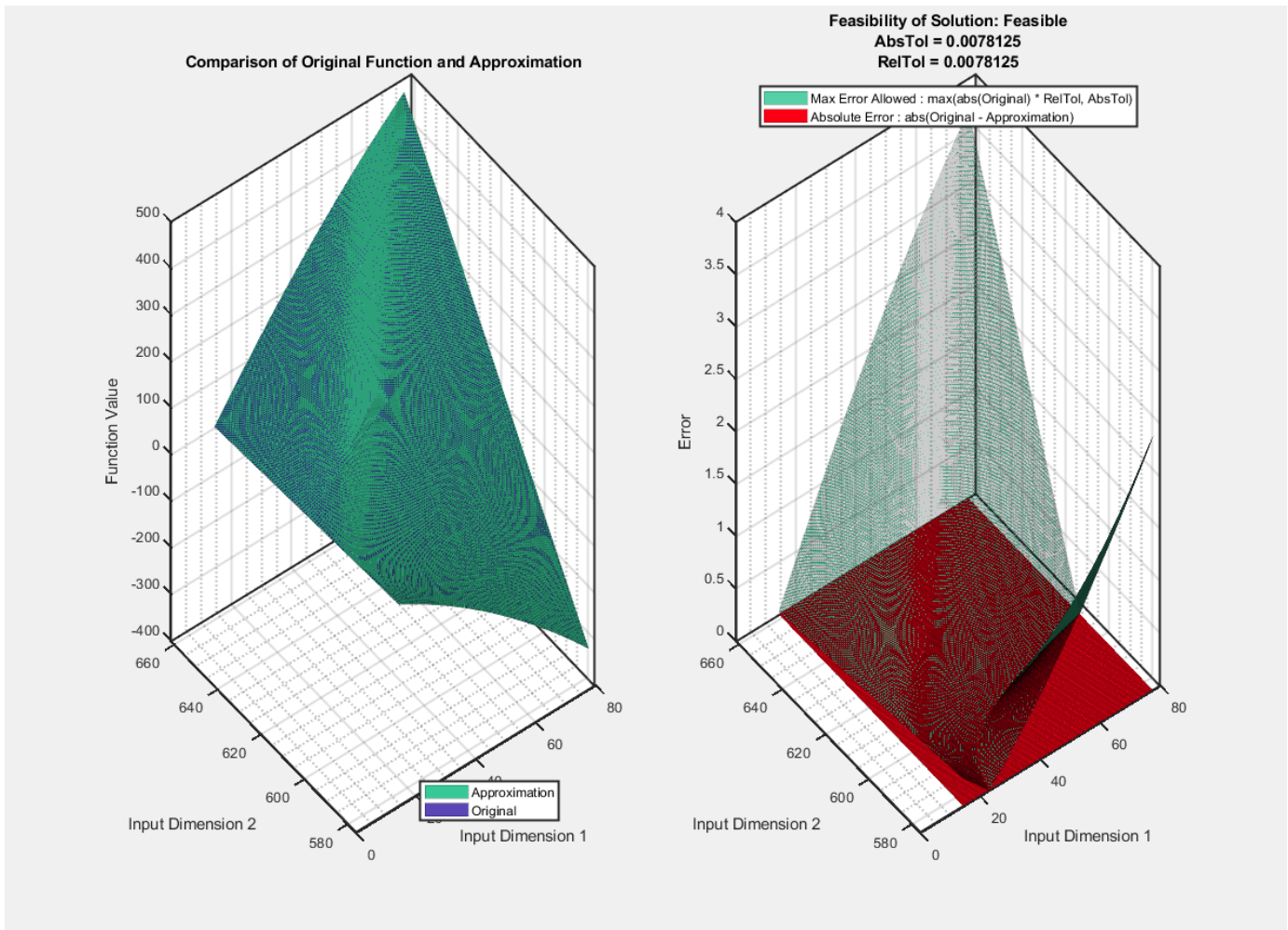
```
solution.TableData
```

```
ans = struct with fields:
    BreakpointValues: {1x2 cell}
    BreakpointDataTypes: [1x2 embedded.numericitype]
    TableValues: [128x12 double]
    TableDataType: [1x1 embedded.numericitype]
    IsEvenSpacing: 1
    Interpolation: Linear
```

Compare Lookup Table Approximation to Original Function

Compare the numerical behavior of the original surface fit function with the optimized lookup table approximation.

```
compare(solution)
```

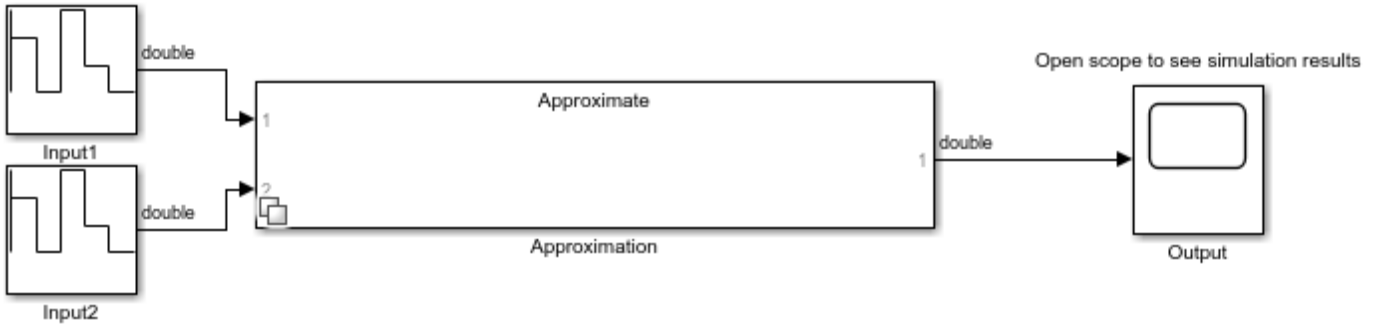



```
ans=1x2 struct array with fields:
    Breakpoints
    Original
    Approximate
```

Generate Subsystem

Use the approximate method to generate a Simulink™ subsystem that contains the lookup table approximation.

```
approximate(solution)
```



See Also

`fit` | `FunctionApproximation.Problem` | `FunctionApproximation.Options`

Related Examples

- “Optimize Lookup Tables for Memory-Efficiency Programmatically” on page 41-19
- “Parametric Fitting” (Curve Fitting Toolbox)

Automatic Data Typing

- “Choosing a Range Collection Method” on page 42-2
- “Best Practices for Fixed-Point Conversion Workflow” on page 42-5
- “Models That Might Cause Data Type Propagation Errors” on page 42-7
- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9
- “Set Up the Model” on page 42-13
- “Prepare System for Conversion” on page 42-14
- “Specify Behavioral Constraints” on page 42-18
- “Collect Ranges” on page 42-20
- “Convert Data Types” on page 42-22
- “Examine Results to Resolve Conflicts” on page 42-26
- “Verify New Settings” on page 42-30
- “Explore Additional Data Types” on page 42-34
- “Restore Model to Original State” on page 42-36
- “Get Proposals for Results with Inherited Types” on page 42-37
- “Rescale a Fixed-Point Model” on page 42-39
- “How the Fixed-Point Tool Proposes Data Types” on page 42-48
- “How Hardware Implementation Settings Affect Data Type Proposals” on page 42-50
- “Propose Data Types For Merged Simulation Ranges” on page 42-54
- “View Simulation Results” on page 42-57
- “Fixed-Point Instrumentation and Data Type Override” on page 42-61
- “Model Configuration Changes Made During Data Type Optimization” on page 42-63

Choosing a Range Collection Method

The Fixed-Point Tool automates the task of specifying fixed-point data types in a Simulink model. You can choose to use an iterative fixed-point conversion process, also known as autoscaling, or you can optimize data types in your model using `fxpopt`. The Fixed-Point Tool also lets you explore the numerical behavior of floating-point vs fixed-point data types in your model.

The tool collects range data for model objects from design minimum and maximum values that objects explicitly specify, from logged minimum and maximum values that occur during simulation, or from the minimum and maximum values derived using static range analysis.

| Method | Advantages | Disadvantages |
|---|---|---|
| Using simulation minimum and maximum values | <ul style="list-style-type: none"> Useful if you know the inputs to use for the model. You do not need to specify any design range information. | <ul style="list-style-type: none"> Not always feasible to collect the full simulation range. Simulation might take a very long time. |
| Using design minimum and maximum values | You can use this method if the model contains blocks that range analysis does not support. However, if possible, use simulation data to propose data types. | <ul style="list-style-type: none"> The design range is often available only on some input and output signals. You can propose data types only for signals with specified design minimum and maximum values. |
| Using derived minimum and maximum values | You do not have to simulate multiple times to ensure that simulation data covers the full intended operating range. | <ul style="list-style-type: none"> Derivation might take a very long time. |

In the Fixed-Point Tool, you can choose between three range collection modes:

- **Simulation ranges** - Collect ranges through simulation. To collect and merge the ranges of multiple simulation runs, you can specify simulation inputs.
- **Derived ranges** - Collect ranges through a static analysis that derives the ranges, also known as *range analysis*.
- **Simulation with Range Analysis** - Collect ranges through simulation and derived range analysis and combine the results.

| Feature | Simulation Ranges | Derived Ranges | Simulation with Range Analysis |
|-----------------------|---|---|---|
| Range coverage | Proposed data types are based on simulation ranges. The proposals provided by the Fixed-Point Tool are as good as the test bench provided. Data type proposals are based on collected minimum and maximum values. | Static range analysis typically delivers a more conservative data type proposal. Data type proposals are based on collected minimum and maximum values. | Proposed data types are based on the union of simulation ranges and derived ranges. Data type proposals are based on collected minimum and maximum values. This option provides the most comprehensive range information. |

| Feature | Simulation Ranges | Derived Ranges | Simulation with Range Analysis |
|---|---|--|--|
| Simulation inputs | Comprehensive set of input signals that exercise the full range of your design. This allows you to collect and merge ranges from multiple simulation input cases. | Ranges reported from derivation are based only on design ranges specified in the model. Simulation inputs are not used to derive ranges. | Ranges are based on the combination of merged simulation ranges and ranges derived from design ranges specified in the model. |
| Design ranges | Simulation ranges are verified against design range specification and violations are reported in the Diagnostic Viewer. | Design ranges must be specified on the model. Data type proposals are based on collected minimum and maximum values. | Simulation ranges are verified against design range specification. To derive ranges, design ranges must be specified on the model. |
| Supported Features | All model objects are supported for instrumentation and range collection. | Range analysis supports a subset of model objects. For more information, see “Unsupported Simulink Software Features” on page 43-27. | Range analysis supports a subset of model objects. For more information, see “Unsupported Simulink Software Features” on page 43-27. |
| Modeling constructs | Ranges always converge during simulation. | Some modeling constructs, such as feedback loops, may require more design range information before converging. | Simulation ranges always converge. Some modeling constructs, such as feedback loops, may require more design range information before derived ranges converge. |
| Tunable parameters with known ranges | You must exercise the full tunable range using simulation inputs. | Design ranges of tunable parameters are reported. | Design ranges of tunable parameters are reported. You can additionally exercise the tunable range using simulation inputs. |

| Feature | Simulation Ranges | Derived Ranges | Simulation with Range Analysis |
|------------------------|---|--|---|
| Simulation mode | Instrumentation data is only collected during Normal mode. No instrumentation data is collected while a model is running in accelerator or rapid accelerator mode. If you know that simulation will take a long time, you may want to derive ranges for your model. | Simulation mode has no affect on range analysis. | Instrumentation data is only collected during Normal mode. No instrumentation data is collected while a model is running in accelerator or rapid accelerator mode. If you know that simulation will take a long time, you may want to derive ranges for your model. |

Based on collected range information, the tool proposes fixed-point data types that maximize precision and cover the range. The Fixed-Point Tool allows you to review the data type proposals and then apply them selectively to objects in your model.

See Also

More About

- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9
- “How Range Analysis Works” on page 43-2

Best Practices for Fixed-Point Conversion Workflow

Using the Fixed-Point Tool, you can prepare a model or subsystem for conversion from floating-point to an equivalent fixed-point representation. The following are modeling best practices for converting a model to fixed point.

Enable Signal Logging

To compare the behavior before and after conversion, enable signal logging for signals of interest in the system under design.

You can specify absolute, relative, and time tolerances for signals in your model that have signal logging enabled. After you simulate with embedded types, the Workflow Browser displays whether the embedded run meets the specified signal tolerances compared to the baseline run created during range collection. You can view the comparison plots in the Simulation Data Inspector.

Back Up Your Simulink Model

Before using the Fixed-Point Tool, back up your Simulink model and associated workspace variables. Backing up your model can provide a baseline for testing and validation.

The Fixed-Point Tool automatically creates a back up of your original model during the **Prepare** stage of the conversion. To restore your model to this state, click the **Restore Original Model** button.

Convert Individual Subsystems

Convert individual subsystems in your model one at a time. This practice facilitates debugging by isolating the source of fixed-point issues.

Do Not Use “Save as” on Referenced Models and MATLAB Function blocks

During the fixed-point conversion process using the Fixed-Point Tool, do not use the “Save as” option to save referenced models or MATLAB Function blocks with a different name. If you do, you might lose existing results for the original model.

Use Lock Output Data Type Setting

You can prevent the Fixed-Point Tool from replacing the current data type. Use the **Lock output data type setting against changes by the fixed-point tools** parameter that is available on many blocks. The default setting allows for replacement. Use this setting when:

- You already know the fixed-point data types that you want to use for a particular block.

For example, the block is modeling a real-world component. Set up the block to allow for known hardware limitations, such as restricting outputs to integer values.

Explicitly specify the output data type of the block and select **Lock output data type setting against changes by the fixed-point tools**.

- You are debugging a model and know that a particular block accepts only certain input signal data types.

Explicitly specify the output data type of upstream blocks and select **Lock output data type setting against changes by the fixed-point tools**.

Save Simulink Signal Objects

If your model contains Simulink signal objects and you accept proposed data types, the Fixed-Point Tool automatically applies the changes to the signal objects. However, the Fixed-Point Tool does not automatically save changes that it makes to Simulink signal objects. To preserve changes, before closing your model, save the Simulink signal objects in your workspace and model.

Do Not Use `clear all`

`clear all` is not supported by fixed-point conversion workflows. Do not use `clear all` in initialization functions (`InitFcn`), or at the MATLAB Command Window when using the Fixed-Point Tool.

See Also

Related Examples

- “Convert Floating-Point Model to Fixed Point” on page 40-2

More About

- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9

Models That Might Cause Data Type Propagation Errors

When the Fixed-Point Tool proposes changes to the data types in your model in the **Iterative Fixed-Point Conversion** workflow, it alerts you to potential issues. If the Fixed-Point Tool alerts you to data type errors, you must diagnose the errors and fix the problems. For more information, see “Examine Results to Resolve Conflicts” on page 42-26.

The Fixed-Point Tool does not detect all potential data type issues. If the tool does not report any issues for your model, it is still possible to experience subsequent data type propagation errors. Before you use the Fixed-Point Tool, back up your model to ensure that you can recover your original data type settings. For more information, see “Best Practices for Fixed-Point Conversion Workflow” on page 42-5.

The following model components are likely to cause data type propagation issues.

| Model Uses... | Fixed-Point Tool Behavior | Data Type Propagation Issue |
|--------------------------------|---|---|
| Simulink parameter objects | The Fixed-Point Tool is not able to detect when a parameter object must be integer only, such as when using a parameter object as a variable for dimensions, variant control, or a Boolean value. | Fixed-Point Tool might propose data types that are inconsistent with the data types for the parameter object or generate proposals that cause overflows. |
| User-defined S-functions | Cannot detect the operation of user-defined S-functions. | <ul style="list-style-type: none"> The user-defined S-function accepts only certain input data types. The Fixed-Point Tool cannot detect this requirement and proposes a different data type upstream of the S-function. Update diagram fails on the model due to data type mismatch errors. The user-defined S-function specifies certain output data types. The Fixed-Point Tool is not aware of this requirement and does not use it for automatic data typing. Therefore, the tool might propose data types that are inconsistent with the data types for the S-function or generate proposals that cause overflows. |
| User-defined masked subsystems | Has no knowledge of the masked subsystem workspace and cannot take this subsystem into account when proposing data types. | Fixed-Point Tool might propose data types that are inconsistent with the requirements of the masked subsystem, particularly if the subsystem uses mask initialization. The proposed data types might cause data type mismatch errors or overflows. |

| Model Uses... | Fixed-Point Tool Behavior | Data Type Propagation Issue |
|----------------------|---|---|
| Linked subsystems | Does not include linked subsystems when proposing data types. | Data type mismatch errors might occur at the linked subsystem boundaries. |

See Also

More About

- “Data Type Propagation Errors After Applying Proposed Data Types” on page 49-25

Iterative Fixed-Point Conversion Using the Fixed-Point Tool

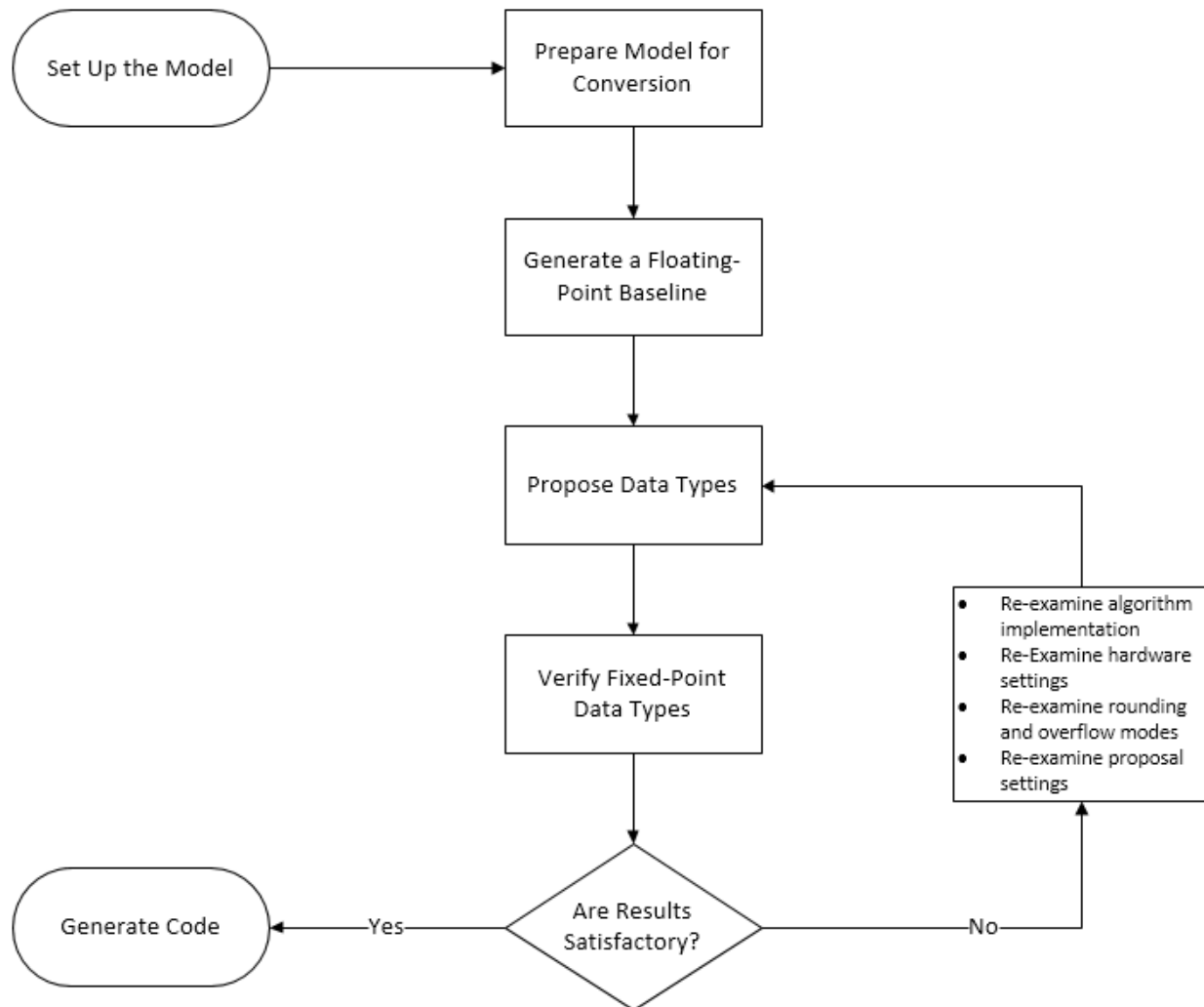
The Fixed-Point Tool is a user interface that automates the task of specifying fixed-point data types in a Simulink model. The tool collects range data for model objects. The range data comes from:

- Design minimum and maximum values that objects specify explicitly on the block
- Logged minimum and maximum values that occur during simulation
- Minimum and maximum values derived using range analysis

Based on these values, the tool in the `Iterative Fixed-Point Conversion` workflow proposes fixed-point data types that maximize precision and cover the range. You can then review the data type proposals and apply them selectively to objects in your model. This process is also known as autoscaling. Using the iterative fixed-point conversion workflow in the Fixed-Point Tool you can:

- Derive range information based on specified design ranges. See “How Range Analysis Works” on page 43-2.
- Propose and apply data types based on simulation data.
- Propose and apply data types based on derived ranges.
- Propose and apply data types based on simulation data from multiple runs. See “Propose Data Types For Merged Simulation Ranges” on page 42-54.
- Propose and apply data types based on simulation data and derived ranges.
- Debug fixed-point models.

Workflow for Automatic Data Typing



The iterative fixed-point conversion workflow for automatic data typing consists of four main stages.

1 “Prepare System for Conversion” on page 42-14

Before you begin conversion, set up the model in Simulink. Then select the system to convert to fixed point. The Fixed-Point Tool will propose data types for the objects in the specified system.

Select whether to collect ranges through simulation, derived range analysis, or simulation with range analysis. You can specify multiple simulation scenarios using a `Simulink.SimulationInput` object. Specify signal tolerances to use to verify the behavior of the converted system.

Automatically prepare the system under design for conversion by clicking the **Prepare** button in the Fixed-Point Tool toolbar. The Fixed-Point Tool analyzes your model and makes configuration recommendations for autoscaling.

2 “Collect Ranges” on page 42-20

Run the simulation or the derivation. When the simulation or derivation is complete, you can examine the ranges of objects in your model using the histograms in the **Visualization of Simulation Data** pane.

3 “Convert Data Types” on page 42-22

The Fixed-Point Tool proposes data types based on the ranges collected in stage two. You can edit the default word length and other proposal settings in the **Settings** menu. To generate proposals, click **Propose Data Types**. If you are satisfied with the proposals, click **Apply Data Types**.

4 “Verify New Settings” on page 42-30

Simulate your model using the newly applied fixed-point data types to examine the behavior of the fixed-point model. You can compare the floating point and fixed-point behavior using the Simulation Data Inspector.

5 “Explore Additional Data Types” on page 42-34

After verification, if you determine that the behavior of the system is not acceptable, you can iterate through the conversion and verification steps until you settle on a design that satisfies your system requirements.

The screenshot displays the Simulink software interface for automatic data typing. The top toolbar contains several key buttons: 'Propose Data Types' (3a), 'Apply Data Types' (3a), and 'Simulate with Embedded Types' (4a). A dropdown menu for 'Simulate with Embedded Types' is set to 'EmbeddedRun'. On the right, the 'RESULT DETAILS' panel shows the 'fxpdemo_feedback/Controller/Numerator Terms : Output' with a table of properties and their specified data types. Below this, 'Range Information' is provided for the simulation. At the bottom, a histogram plot titled 'Histograms of all results in the model' shows the distribution of data values across different categories: Overflows (red), Representable (grey), In-Range (blue), and Underflows (yellow). A legend in the bottom right of the plot identifies these categories. Orange callout boxes (1a, 1b, 2, 3a, 3a, 4a, 4b, 5) highlight specific UI elements and workflow steps.

| Name | CompiledDT | SpecifiedDT | ProposedDT | ActualDT | SimMin | SimMax |
|-------------------|----------------|-----------------------|------------|----------|------------------|------------------|
| Combine Terms... | fixdt(1,32,28) | Inherit: Inherit v... | | | -6.4669937491... | 4.33369296789... |
| Combine Terms... | fixdt(1,32,28) | fixdt(1,32,28) | | | -2.4210555553... | 4.33369296789... |
| Denominator Te... | fixdt(1,32,27) | fixdt(1,32,27) | | | -8.5294544696... | 5.40470331907... |
| Denominator Te... | fixdt(1,32,28) | fixdt(1,32,28) | | | -6.4669937491... | 3.48770594596... |
| Denominator Te... | fixdt(1,32,27) | fixdt(1,32,27) | | | -8.5294544696... | 5.40470331907... |
| Down Cast | fixdt(1,16,12) | fixdt(1,16,12) | | | -2.421142578125 | 4.33349609375 |
| Numerator Ter... | fixdt(1,32,28) | fixdt(1,32,28) | | | -5.67724609375 | 5.700439453125 |
| Numerator Ter... | fixdt(1,32,29) | fixdt(1,32,29) | | | -3.3876037597... | 3.5439453125 |
| Numerator Ter... | fixdt(1,32,28) | fixdt(1,32,28) | | | -5.67724609375 | 5.700439453125 |
| Up Cast | fixdt(1,16,12) | fixdt(1,16,12) | | | -2 | 4 |

| Property | Specified Data Type |
|-----------|-----------------------|
| Data Type | fixdt(1,32,29) |
| Minimum | -4 |
| Maximum | 3.999999998137355 |
| Precision | 1.862645149230957e-09 |

| Property | Minimum | Maximum |
|------------|------------------|--------------|
| Simulation | -3.3876037597... | 3.5439453125 |

See Also

Related Examples

- “Rescale a Fixed-Point Model” on page 42-39

Set Up the Model

Before using the Fixed-Point Tool to generate data type proposals for your model, set up your model in Simulink.

- 1 If you are using design minimum and maximum range information, add this information to the blocks. To autoscale using derived data, you **must** specify design minimum and maximum values on at least the model inputs. The range analysis tries to narrow the derived range by using all the specified design ranges in the model. The more design range information you specify, the more likely the range analysis is to succeed. As the analysis is performed, it derives new range information for the model and then attempts to use this new information together with the specified ranges to derive ranges for the remaining objects in the model. For this reason, the analysis results might depend on block priorities because these priorities determine the order in which the software analyzes the blocks.

You specify a design range for model objects using parameters such as **Output minimum** and **Output maximum**. For a list of blocks in which you can specify these values, see “Blocks That Allow Signal Range Specification”.

- 2 Enable signal logging.

To view simulation results using the Simulation Data Inspector, you must enable signal logging for the system you want to convert to fixed point. You can choose to plot results using the Simulation Data Inspector only for signals that have signal logging enabled.

a In the Simulink Editor, select one or more signals.

b In the **Signal** tab of the Simulink Editor, click **Log Signals**.

- 3 You can choose to lock some blocks against automatic data typing by selecting the block's **Lock output data type setting against changes by the fixed-point tools** parameter. If you select this parameter, the tool does not propose data types for the block.
- 4 Update the diagram to perform parameter range checking for all blocks in the model.

If updating the diagram fails, use the error messages to fix the errors in your model. After fixing the errors, update the diagram again. If you cannot fix the errors, restore your backup model.

To learn about the next step in the conversion process, see “Prepare System for Conversion” on page 42-14.

See Also

More About

- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9

Prepare System for Conversion

Before using the Fixed-Point Tool to generate data type proposals for your model, set up your model in Simulink and prepare the system for conversion.

Set Up the Model

Specify Design Ranges

If you are using design minimum and maximum range information, add this information to the blocks. To autoscale using derived data, you **must** specify design minimum and maximum values on at least the model inputs. The range analysis tries to narrow the derived range by using all the specified design ranges in the model. The more design range information you specify, the more likely the range analysis is to succeed. As the analysis is performed, it derives new range information for the model and then attempts to use this new information together with the specified ranges to derive ranges for the remaining objects in the model. For this reason, the analysis results might depend on block priorities because these priorities determine the order in which the software analyzes the blocks.

You specify a design range for model objects using parameters such as **Output minimum** and **Output maximum**. For a list of blocks in which you can specify these values, see “Blocks That Allow Signal Range Specification”.

Enable Signal Logging

To view simulation results using the Simulation Data Inspector, you must enable signal logging for the system you want to convert to fixed point. You can choose to plot results using the Simulation Data Inspector only for signals that have signal logging enabled.

- 1 In the Simulink Editor, select one or more signals.
- 2 In the **Signal** tab of the Simulink Editor, click **Log Signals**.

Use Lock Output Data Type Setting

You can choose to lock some blocks against automatic data typing by selecting the block's **Lock output data type setting against changes by the fixed-point tools** parameter. If you select this parameter, the tool does not propose data types for the block.

Update the Diagram

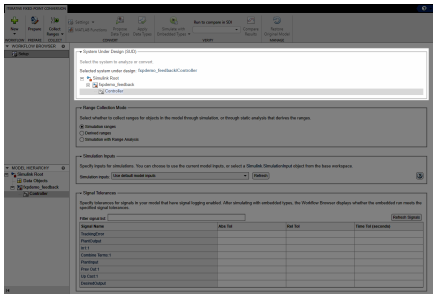
Update the diagram to perform parameter range checking for all blocks in the model.

If updating the diagram fails, use the error messages to fix the errors in your model. After fixing the errors, update the diagram again. If you cannot fix the errors, restore your backup model.

Select the System Under Design

To open the Fixed-Point Tool, in your model, in the **Apps** gallery, select **Fixed-Point Tool**. Alternatively, use the `fxptdlg` function.

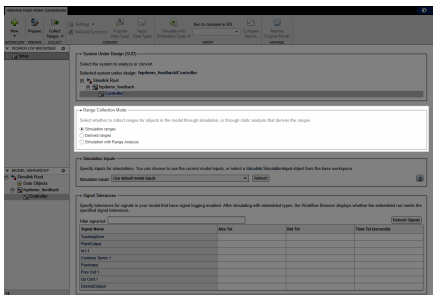
In the Fixed-Point Tool, click **New** and select **Iterative Fixed-Point Conversion**.



Select the system or subsystem you want to convert to fixed point. Convert individual subsystems in your model one at a time. This practice facilitates debugging by isolating the source of numerical issues.

In the main working area, under **System Under Design (SUD)**, use the drop-down menu to select the system or subsystem you want to convert.

Set Range Collection Method



You can collect ranges through simulation, derived range analysis, or by using simulation combined with derived range analysis. Using simulation-based range collection, the Fixed-Point Tool can be configured to perform a global override of the fixed-point data types with double-precision or single-precision data types, thereby avoiding quantization effects. This setting provides a floating-point benchmark that represents the ideal output. You can also collect benchmark range data using the current data type override set on the model.

If you collect ranges through simulation, you can choose to specify additional simulation inputs. During the range collection simulation, the Fixed-Point Tool captures the minimum and maximum values from each specified simulation scenario. For more information, see “Specify Simulation Input” on page 42-16.

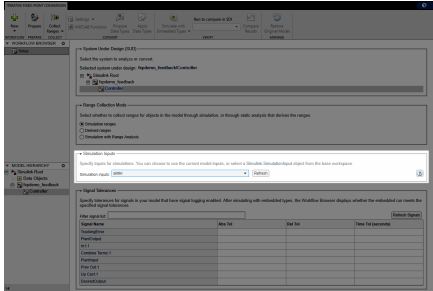
Using derived range analysis, the Fixed-Point Tool uses design ranges specified on blocks to analyze and derive static ranges for other objects in your model. The tool uses all design range information specified on the model to derive ranges for objects in the system under design. If you choose to collect ranges for objects in your model through derived range analysis, you do not need to simulate the model. However, to compare floating-point and fixed-point behavior using the Simulation Data Inspector, simulation is required.

Using simulation with range analysis, the Fixed-Point Tool uses the union of the ranges collected through simulation and derived range analysis.

Under **Range Collection Mode**, select the method that you want to use to collect ranges. The Fixed-Point Tool uses these collected ranges to later generate data type proposals.

For more information on deciding which method of range collection is right for your application, see “Choosing a Range Collection Method” on page 42-2.

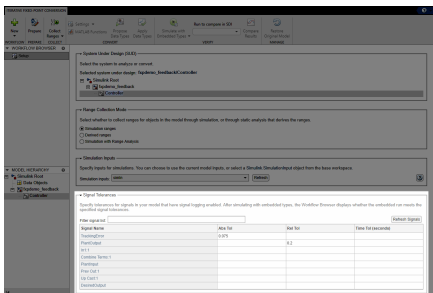
Specify Simulation Input



If you choose to collect ranges through simulation, you must specify the simulation input for your system. Under **Simulation Inputs**, select whether to use the default model inputs to simulate the model for range collection, or select a `Simulink.SimulationInput` object from the base workspace to specify one or more simulation scenarios.

If the `SimulationInput` object that you select contains more than one simulation scenario, the Fixed-Point Tool proposes data types based on the merged ranges from all simulation scenarios. Because the proposals provided by the Fixed-Point Tool are as good as the test bench provided, a comprehensive set of input signals that exercise the full range of your design will result in more accurate data type proposals for your system. For an example, see “Propose Data Types For Merged Simulation Ranges” on page 42-54.

Edit Signal Tolerances



You can specify absolute, relative, and time tolerances for signals in your model that have signal logging enabled. After converting your system, when you simulate the embedded run, the **Workflow Browser** displays whether the embedded run meets the specified signal tolerances compared to the baseline run established during range collection. You can view the comparison plots in the Simulation Data Inspector.

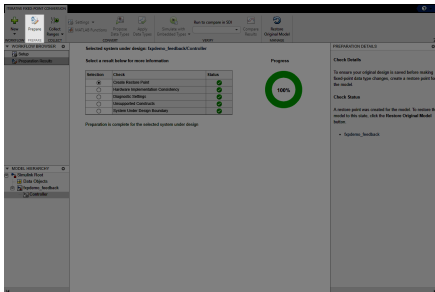
Specify signal tolerances in the table under **Signal Tolerances**. The table contains all signals in the model with signal logging enabled. In the boxes to the right of the signal for which you want to register a tolerance, enter the tolerances for the signal. You can specify any of the following types of tolerances.

- **Abs Tol** - Absolute value of the maximum acceptable difference between the original signal, and the signal in the converted design.

- **Rel Tol** – Maximum relative difference, specified as a percentage, between the original output, and the output of the new design. For example, a value of $1e-2$ indicates a maximum difference of one percent between the original signal values, and the signal values of the converted design.
- **Time Tol (seconds)** – Time interval, in which the maximum and minimum values define the upper and lower values to compare against.

For more information, see “Specify Behavioral Constraints” on page 42-18.

Prepare the System for Conversion



Click the **Prepare** button. The Fixed-Point Tool creates a backup version of the model and checks the system under design and the model containing the system under design for compatibility with the conversion process.

When possible, the Fixed-Point Tool automatically changes settings that are not compatible. In cases where the tool is not able to automatically change the settings, it notifies you of the changes you must make manually to help the conversion process be successful. For more information about preparation checks, see “Use the Fixed-Point Tool to Prepare a System for Conversion” on page 38-2.

To learn about the next step in the conversion process, see “Collect Ranges” on page 42-20.

See Also

More About

- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9
- “Best Practices for Fixed-Point Conversion Workflow” on page 42-5

Specify Behavioral Constraints

To determine if the numerical behavior of a new fixed-point implementation is acceptable, define constraints by setting signal tolerances, by using one or more model verification blocks, or both.

In the `Optimized Fixed-Point Conversion` workflow of the Fixed-Point Tool, or when using `fxpopt` at the command line, you must specify at least one behavioral constraint. Data types are optimized to meet all specified constraints.

In the `Iterative Fixed-Point Conversion` workflow of the Fixed-Point Tool, or when using `DataTypeWorkflowConverter` at the command line, you can specify behavioral constraints to verify the numerical behavior of the model with embedded types. After simulating with embedded types, the **Workflow Browser** indicates whether the embedded run meets the specified signal tolerances compared to the range collection run. For more information, see “Verify New Settings” on page 42-30.

Specify Signal Tolerances

You can specify tolerances for signals in your model that have signal logging enabled. To enable signal logging:

- In the Simulink Editor, select one or more signals.
- In the **Signal** tab of the Simulink Editor, click **Log Signals**.

In the Fixed-Point Tool, specify individual signal tolerances in the table under **Signal Tolerances**. The table contains all signals in the model with signal logging enabled. If you log additional signals after opening the Fixed-Point Tool, click **Refresh Signals** to update the **Signal Tolerances** table. At the command line, specify tolerances using the `addTolerance` method.

You can specify any of the following types of tolerances:

- **Abs Tol** — Absolute value of the maximum acceptable difference between the original signal and the signal in the converted design.
- **Rel Tol** — Maximum relative difference, specified as a percentage, between the original signal and the signal in the converted design. For example, a value of `1e-2` indicates a maximum relative difference of one percent.
- **Time Tol (seconds)** — Time interval, in which the maximum and minimum values define the upper and lower values to compare against.

Enter signal tolerances using any valid MATLAB expression that returns a finite, non-negative value.

You can define a tolerance band using any combination of absolute, relative, and time tolerance values. When you specify the tolerance for your signal using multiple types of tolerances, the overall tolerance band is computed by selecting the most lenient tolerance result for each data point. For more information about how tolerances are computed, see “Tolerance Computation”.

Use Model Verification Blocks

You can use enabled “Model Verification” blocks to specify constraints on the behavior of your system. For examples of data type optimization using model verification blocks, see “Optimize Data Types Using Multiple Simulation Scenarios” on page 40-20 and “Image Denoising Using Fixed-Point Quantized Restricted Boltzmann Machine Algorithm” on page 40-33.

See Also

“Optimize Fixed-Point Data Types for a System” on page 40-14

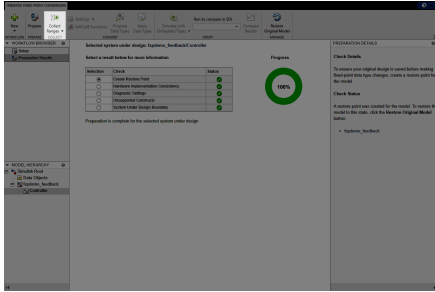
More About

- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9

Collect Ranges

After preparing the system under design for conversion as described in “Prepare System for Conversion” on page 42-14, collect ranges for the objects in your model.

Collect Ranges



To collect ranges, click the **Collect Ranges** button.

If you selected to collect ranges via simulation, the Fixed-Point Tool simulates the model with instrumentation to collect minimum and maximum values for each object in your model. By default, the Fixed-Point Tool uses the current data type override set on the model. You can also choose to override data types in your model with doubles or singles during the range collection simulation. The tool displays the results of the simulation in the spreadsheet and highlights any simulation results that have issues, such as overflows due to wrap or saturations.

Note Data type override does not apply to Boolean or enumerated data types.

If you defined a `Simulink.SimulationInput` object with multiple simulation scenarios, the **Workflow Browser** shows the results of each simulation, as well as the results of all simulation scenarios merged.

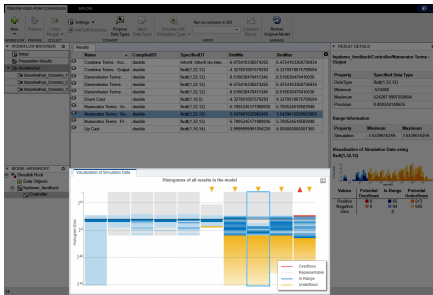
If you opted to collect ranges via range analysis, the Fixed-Point Tool uses the specified design ranges to derive ranges for the remaining objects in the system under design.

If you chose to collect ranges via simulation with range analysis, the Fixed-Point Tool uses the union of ranges collected via simulation and derivation.

If the analysis successfully derives range data for the model, the Fixed-Point Tool displays the derived minimum and maximum values for the blocks in the selected system. Before proposing data types, review the results.

If the analysis fails, examine the error messages and resolve the issues. See “Resolve Range Analysis Issues” on page 49-27.

Explore Collected Ranges



Using the Visualization of Simulation Data pane, you can view a summary of histograms of the bits used by each object in your model. Each column in the data type visualization represents a histogram for one object in your model. Each bin in a histogram corresponds to a bit in the binary word.

Selecting a column highlights the corresponding model object in the spreadsheet of the Fixed-Point Tool, and populates the **Result Details** pane with more detailed information about the selected result.

You can use the data type visualization to see a summary of the ranges of objects in your model and to spot sources of overflows, underflows, and inefficient data types. Using the **Explore** tab of the Fixed-Point Tool, you can sort and filter results in the tool based on additional criteria.

To learn about the next step in the conversion process, see “Convert Data Types” on page 42-22.

See Also

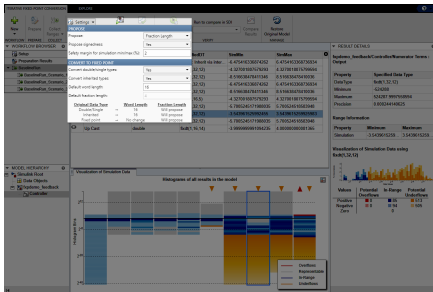
More About

- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9

Convert Data Types

After collecting ranges as described in “Collect Ranges” on page 42-20, propose and apply data types for objects in your model based on the collected ideal ranges stored in the baseline run. The Fixed-Point Tool proposes a data type for all objects in the system under design whose **Lock output data type setting against changes by the fixed-point tools** parameter is cleared.

Edit Proposal Settings

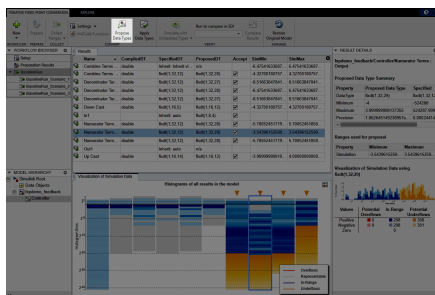


In the **Convert** section of the toolstrip, under the **Settings** menu, configure the settings that the Fixed-Point Tool uses to generate data type proposals for objects in your system under design.

| Setting | Description |
|---------------------------|---|
| Propose | <p>Select whether to propose fraction lengths or word lengths for objects in the system under design.</p> <ul style="list-style-type: none"> When you select Word Length, the Fixed-Point Tool uses range information and the specified Default fraction length value to propose word lengths for the objects in your model. When you select Fraction Length, the Fixed-Point Tool uses the range information and the specified Default word length value to propose best-precision fraction lengths for the objects in your model. |
| Propose signedness | Select whether to use the collected range information to propose signedness. |

| Setting | | Description |
|-------------------------------|---|---|
| | Safety margin for simulation min/max (%) | Specify a safety margin to apply to collected simulation ranges. The Fixed-Point Tool will add the specified amount to the collected ranges and base proposals on this larger range. The default value for this setting is two percent. |
| Convert to Fixed Point | Convert double/single types | Select whether to generate proposals for results that currently specify a double or single data type. |
| | Convert inherited types | Select whether to generate data type proposals for results that currently specify an inherited data type. |
| | Default word length | Select the default word length to use for proposals. This setting is enabled only when the Propose setting is set to Fraction Length . The default value for this setting is 16. |
| | Default fraction length | Select the default fraction length to use for proposals. This setting is enabled only when the Propose setting is set to Word Length . The default value for this setting is 4. |

Propose Data Types



When proposing data types, the Fixed-Point Tool uses the following types of range data:

- Design minimum or maximum values — You specify a design range for model objects using parameters such as **Output minimum** and **Output maximum**. For a list of blocks for which you can specify these values, see “Blocks That Allow Signal Range Specification”.
- Simulation minimum or maximum values — When simulating a system with instrumentation enabled, the Fixed-Point Tool logs the minimum and maximum values generated by model objects.

For more information about instrumentation settings, see “Fixed-Point Instrumentation and Data Type Override” on page 42-61.

If you specified multiple simulation scenarios through a `Simulink.SimulationInput` object, the Fixed-Point Tool proposes data types based on the merged ranges of all simulations.

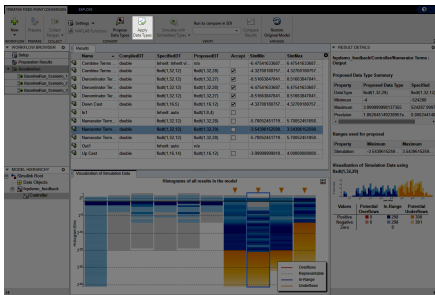
- Derived minimum or maximum values — When deriving minimum and maximum values for a selected system, the Fixed-Point Tool uses the design minimum and maximum values that you specify on the blocks to derive range information for signals in your model. For more information, see “How Range Analysis Works” on page 43-2.

The Fixed-Point Tool uses all available range data to calculate data type proposals.

To generate proposals, click the **Propose data types** button .

Apply Proposed Data Types

After reviewing the data type proposals, apply the proposed data types to your model.





The Fixed-Point Tool allows you to apply data type proposals selectively to objects in your model. In the spreadsheet, use the **Accept** check box to specify the proposals that you want to assign to model objects.

- The Fixed-Point Tool applies the proposed data type to this object. By default, the tool selects the **Accept** check box when a proposal differs from the current data type of the object.
- The Fixed-Point Tool ignores the proposed data type and leaves the current data type intact for this object.

No proposal exists for this object, for example, the object is locked against automatic data typing.

- 1 Examine each result. For more information about a particular result, select the result and examine the **Result Details** pane.

This pane also describes potential issues or errors and suggests methods for resolving them.


Results for which the data type proposal may cause issues, are marked with a warning () or an error () icon. For more detail on the information contained in the **Result Details** pane, see “Examine Results to Resolve Conflicts” on page 42-26.

- 2 If you do not want to accept the proposal for a result, on the spreadsheet, clear the **Accept** check box for that result.

Before applying proposals to your model, you can customize them. In the spreadsheet, click a **ProposedDT** cell and edit the data type expression. Some results belong to a data type group in

which they must all share the same data type. In these cases, the Fixed-Point Tool will report an error unless all results in the data type group share the same data type.

3

To write the proposed data types to the model, click the **Apply Data Types** button .

To complete the next step in the conversion process, see “Verify New Settings” on page 42-30.

See Also

More About

- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9

Examine Results to Resolve Conflicts

After proposing data types with the Fixed-Point Tool as described in “Convert Data Types” on page 42-22, you can examine each proposal using the **Result Details** pane. This pane displays the rationale for the proposed data types and a histogram plot of the signal. This tab also describes potential issues or errors and suggests methods for resolving them. To view the details, in the **Results** spreadsheet, select an object that has a proposed data type. The **Result Details** pane will update with information related to the selected result.

RESULT DETAILS

fxpdemo_feedback/Controller/In1

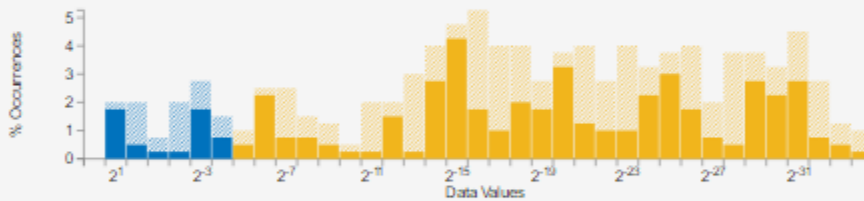
Proposed Data Type Summary

| Property | ProposedDT | SpecifiedDT |
|-----------|--------------|---------------|
| DataType | fixdt(1,8,4) | Inherit: auto |
| Minimum | -8 | |
| Maximum | 7.9375 | |
| Precision | 0.0625 | |

Ranges used for proposal

| Property | Minimum | Maximum |
|-------------------|---------|--------------------|
| Shared Simulation | -2 | 3.9999999999711746 |

Visualization of Simulation Data



| | Potential Overflows | In-Range | Potential Underflows |
|-----------------|---------------------|----------|----------------------|
| Positive Values | 0 | 21 | 178 |
| Negative Values | 0 | 23 | 177 |

Number of times zero occurred: 0

Proposal Details


- There is a requirement for the data type of this result to match the data type of other results.
 - [Highlight Elements Sharing Same Data Type](#)

Proposed Data Type Summary

The **Proposed Data Type Summary** section describes how the proposal differs from the currently specified data type of the object. For cases when the Fixed-Point Tool does not propose data types, this section provides a rationale. For example, the data type might be locked against changes by the fixed-point tools.

This section of the **Result Details** pane also informs you if the selected object must use the same data type as other objects in the model because of data type propagation rules. For example, the inputs to a Merge block must have the same data type. Therefore, the outputs of blocks that connect to these inputs must use the same data type. Similarly, blocks that are connected by the same element of a virtual bus must use the same data type.

Click **Highlight Elements Sharing Same Data Type** to highlight the objects that share data types in the model. To clear this highlighting, right-click in the model and select **Remove Highlighting**.

The Fixed-Point Tool allocates a tag to objects that must use the same data type. The tool displays this tag in the **DTGroup** column for the object. To view the **DTGroup** column, click the add column button  and select **DTGroup**.

Some Simulink blocks accept only certain data types on some ports. This section of the **Result Details** pane also informs you when a block that connects to the selected object has data type constraints that affect the proposed data type of the selected object.

The **Proposed Data Type Summary** section also provides a table with the proposed data type information:

| Item | Description |
|----------------------------|---|
| Proposed Data Type | The data type that the Fixed-Point Tool proposes for this object and the minimum and maximum values that the proposed data type can represent |
| Specified Data Type | The data type that an object specifies |

Needs Attention

This section lists potential issues and errors associated with the data type proposals, describes the issues, and suggests methods for resolving them.



Indicates a warning message



Indicates an error message

Range Information

This section provides a table with model object attributes that influence the data type proposal.

| Item | Description |
|---------------|--|
| Design | Design maximum and minimum values that an object specifies, such as its Output maximum and Output minimum parameters |

| Item | Description |
|------------|---|
| Simulation | The maximum and minimum values that occur during simulation |

Shared Values

When proposing data types, the Fixed-Point Tool attempts to satisfy data type requirements that model objects impose on one another. For example, the Sum block has an option that requires all its inputs to have the same data type. As a result, the table might also list attributes of other model objects that affect the proposal for the selected object. In such cases, the table displays these types of shared values:

- **Initial Values** — Some model objects have parameters that allow you to specify the initial values of their signals. For example, the Constant block has a **Constant value** parameter that initializes the block output signal. The Fixed-Point Tool uses initial values to propose data types for model objects whose design and simulation ranges are unavailable. With data type dependencies, the tool determines how initial values impact the proposals for neighboring objects.
- **Model-Required Parameters** — Some model objects require you to specify numeric parameters to compute the value of their outputs. For example, the **Table data** parameter of an n-D Lookup Table block specifies values that the block requires to perform a lookup operation and generate output. When proposing data types, the Fixed-Point Tool considers how this parameter value required by the model impacts the proposals for neighboring objects.

Examine the Results and Resolve Conflicts

- 1 In the **Results** spreadsheet, click the column header of the column containing the block icons. This action sorts the results so any results that contain conflicts with proposed data types appear at the top of the list.

Potential issues for each object appear coded by color in the list.




The proposed data type poses no issues for this object.



The proposed data type poses potential issues for this object.



The proposed data type will introduce data type errors if applied to this object.

- 2 Review and fix each error. Select the result with the error, then double-click the block icon in the spreadsheet to highlight the result in the Simulink editor. Use the information in the **Needs Attention** section of the **Result Details** pane to resolve the conflict.
- 3 Review the **Result Details** pane for the warnings and correct the problem, if necessary.
- 4 If you have changed the Simulink model, the baseline data, restore point, and preparation checks are not up to date. Start a new analysis of the updated data by clicking the **New** button and selecting **Iterative Fixed-Point Conversion**. Review the **Setup** pane, click **Prepare** to create a new restore point, then click the **Collect Ranges** button to rerun the simulation, or to derive new ranges. To generate new data type proposals, click **Propose Data Types**.
- 5 To generate a proposal, click **Propose Data Types** .

You are now ready to apply the proposed data types to the model. For more information, see “Apply Proposed Data Types” on page 42-24.

Verify New Settings

After applying proposed data types to your model as described in “Convert Data Types” on page 42-22, simulate the model using the applied fixed-point data types, and compare the fixed-point behavior of the system with the floating-point behavior.

Simulate Using Embedded Types

In the **Verify** section of the toolstrip, click the **Simulate with Embedded Types** button. The Fixed-Point Tool simulates the model using the new fixed-point data types. It logs minimum and maximum values, overflow data for all objects in the system under design. The tool stores the run information in a new run named EmbeddedRun. To edit the default name for the embedded run, under the **Simulate with Embedded Types** menu, type a new name in the **Run name** field.

If you specified multiple simulation scenarios using a `Simulink.SimulationInput` object, the tool simulates the model using the fixed-point data types for each simulation scenario.

Examine Visualization

The screenshot displays the MATLAB Fixed-Point Designer interface. The top toolbar includes buttons for 'New', 'Prepare', 'Collect Ranges', 'Settings', 'MATLAB Functions', 'Propose Data Types', 'Apply Data Types', 'Simulate with Embedded Types', 'Compare Results', and 'Restore Original Model'. The 'Simulate with Embedded Types' button is highlighted, and a dropdown menu shows 'EmbeddedRun' selected.

The main workspace is divided into several panes:

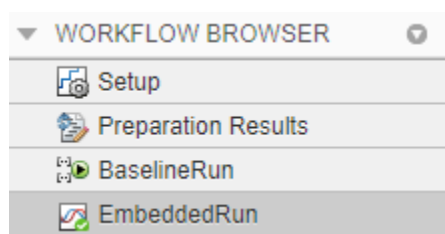
- Workflow Browser:** Shows a hierarchy of runs including 'BaselineRun', 'EmbeddedRun', and three 'EmbeddedRun_Scenario' sub-runs.
- Results Table:** A table with columns: Name, CompiledDT, SpecifiedDT, ProposedDT, Accept, SimMin, and SimMax. It lists various data types like 'Combine Terms', 'Denominator Te...', 'Down Cast', 'Numerator Term...', and 'Up Cast'.
- Result Details:** Shows properties for 'fxpdemo_feedback/Controller/Numerator Terms : Output', including Data Type (fixdt(1,32,29)), Minimum (-4), Maximum (3.999999998137355), and Precision (1.862645149230957e-09).
- Visualization of Simulation Data:** A histogram titled 'Histograms of all results in the model'. The y-axis is 'Histogram Bins' on a log scale from 2^{-28} to 2^2 . The x-axis is 'Data Values' on a log scale from 2^{-28} to 2^{28} . The histogram shows a distribution of values, with a legend indicating 'Overflows' (red), 'Representable' (grey), 'In-Range' (blue), and 'Underflows' (yellow).
- Range Information:** A table showing 'Property', 'Minimum', and 'Maximum' for 'Simulation', with values -3.5098876953... and 3.5993194580... respectively.
- Visualization of Simulation Data using fixdt(1,32,29):** A small histogram showing the distribution of values for the specified data type.
- Values Summary:** A table showing 'Values' (Positive: 287, Negative: 287, Zero: 623) and 'In-Range' (Positive: 287, Negative: 287, Zero: 623).

After simulating with embedded types, the **Visualization of Simulation Data** pane displays the new run data. Examine the histogram visualization to view the dynamic range of the objects in your model using the newly applied fixed-point data types.




Using the **Explore** tab of the Fixed-Point Tool, you can also sort and filter the results according to different criteria.

Compare Results

The **Workflow Browser** indicates whether the embedded run meets the specified signal tolerances compared to the range collection run. If there were multiple simulation scenarios, the tool indicates whether each scenario met the required tolerances.



The **Workflow Browser** displays one of the following.

| Icon | Status | Description |
|---|--------|---|
|  | Pass | All signals with a specified tolerance are within the specified tolerances in all embedded runs. |
|  | Warn | One of the following conditions occurred: <ul style="list-style-type: none"> No signals logged or no tolerances set in the model. Unable to compare the signals because the signals don't exist in both the range collection and the verification runs. The range collection run is no longer available. The range collection run used for data type proposals is two merged simulations. |
|  | Fail | One or more signals with a specified tolerance are not within the specified tolerances in any of the embedded runs. |

To compare the ideal results stored in `BaselineRun` with the fixed-point results, select the embedded run from the **Run to compare in SDI** drop down menu. Then click **Compare Results** to

open the Simulation Data Inspector. Alternatively, you can right-click the embedded run name in the Workflow Browser and select Open SDI.

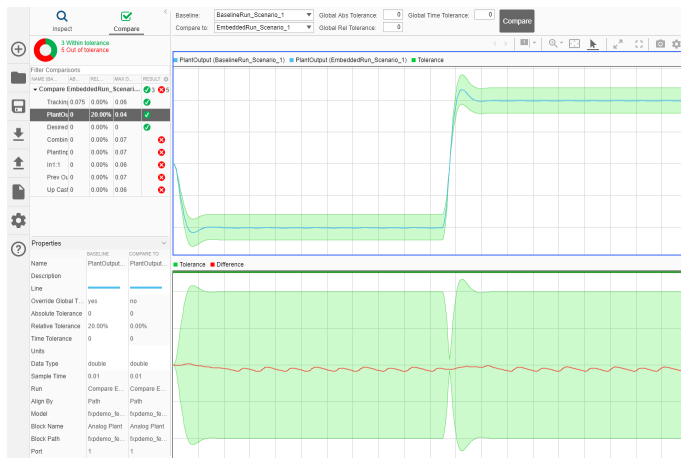
The screenshot shows the MATLAB/Simulink interface with the Simulation Data Inspector (SDI) workflow browser on the left. The 'EmbeddedRun' folder is expanded, showing three scenarios: EmbeddedRun_Scenario_1, EmbeddedRun_Scenario_2, and EmbeddedRun_Scenario_3. A context menu is open over EmbeddedRun_Scenario_1, with 'Run to compare in SDI' selected. The main workspace displays a table of simulation results for the selected scenario.

| Name | CompiledDT | SpecifiedDT | SimMin | SimMax |
|-------------------|----------------|------------------------|-------------------|------------------|
| Combine Terms ... | fixdt(1,32,28) | Inherit: Inherit vi... | | |
| Combine Terms ... | fixdt(1,32,28) | fixdt(1,32,28) | -4.3019438087... | 4.34968733787... |
| Denominator Te... | fixdt(1,32,27) | fixdt(1,32,27) | -8.56116962432... | 8.46746575832... |
| Denominator Te... | fixdt(1,32,27) | fixdt(1,32,27) | -6.5716950893... | 6.54254439473... |
| Down Cast | fixdt(1,16,12) | fixdt(1,16,12) | -8.56116962432... | 8.46746575832... |
| Numerator Term... | fixdt(1,32,28) | fixdt(1,32,28) | -4.302001953125 | 4.349609375 |
| Numerator Term... | fixdt(1,32,28) | fixdt(1,32,28) | -5.7659530639... | 5.78950881958... |
| Numerator Term... | fixdt(1,32,29) | fixdt(1,32,29) | -3.5098876953... | 3.59931945800... |
| Numerator Term... | fixdt(1,32,28) | fixdt(1,32,28) | -5.7659530639... | 5.78950881958... |
| Up Cast | fixdt(1,16,12) | fixdt(1,16,12) | -4 | 4.0625 |

Below the table, the 'Visualization of Simulation Data' section shows histograms of all results in the model. The y-axis is labeled 'Histogram Bins' and ranges from 2^{-28} to 2^2 . The x-axis is labeled 'Data Values' and ranges from 2^1 to 2^{11} . A legend indicates four categories: Overflows (red), Representable (grey), In-Range (blue), and Underflows (orange). The 'In-Range' category is highlighted with a blue box.

On the right, the 'RESULT DETAILS' section shows properties for 'fxpdemo_feedback/Controller/Numerator Terms : Output'. The 'Specified Data Type' is 'fixdt(1,32,29)'. Other properties include Minimum (-4), Maximum (3.999999998137355), and Precision (1.862645149230957e-09). Below this, 'Range Information' shows Minimum (-3.5098876953...) and Maximum (3.59931945800...). A 'Visualization of Simulation Data using fixdt(1,32,29)' histogram shows the distribution of values, with a legend for 'Values' (Positive: 287, Negative: 287, Zero: 623) and 'In-Range' (287, 287, 623).

The Simulation Data Inspector displays the comparison plots for the logged signals.



Note This step requires that you run a simulation during the “Collect Ranges” on page 42-20 phase of the conversion. If you use range analysis to collect ideal ranges for your system under design and do not run a simulation, you will not have a baseline run to compare to at this step.

If the behavior of the converted system does not meet your requirements, you can propose new data types after applying new proposal settings. For more information, see “Explore Additional Data Types” on page 42-34.

See Also

More About

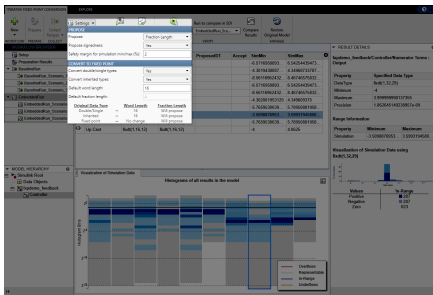
- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9

Explore Additional Data Types

After you simulate your model using embedded types and compare the floating-point and fixed-point behavior of your system, determine if the new fixed-point behavior is satisfactory. If the behavior of the system using the newly applied fixed-point data types is not acceptable, you can iterate through the process until you find settings that work for your system.

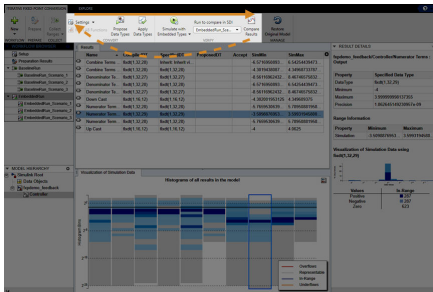
Edit Proposal Settings

In the **Convert** section of the toolstrip, under the **Settings** menu, alter the proposal settings that the Fixed-Point Tool uses to generate data type proposals for objects in your system under design.



Propose, Apply, Simulate, Compare

Click the **Propose Data Types** button to generate data type proposals using the new settings. After examining the new proposals in the spreadsheet, click the **Apply Data Types** button.



Iterate

Simulate the model using the newly applied data types and compare the behavior as described in “Verify New Settings” on page 42-30. Continue to iterate through this process (edit proposal settings, propose data types, apply data types, verify system behavior) until you find settings for which your system's fixed-point behavior is acceptable.

Restore Model to Original State

During the “Prepare System for Conversion” on page 42-14 step, the Fixed-Point Tool creates a restore point for your model. After the conversion process, if you want to restore your model to its state at the start of the conversion process, in the Fixed-Point Tool, click **Restore Original Model**.

The Fixed-Point Tool closes your model and reopens the model in its original state. Any changes made to your model after the preparation stage of the conversion are removed.

See Also

More About

- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9

Restore Model to Original State

During the “Prepare System for Conversion” on page 42-14 step, the Fixed-Point Tool creates a restore point for your model. After the conversion process, if you want to restore your model to its state at the start of the conversion process, in the Fixed-Point Tool, click **Restore Original Model**. The Fixed-Point Tool closes your model and reopens the model in its original state. Any changes made to your model after the preparation stage of the conversion are removed.

See Also

More About

- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9

Get Proposals for Results with Inherited Types

Blocks can inherit data types from a variety of sources, including signals to which they are connected and particular block parameters. The following table lists the types of inheritance rules that a block might specify.

| Inheritance Rule | Description |
|---------------------------------------|--|
| Inherit: Inherit via back propagation | Simulink automatically determines the output data type of the block during data type propagation. In this case, the block uses the data type of a downstream block or signal object. |
| Inherit: Same as input | The block uses the data type of its sole input signal for its output signal. |
| Inherit: Same as first input | The block uses the data type of its first input signal for its output signal. |
| Inherit: Same as second input | The block uses the data type of its second input signal for its output signal. |
| Inherit: Inherit via internal rule | The block uses an internal rule to determine its output data type. The internal rule chooses a data type that optimizes numerical accuracy, performance, and generated code size, while taking into account the properties of the embedded target hardware. It is not always possible for the software to optimize efficiency and numerical accuracy at the same time. |

How to Get Proposals for Objects That Use an Inherited Output Data Type

To enable proposals for results that specify an inherited output data type, in the Fixed-Point Tool, in the **Convert** section of the toolstrip, under **Settings**, set the **Convert inherited types** setting to **Yes**.

For objects that specify an inherited output data type, the Fixed-Point Tool proposes a new data type based on collected ranges and the specified proposal settings.

When the Fixed-Point Tool Will Not Propose for Inherited Data Types

The Fixed-Point Tool proposes data types only for the **Output data type** parameter of a block or model object. It will not propose for other block data types, such as the **Accumulator data type** of a Sum block, or the **Gain** parameter in a Gain block.

The Fixed-Point Tool will also not propose for the following model objects if they use an inherited output data type.

- Signal objects
- Stateflow charts
- Bus objects

- MATLAB variables

See Also

More About

- “Data Type Propagation Errors After Applying Proposed Data Types” on page 49-25

Rescale a Fixed-Point Model

In this section...

“About the Feedback Controller Example Model” on page 42-39

“Explore the Numerical Behavior of the Model” on page 42-43

“Propose Fraction Lengths Using Simulation Range Data” on page 42-45

This example shows you how to use the Fixed-Point Tool to refine the scaling of fixed-point data types associated with the feedback controller model. Although the tool enables multiple workflows for converting a digital controller described in ideal double-precision numbers to one realized in fixed-point numbers, this example uses the following approach:

- Range Collection — Use the range collection workflow to explore the numerical behavior of the model.

Perform a global override of the fixed-point data types using double-precision numbers. The Simulink software logs the simulation results and the Fixed-Point Tool displays them.

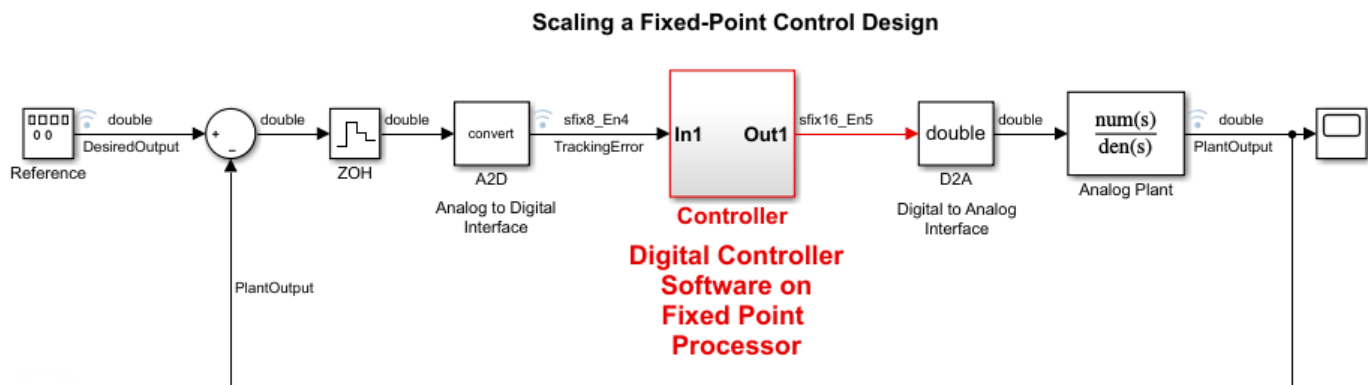
Run an initial simulation using a reasonable guess at the fixed-point word size and scaling, then compare the simulation results to the double-precision run. This task illustrates how difficult it is to guess the best scaling.

- Propose Fraction Lengths Using Simulation Range Data — Use the iterative fixed-point conversion workflow to autoscale the model.

The Fixed-Point Tool uses double-precision simulation results to propose fixed-point scaling for appropriately configured blocks. The Fixed-Point Tool allows you to accept and apply the scaling proposals selectively. Afterward, you determine the quality of the results by examining the input and output of the model's analog plant.

About the Feedback Controller Example Model

To open the Simulink feedback design model for this tutorial, at the MATLAB command line, type `fxpdemo_feedback`.



The model consists of the following blocks and subsystems:

- **Reference**

This Signal Generator block generates a continuous-time reference signal. It is configured to output a square wave.

- **Sum**

This Add block subtracts the plant output from the reference signal.

- **ZOH**

The Zero-Order Hold block samples and holds the continuous signal. This block is configured so that it quantizes the signal in time by 0.01 seconds.

- **Analog to Digital Interface**

The analog to digital (A/D) interface consists of a Data Type Conversion block that converts a `double` to a fixed-point data type. It represents any hardware that digitizes the amplitude of the analog input signal. In the real world, its characteristics are fixed.

- **Controller**

The digital controller is a subsystem that represents the software running on the hardware target. Refer to “Digital Controller Realization” on page 42-41.

- **Digital to Analog Interface**

The digital to analog (D/A) interface consists of a Data Type Conversion block that converts a fixed-point data type into a `double`. It represents any hardware that converts a digitized signal into an analog signal. In the real world, its characteristics are fixed.

- **Analog Plant**

The analog plant is described by a transfer function, and is controlled by the digital controller. In the real world, its characteristics are fixed.

- **Scope**

The model includes a Scope block that displays the plant output signal.

Simulation Setup

To set up this kind of fixed-point feedback controller simulation:

- 1 Identify all design components.

In the real world, there are design components with fixed characteristics (the hardware) and design components with characteristics that you can change (the software). In this feedback design, the main hardware components are the A/D hardware, the D/A hardware, and the analog plant. The main software component is the digital controller.

- 2 Develop a theoretical model of the plant and controller.

For the feedback design in this tutorial, the plant is characterized by a transfer function.

The digital controller model in this tutorial is described by a z -domain transfer function and is implemented using a direct-form realization.

- 3 Evaluate the behavior of the plant and controller.

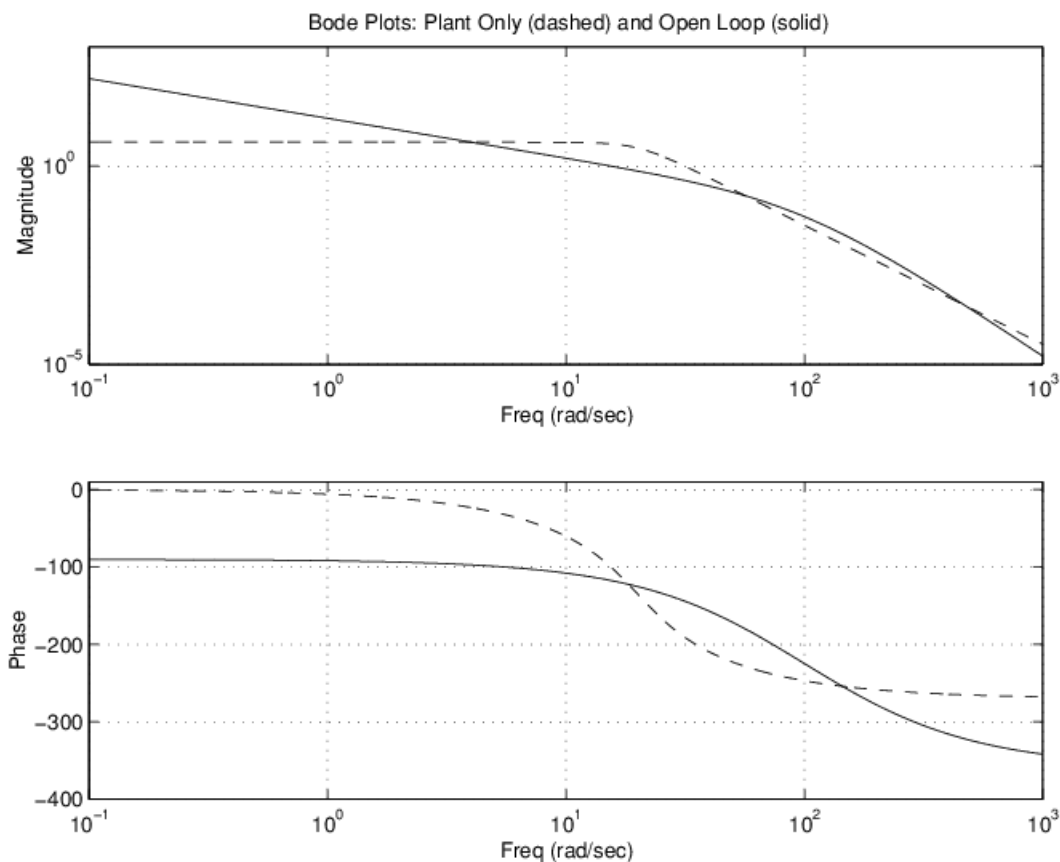
You evaluate the behavior of the plant and the controller with a Bode plot. This evaluation is idealized, because all numbers, operations, and states are double-precision.

4 Simulate the system.

You simulate the feedback controller design using Simulink and Fixed-Point Designer software. In a simulation environment, you can treat all components (software *and* hardware) as though their characteristics are not fixed.

Idealized Feedback Design

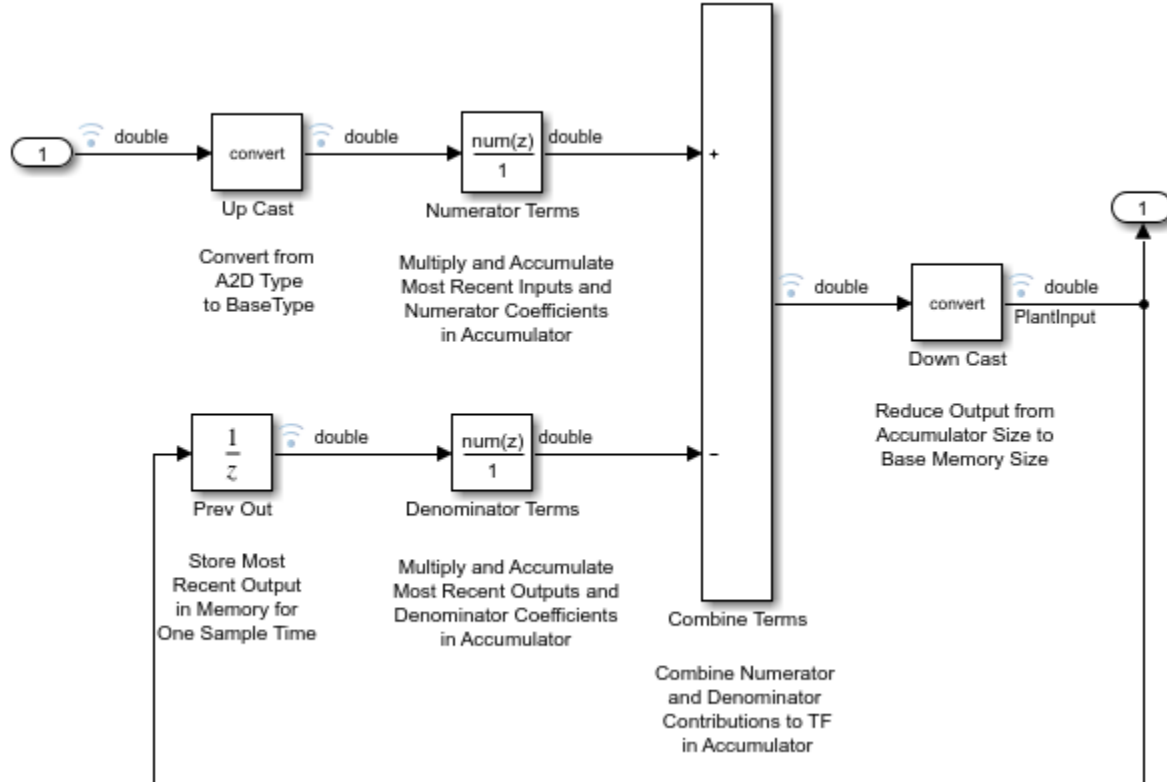
Open loop (controller and plant) and plant-only Bode plots for the “Scaling a Fixed-Point Control Design” model are shown in the following figure. The open loop Bode plot results from a digital controller described in the idealized world of continuous time, double-precision coefficients, storage of states, and math operations.



The Bode plots were created using workspace variables produced by a script named `preload_feedback.m`.

Digital Controller Realization

In this simulation, the digital controller is implemented using the fixed-point direct form realization shown in the following diagram. The hardware target is a 16-bit processor. Variables and coefficients are generally represented using 16 bits, especially if these quantities are stored in ROM or global RAM. Use of 32-bit numbers is limited to temporary variables that exist briefly in CPU registers or in a stack.



The digital controller realization consists of these blocks:

- **Up Cast**

Up Cast is a Data Type Conversion block that connects the A/D hardware with the digital controller. It pads the output word size of the A/D hardware with trailing zeros to a 16-bit number (the base data type).

- **Numerator Terms and Denominator Terms**

Each of these Discrete FIR Filter blocks represents a weighted sum carried out in the CPU target. The word size and precision in the calculations reflect those of the accumulator. Numerator Terms multiplies and accumulates the most recent inputs with the FIR numerator coefficients. Denominator Terms multiplies and accumulates the most recent delayed outputs with the FIR denominator coefficients. The coefficients are stored in ROM using the base data type. The most recent inputs are stored in global RAM using the base data type.

- **Combine Terms**

Combine Terms is a Add block that represents the accumulator in the CPU. Its word size and precision are twice that of the RAM (double bits).

- **Down Cast**

Down Cast is a Data Type Conversion block that represents taking the number from the CPU and storing it in RAM. The word size and precision are reduced to half that of the accumulator when converted back to the base data type.

- **Prev Out**

Prev Out is a Unit Delay block that delays the feedback signal in memory by one sample period. The signals are stored in global RAM using the base data type.

Direct Form Realization

The controller directly implements this equation:

$$y(k) = \sum_{i=0}^N b_i u(k-1) - \sum_{i=1}^N a_i y(k-1),$$

where:

- $u(k-1)$ represents the input from the previous time step.
- $y(k)$ represents the current output, and $y(k-1)$ represents the output from the previous time step.
- b_i represents the FIR numerator coefficients.
- a_i represents the FIR denominator coefficients.

The first summation in $y(k)$ represents the multiplication and accumulation of the most recent inputs and numerator coefficients in the accumulator. The second summation in $y(k)$ represents the multiplication and accumulation of the most recent outputs and denominator coefficients in the accumulator. Because the FIR coefficients, inputs, and outputs are all represented by 16-bit numbers (the base data type), any multiplication involving these numbers produces a 32-bit output (the accumulator data type).

Explore the Numerical Behavior of the Model

Initial guesses for the scaling of each block are already specified in each block mask in the model. This task illustrates the difficulty of guessing the best fixed-point scaling. In this example, you compare the behavior of the model with an idealized floating-point version using the range collection workflow in the Fixed-Point Tool.

- 1 Open the `fxpdemo_feedback` model.
- 2 Open the Fixed-Point Tool. In the **Apps** gallery, select **Fixed-Point Tool**.
- 3 In the Fixed-Point Tool, click **New**, and select **Range Collection**.

You can use the range collection workflow to explore the numerical behavior of your model, and compare it to an idealized, floating-point version.

- 4 Under **System Under Design (SUD)**, select the subsystem you want to analyze. In this example, select **Controller**.
- 5 Under **Range Collection Mode**, select **Simulation Ranges** as the range collection method.
- 6 Under **Simulation Inputs**, use the default model inputs for simulation.
- 7 Click the **Collect Ranges** button arrow and select **Double precision**. Click the **Collect Ranges** button to start the simulation.

The Simulink software simulates the `fxpdemo_feedback` model in data type override mode and stores the results in `BaselineRun`. Data type override enables you to perform a global override of the fixed-point data types with double-precision data types, thereby avoiding quantization effects. In the **Results** spreadsheet, the Fixed-Point Tool displays the run results. The compiled data type (**CompiledDT**) column for `BaselineRun` shows that the blocks in the model used a `double` data type during simulation.

- 8 Next, simulate the system using the fixed-point data types specified in the model. Click the **Settings** button arrow and select **Specified data types**. Click **Simulate with Embedded Types**.

The Fixed-Point Tool simulates the model using the currently specified fixed-point data types and stores the range information in **EmbeddedRun**. You can view the collected ranges in the **SimMin** and **SimMax** columns of the spreadsheet.

The Fixed-Point Tool highlights the row containing the **Up Cast** block to indicate that there is an issue with this result. The **Result Details** pane shows that the block saturated 23 times, which indicates a poor guess for its scaling.

Tip You can use the **Explore** tab to explore and filter results.

- 9 Click **Compare Results** to open the Simulation Data Inspector.
- 10 In the Simulation Data Inspector, select **PlantOutput** as the signal to compare.

Simulation Data Inspector plots the signal associated with the plant output for the **BaselineRun** and the **EmbeddedRun**.



The plot of the plant output signal for **EmbeddedRun** reflects the initial guess at scaling. The Bode plot design sought to produce a well-behaved linear response for the closed-loop system, represented by the ideal **BaselineRun**. However, the response of the **EmbeddedRun** is nonlinear.

Significant quantization effects cause the nonlinear features. An important part of fixed-point design is finding a scaling that reduces quantization effects to acceptable levels.

Propose Fraction Lengths Using Simulation Range Data

Using automatic data typing, you can maximize the precision of the output data type while spanning the full simulation range. The iterative fixed-point conversion workflow in the Fixed-Point Tool lets you maximize the precision of the output data types while spanning the full simulation range. This process is known as autoscaling.

Because no design range information is supplied in this example, the Fixed-Point Tool uses simulation range data for proposing data types. The **Safety margin for simulation min/max (%)** parameter value multiplies the “raw” simulation values. Setting this parameter to a value greater than 1 decreases the likelihood that an overflow will occur when fixed-point data types are being used. For more information about how the Fixed-Point Tool calculates data type proposals, see “How the Fixed-Point Tool Proposes Data Types” on page 42-48.

Because of the nonlinear effects of quantization, a fixed-point simulation produces results that are different from an idealized, doubles-based simulation. Signals in a fixed-point simulation can cover a larger or smaller range than in a doubles-based simulation. If the range increases enough, overflows or saturations could occur. A safety margin decreases this likelihood, but it might also decrease the precision of the simulation.

Note When the maximum and minimum simulation values cover the full, intended operating range of your design, the Fixed-Point Tool yields meaningful automatic data typing results.

Autoscale the Controller subsystem. This subsystem represents software running on the target, and requires optimization.

- 1 In the Fixed-Point Tool, click **New**, and select **Iterative Fixed-Point Conversion**.

Tip Switching workflows in the Fixed-Point tool clears the settings and any data collected during the active workflow. The model remains in its current state.

- 2 Under **System Under Design (SUD)**, select the Controller subsystem as the system to analyze and convert.
- 3 Under **Range Collection Mode**, select **Simulation ranges**.
- 4 Under **Simulation Inputs**, use the default model inputs for simulation.
- 5 Click **Prepare** to create a restore point and automatically prepare the system under design for conversion.
- 6 Click the **Collect Ranges** button arrow and select **Double precision**. Click **Collect Ranges** to start the simulation.

The Simulink software simulates the `fxpdemo_feedback` model in data type override mode and stores the results in `BaselineRun_2`.

- 7 In the **Convert** section, click the **Settings** button. Set the **Safety margin for simulation min/max (%)** parameter to 20. Use the default settings for all other parameters.
- 8 Click **Propose Data Types**.

The Fixed-Point Tool analyzes the scaling of all fixed-point blocks whose **Lock output data type setting against changes by the fixed-point tools** parameter is cleared.

The Fixed-Point Tool uses the minimum and maximum values stored in `BaselineRun_2` to propose each block's data types such that the precision is maximized while the full range of simulation values is spanned. The tool displays the proposed data types in the **Results** spreadsheet.

- 9 Review the scaling that the Fixed-Point Tool proposes. You can choose to accept the scaling proposal for each block. In the **Results** spreadsheet, select the corresponding **Accept** check box. By default, the Fixed-Point Tool accepts all scaling proposals that differ from the current scaling. For this example, ensure that the **Accept** check box is selected for each of the Controller subsystem's blocks.
- 10 Click the **Apply Data Types** button.

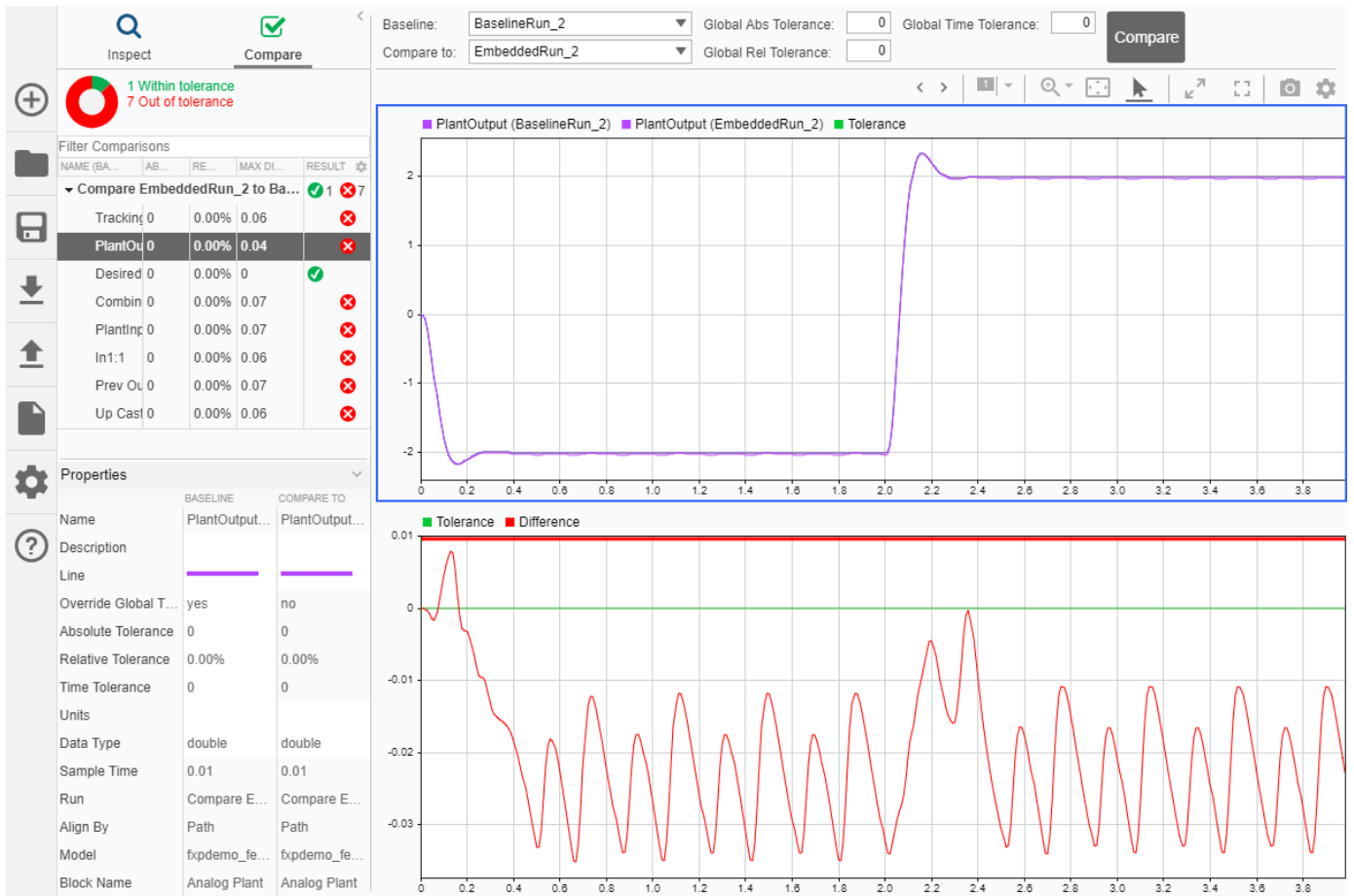
The Fixed-Point Tool applies the scaling proposals that you accepted in the previous step to the blocks in the `Controller` subsystem.

- 11 In the **Verify** section, click the **Simulate with Embedded Types** button.

Simulink simulates the `fxpdemo_feedback` model using the new scaling that you applied. Information about this simulation is stored in a run named `EmbeddedRun_2`. Afterward, the Fixed-Point Tool displays information about blocks that logged fixed-point data. The compiled data type (**CompiledDT**) column for `EmbeddedRun_2` shows that the Controller subsystem's blocks used fixed-point data types with the new scaling.

- 12 Click **Compare Results** to open the Simulation Data Inspector.
- 13 In the Simulation Data Inspector, select `PlantOutput` as the signal to compare.

Simulation Data Inspector plots the signal associated with the plant output for `BaselineRun_2` and `EmbeddedRun_2`, as well as their difference.



The plant output signal represented by the fixed-point run achieves a steady state, but a small limit cycle is present because of non-optimal A/D design.

See Also

Related Examples

- “How Hardware Implementation Settings Affect Data Type Proposals” on page 42-50

More About

- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9

How the Fixed-Point Tool Proposes Data Types

To use the Fixed-Point Tool to propose word lengths, you must specify the fraction length requirements for data types in the model. Select the fraction lengths based on the precision required for the system that you are modeling. If you do not specify fraction lengths, the Fixed-Point Tool sets the **Default fraction length** to 4. The Fixed-Point Tool uses these specified fraction lengths to recommend the minimum word length for objects in the selected model or subsystem to avoid overflow for the collected range information.

The proposed word length is based on:

- Design range information and range information that the Fixed-Point Tool collects. This collected range information can be either simulation range data, derived range data, or simulation with derived range data.
- The signedness and fraction lengths of data types that you specify on blocks, signal objects.
- The production hardware implementation settings specified in the Configuration Parameters.

How the Fixed-Point Tool Uses Range Information

The Fixed-Point Tool determines whether to use different types of range information based on its availability and on the Fixed-Point Tool setting.

Design range information always takes precedence over both simulation and derived range data. When there is no design range information, the Fixed-Point Tool uses either simulation or derived range data. If you specify a safety margin, the Fixed-Point Tool takes the safety margin into account.

For example, if a signal has a design range of $[-10, 10]$, the Fixed-Point Tool uses this range for the proposal and ignores all simulation and derived range information.

If the signal has no specified design information, but does have a simulation range of $[-8, 8]$ and a derived range of $[-2, 2]$, the proposal uses the union of the ranges, $[-8, 8]$. If you specify a safety margin of 50%, the proposal uses a range of $[-12, 12]$.

How the Fixed-Point Tool Uses Target Hardware Information

The Fixed-Point Tool calculates the ideal word length and then checks this length against the production hardware implementation settings for the target hardware.

- If the target hardware is an FPGA/ASIC, then the Fixed-Point Tool proposes the ideal word length. If the ideal word length is greater than 128, then the Fixed-Point Tool proposes 128.
- If the target hardware is an embedded processor, then the Fixed-Point Tool rounds the ideal word length up and proposes the nearest supported data type of your processor.

How to Get Proposals for Objects That Use an Inherited Output Data Type

Blocks can inherit data types from a variety of sources, including signals to which they are connected and particular block parameters. The following table lists the types of inheritance rules that a block might specify.

| Inheritance Rule | Description |
|---------------------------------------|--|
| Inherit: Inherit via back propagation | Simulink automatically determines the output data type of the block during data type propagation. In this case, the block uses the data type of a downstream block or signal object. |
| Inherit: Same as input | The block uses the data type of its sole input signal for its output signal. |
| Inherit: Same as first input | The block uses the data type of its first input signal for its output signal. |
| Inherit: Same as second input | The block uses the data type of its second input signal for its output signal. |
| Inherit: Inherit via internal rule | The block uses an internal rule to determine its output data type. The internal rule chooses a data type that optimizes numerical accuracy, performance, and generated code size, while taking into account the properties of the embedded target hardware. It is not always possible for the software to optimize efficiency and numerical accuracy at the same time. |

To enable proposals for results that specify an inherited output data type, in the Fixed-Point Tool, in the **Convert** section of the toolstrip, under **Settings**, set the **Convert inherited types** setting to Yes.

For objects that specify an inherited output data type, the Fixed-Point Tool proposes a new data type based on collected ranges and the specified proposal settings.

When the Fixed-Point Tool Will Not Propose for Inherited Data Types

The Fixed-Point Tool proposes data types only for the **Output data type** parameter of a block or model object. It will not propose for other block data types, such as the **Accumulator data type** of a Sum block, or the **Gain** parameter in a Gain block.

The Fixed-Point Tool will also not propose for the following model objects if they use an inherited output data type.

- Signal objects
- Stateflow charts
- Bus objects
- MATLAB variables

See Also

Related Examples

- “How Hardware Implementation Settings Affect Data Type Proposals” on page 42-50

How Hardware Implementation Settings Affect Data Type Proposals

In this section...

“Open the Model and Specify Hardware Implementation Settings” on page 42-50

“Propose Word Lengths Based on Simulation Data” on page 42-51

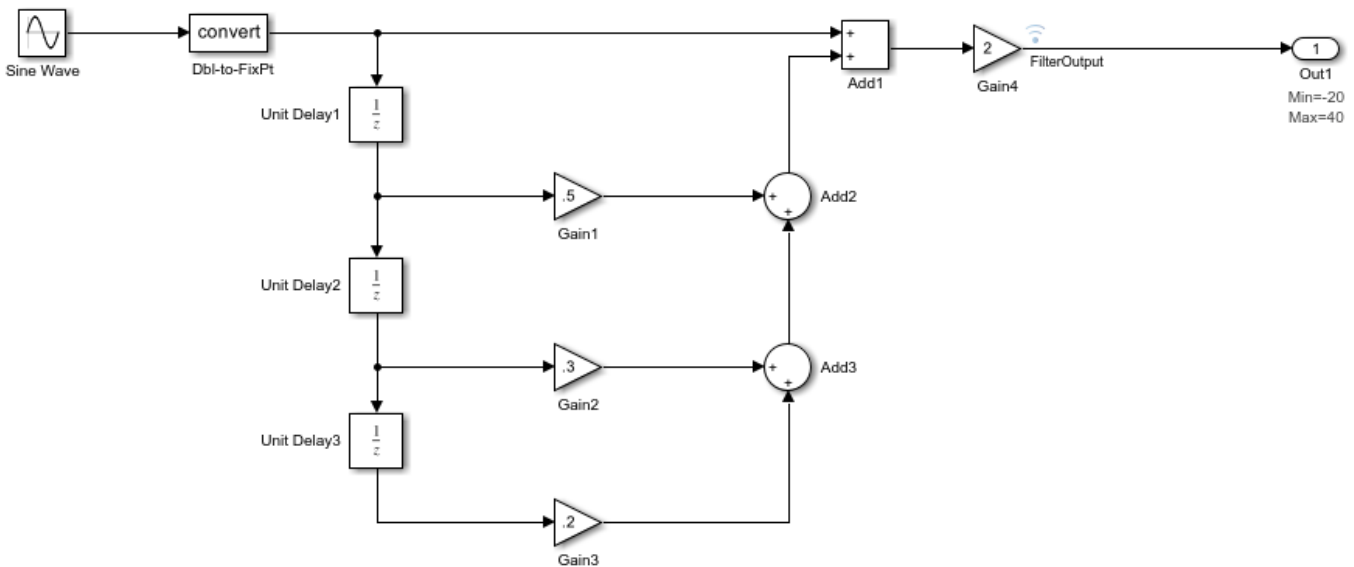
This example shows how to use the Fixed-Point Tool to propose word lengths for a model that implements a simple moving average algorithm. The model already uses fixed-point data types, but they are not optimal. Simulate the model and propose data types based on simulation data. To see how the target hardware affects the word length proposals, first set the target hardware to an embedded processor and propose word lengths. Then, set the target hardware to an FPGA and propose word lengths.

Open the Model and Specify Hardware Implementation Settings

In the Configuration Parameters dialog box, on the **Hardware Implementation** pane, you can specify the **Device vendor** and **Device type** of your target hardware. The Fixed-Point Tool uses this information when it proposes fixed-point data types for objects in your model. For example, if you specify the target hardware to be an embedded processor, the tool will propose standard word lengths appropriate for the target.

Open the `ex_moving_average` example.

```
open_system('ex_moving_average')
```



Copyright 2011-2012 The MathWorks, Inc.

Verify that the target hardware is an embedded processor. In the Configuration Parameters dialog box, on the **Hardware Implementation** pane, set the **Device vendor** to Custom Processor. Close the Configuration Parameters dialog box.

Propose Word Lengths Based on Simulation Data

Some blocks in the model already have specified fixed-point data types.

| Block | Data Type Specified on Block |
|-----------|------------------------------|
| Dbl2Fixpt | fixdt(1,16,10) |
| Gain1 | fixdt(1,32,17) |
| Gain2 | fixdt(1,32,17) |
| Gain3 | fixdt(1,32,17) |
| Gain4 | fixdt(1,16,1) |
| Add1 | fixdt(1,32,17) |
| Add2 | fixdt(1,32,17) |
| Add3 | fixdt(1,32,17) |

Use the iterative fixed-point conversion workflow in the Fixed-Point Tool to see how the target hardware affects word length proposals.

- 1 In the model, in the **Apps** gallery, select **Fixed-Point Tool**.
- 2 In the Fixed-Point Tool, click **New**, and select Iterative Fixed-Point Conversion.
- 3 In the Fixed-Point Tool, under **System Under Design (SUD)**, select `ex_moving_average` as the system to convert.
- 4 Under **Range Collection Mode**, select **Simulation ranges** as the method of range collection. This configures the model to collect ranges using idealized floating-point data types.
- 5 In the toolstrip, click **Prepare** to prepare the system for conversion.
- 6 Expand the **Collect Ranges** button arrow and select `Double precision`. Click **Collect Ranges** to start the simulation.

The Fixed-Point Tool stores the simulation data in a run titled `BaselineRun`. You can examine the range information of the blocks in the model in the spreadsheet.

- 7 In the **Convert** section of the toolstrip you can configure the data type proposal settings for the blocks. Click the **Settings** button arrow. In the Settings dialog, next to **Propose**, select `Word Length`.
- 8 Click **Propose Data Types**.

The Fixed-Point Tool uses available range data to calculate data type proposals according to the following rules:

- Design minimum and maximum values take precedence over the simulation range.
- The tool observes the simulation range because you selected **Simulation ranges** as the range collection method.

The **Safety margin for simulation min/max (%)** parameter specifies a range that differs from that defined by the simulation range. In this example, the default safety margin is used.

The Fixed-Point Tool analyzes the data types of all fixed-point blocks whose **Lock output data type setting against changes by the fixed-point tools** parameter is cleared.

For each object in the model, the Fixed-Point Tool proposes the minimum word length that avoids overflow for the collected range information. Because the target hardware is a 16-bit embedded processor, the Fixed-Point tool proposes word lengths based on the number of bits used by the processor for each data type. For more information, see “How the Fixed-Point Tool Uses Target Hardware Information” on page 42-48.

The tool proposes smaller word lengths for Gain4 and Gain4:Gain. The tool calculated that their ideal word length is less than or equal to the character bit length for the embedded processor (8), so the tool rounds up the word length to 8.

The screenshot displays the MATLAB Fixed-Point Tool interface. The main window is titled "ITERATIVE FIXED-POINT CONVERSION" and "EXPLORE". The "RESULTS" pane shows a table of conversion results for various blocks in the model. The "Histograms of all results in the model" pane shows a visualization of the data distribution for the proposed data types.

| Name | CompiledDT | SpecifiedDT | ProposedDT | Accept | SimMin | SimMax |
|-----------------------|----------------------|----------------------|----------------|-------------------------------------|--------|-----------------|
| Add3 : Accumulator | double | Inherit: Inherit ... | n/a | <input type="checkbox"/> | 0 | 2.5167860950... |
| Add3 : Output | double | fixdt(1,32,17) | fixdt(1,32,17) | <input type="checkbox"/> | 0 | 2.5167860950... |
| Data Type Conversion1 | double | fixdt(1,16,10) | fixdt(0,16,10) | <input checked="" type="checkbox"/> | 0 | 5.0488259088... |
| Gain1 | double | fixdt(1,32,17) | fixdt(0,32,17) | <input checked="" type="checkbox"/> | 0 | 2.5212423076... |
| Gain1 : Gain | Inherit: Inherit ... | Inherit: Inherit ... | n/a | <input type="checkbox"/> | | |
| Gain2 | double | fixdt(1,32,17) | fixdt(0,32,17) | <input checked="" type="checkbox"/> | 0 | 1.5108372258... |
| Gain2 : Gain | Inherit: Inherit ... | Inherit: Inherit ... | n/a | <input type="checkbox"/> | | |
| Gain3 | double | fixdt(1,32,17) | fixdt(0,32,17) | <input checked="" type="checkbox"/> | 0 | 1.0059488692... |
| Gain3 : Gain | Inherit: Inherit ... | Inherit: Inherit ... | n/a | <input type="checkbox"/> | | |
| Gain4 | double | fixdt(1,16,1) | fixdt(1,8,1) | <input checked="" type="checkbox"/> | 0 | 20.173708623... |
| Gain4 : Gain | fixdt(1,16,0) | fixdt(0,8,0) | fixdt(0,8,0) | <input checked="" type="checkbox"/> | | |
| Out1 | Inherit: auto | Inherit: auto | n/a | <input type="checkbox"/> | | |

The "Histograms of all results in the model" pane shows a visualization of the data distribution for the proposed data types. The x-axis represents the data value, and the y-axis represents the histogram bins. The legend indicates the following categories:

- Overflows (Red)
- Representable (Grey)
- In-Range (Blue)
- Underflows (Yellow)

The "RESULT DETAILS" pane shows the proposed data type summary for the selected block, "ex_moving_average/Gain4".

| Property | Proposed Data Type | Specific |
|-----------|--------------------|---------------|
| Data Type | fixdt(1,8,1) | fixdt(1,16,1) |
| Minimum | -64 | -16384 |
| Maximum | 63.5 | 16383.5 |
| Precision | 0.5 | 0.5 |

The "Ranges used for proposal" table shows the minimum and maximum values for the proposed data type.

| Property | Minimum | Maximum |
|----------------|---------|-----------------|
| Shared Design | -20 | 40 |
| Shared Simu... | 0 | 20.173708623... |
| Simulation | 0 | 20.173708623... |

The "Visualization of Simulation Data using fixdt(1,8,1)" pane shows a histogram of the data distribution for the proposed data type. The x-axis represents the data value, and the y-axis represents the percentage of data points. The legend indicates the following categories:

- Positive (Red)
- Negative (Blue)
- Zero (Yellow)
- Potential Overflows (Red)
- In-Range (Blue)
- Potential Underflows (Yellow)

The "Proposal Details" table shows the potential overflows, in-range, and underflows for the proposed data type.

| Values | Potential Overflows | In-Range | Potential Underflows |
|----------|---------------------|----------|----------------------|
| Positive | 0 | 501 | 11 |
| Negative | 0 | 0 | 0 |
| Zero | 0 | 1 | 0 |

- 9 To see how the target hardware affects the word length proposal, change the target hardware to FPGA/ASIC.
 - a In the Configuration Parameters dialog box, on the **Hardware Implementation** pane, set **Device vendor** to ASIC/FPGA.
 - b Click **Apply** and close the Configuration Parameters dialog box.
- 10 In the Fixed-Point Tool, click **Propose data types** again.

Because the target hardware is an FPGA, there are no constraints on the word lengths that the Fixed-Point Tool proposes. The word length for Gain4:Gain is now 2.

The screenshot displays the MATLAB Fixed-Point Designer interface. The main window shows a table of data type proposals for various components in the model. The 'Gain4' component is highlighted, showing a proposed data type of `fixdt(1,8,1)`. The interface includes a 'Propose Data Types' button and a 'Result Details' panel on the right.

| Name | CompiledDT | SpecifiedDT | ProposedDT | Accept | SimMin | SimMax |
|----------------------|------------|-----------------------------|-----------------------------|-------------------------------------|--------|-----------------|
| Data type Conversion | double | <code>fixdt(1,16,1)</code> | <code>fixdt(0,15,1)</code> | <input checked="" type="checkbox"/> | 0 | 5.4466299088... |
| Gain1 | double | Inherit: Inherit... | <code>fixdt(0,19,17)</code> | <input checked="" type="checkbox"/> | 0 | 2.5212423076... |
| Gain1 : Gain | | | | | | |
| Gain2 | double | <code>fixdt(1,32,17)</code> | <code>fixdt(0,18,17)</code> | <input checked="" type="checkbox"/> | 0 | 1.5108372258... |
| Gain2 : Gain | | | | | | |
| Gain3 | double | <code>fixdt(1,32,17)</code> | <code>fixdt(0,18,17)</code> | <input checked="" type="checkbox"/> | 0 | 1.0059488692... |
| Gain3 : Gain | | | | | | |
| Gain4 | double | <code>fixdt(1,16,1)</code> | <code>fixdt(1,8,1)</code> | <input checked="" type="checkbox"/> | 0 | 20.173708623... |
| Gain4 : Gain | | | | | | |
| Out1 | | Inherit: auto | n/a | | | |
| Unit Delay1 | | | n/a | | | |
| Unit Delay2 | | | n/a | | | |

The 'Result Details' panel for `ex_moving_average/Gain4` shows the following summary:

| Property | Proposed Data Type | Specific |
|-----------|---------------------------|----------------------------|
| Data Type | <code>fixdt(1,8,1)</code> | <code>fixdt(1,16,1)</code> |
| Minimum | -64 | -16384 |
| Maximum | 63.5 | 16383.5 |
| Precision | 0.5 | 0.5 |

Below the table, there are histograms for 'Histograms of all results in the model' and 'Visualization of Simulation Data using `fixdt(1,8,1)`'. The histograms show the distribution of data values across different components, with a legend indicating 'Overflows', 'Representable', 'In-Range', and 'Underflows'.

See Also

Related Examples

- “Rescale a Fixed-Point Model” on page 42-39

More About

- “How the Fixed-Point Tool Proposes Data Types” on page 42-48

Propose Data Types For Merged Simulation Ranges

In this section...

“Set up the Model” on page 42-54

“Open the Fixed-Point Tool and Prepare the System for Conversion” on page 42-55

“Collect Ranges and Convert to Fixed-Point” on page 42-55

“Verify Fixed-Point Behavior” on page 42-56

This example shows how to use the Fixed-Point Tool to propose fraction lengths for a model based on the minimum and maximum values captured over multiple simulations. In this example, you define a `Simulink.SimulationInput` object in the base or model workspace to specify the simulation scenarios to use for range collection. The Fixed-Point Tool merges the results from two simulation runs and proposes a data type based on the merged ranges. Merging results allows you to autoscale your model over the complete simulation range.

When converting a system based on multiple simulation scenarios, structurally altering the contents of the system under design during the conversion process could lead to errors. When defining the simulation scenarios avoid making any of the following changes to the system under design:

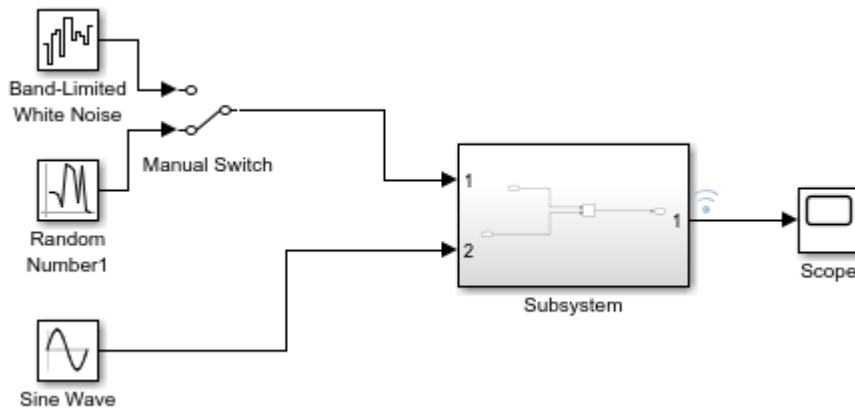
- Add or delete a block in the system under design
- Add another input to the system under design
- Change a block type in the system under design

Set up the Model

This example uses the `ex_merge_ranges` model. The model contains a sine wave input and two alternate noise sources, band-limited white noise and random uniform noise. In this example, define a `Simulink.SimulationInput` object and collect ranges using the Band-Limited White Noise source and the Random Number 1 source. Propose data types for the model based on the merged simulation ranges.

Open the model.

```
model = 'ex_merge_ranges';  
open_system(model);
```

Define the `Simulink.SimulationInput` object. The first object sets the Manual Switch block to the Band-Limited White Noise source, the second `SimulationInput` object sets the Manual Switch block to the Random Number source.

```
simIn(1) = Simulink.SimulationInput(model);
simIn(2) = Simulink.SimulationInput(model);

simIn(1) = simIn(1).setBlockParameter('ex_merge_ranges/Manual Switch', 'sw', '0');
simIn(2) = simIn(2).setBlockParameter('ex_merge_ranges/Manual Switch', 'sw', '1');
```

Open the Fixed-Point Tool and Prepare the System for Conversion

- 1 In the **Apps** gallery of the `ex_merge_ranges` model, select **Fixed-Point Tool**.
- 2 In the Fixed-Point Tool, click **New**, and select Iterative Fixed-Point Conversion.
- 3 Under **System Under Design**, select Subsystem.
- 4 Under **Range Collection Mode**, select **Simulation Ranges** as the range collection method.
- 5 Under **Simulation Inputs**, select the `Simulink.SimulationInput` object, `simIn` that you defined in the base workspace.
- 6 Set the absolute tolerance of the Subsystem: 1 signal to 0.1, or 10%.
- 7 In the toolstrip, click the **Prepare** button.

Collect Ranges and Convert to Fixed-Point

- 1 Expand the **Collect Ranges** button arrow and select Double precision. Click **Collect Ranges**.


Simulink simulates the `ex_merge_ranges` model twice, once using the Band-Limited White Noise source block and once using the Random Number source block.

You can view the ranges of each simulation individually by selecting the simulation in the **Workflow Browser**. In this example the `BaselineRun_Scenario_1` simulation had a **SimMin** value of -3.5821 and a **SimMax** value of 2.7598. The `BaselineRun_Scenario_2` simulation had a **SimMin** value of -2.5317 and a **SimMax** value of 3.1542.

Selecting the `BaselineRun` node in the **Workflow Browser** shows the merged ranges from the two simulation scenarios.


| WORKFLOW BROWSER | | Results | | | |
|-------------------|------------|-------------------------------|--------------------|--------------------|--|
| Name | CompiledDT | SpecifiedDT | SimMin | SimMax | |
| Add : Accumulator | double | Inherit: Inherit via inter... | -3.582183763440959 | 3.1542109234550813 | |
| Add : Output | double | Inherit: Inherit via inter... | -3.582183763440959 | 3.1542109234550813 | |

2

In the **Convert** section of the toolbar, click the **Propose Data Types** button .


The Fixed-Point Tool uses the merged minimum and maximum values to propose fraction lengths for each block. These values ensure maximum precision while spanning the full range of simulation values. The tool displays the proposed data types in the spreadsheet.

3

Click the **Apply Data Types** button  to write the proposed data types to the model.

Verify Fixed-Point Behavior

1

In the **Verify** section of the toolbar, click the **Simulate with Embedded Types** button . The Fixed-Point Tool simulates the model using the same `Simulink.SimulationInput` scenarios that were used to collect ranges and verifies whether each scenario met the specified tolerances.

The **Workflow Browser** indicates whether the verification runs met the tolerances. In this example, both simulation scenarios met the specified tolerances.

| WORKFLOW BROWSER | |
|------------------------|--|
| Setup | |
| Preparation Results | |
| BaselineRun | |
| EmbeddedRun | |
| EmbeddedRun_Scenario_1 | |
| EmbeddedRun_Scenario_2 | |

2 To view the simulation data for an individual run, right-click on the run in the **Workflow Browser**.

See Also

More About

- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9
- How to Convert Models to Fixed Point Using Multiple Simulation Scenarios

View Simulation Results

| In this section... |
|--|
| “Compare Runs” on page 42-57 |
| “Histogram Plot of Signal” on page 42-58 |

The Fixed-Point Tool uses the Simulation Data Inspector tool plotting capabilities that enable you to plot logged signals for graphical analysis. Using the Simulation Data Inspector to inspect and compare data after converting your floating-point model to fixed point facilitates tracking numerical error propagation.

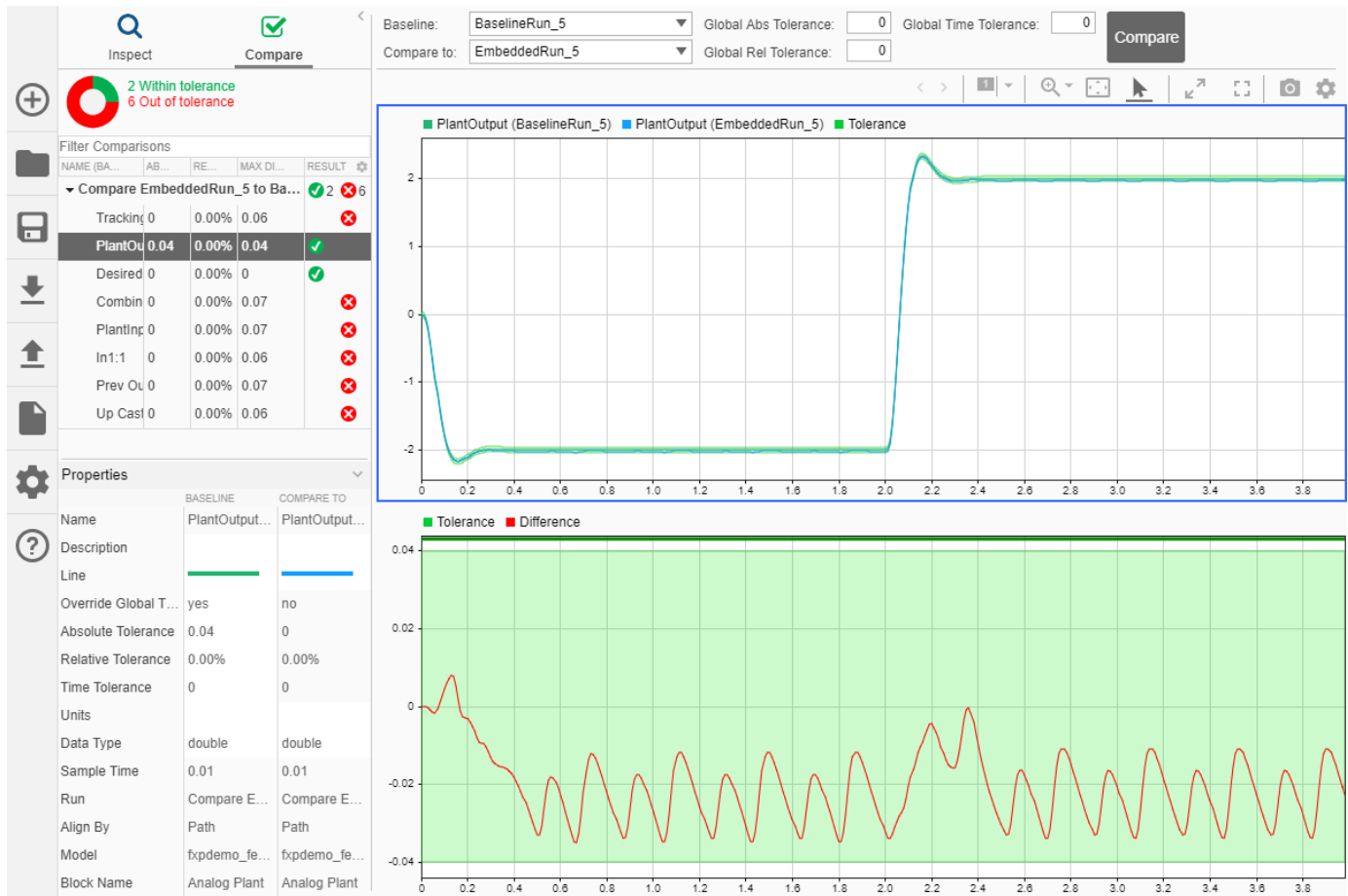
Use the Simulation Data Inspector to:

- Plot multiple signals in one or more axes
- Compare a signal in different runs
- Compare all logged signal data from different runs
- Export signal logging results to a MAT-file
- Specify tolerances for signal comparison
- Create a report of the current view and data in the Simulation Data Inspector

Compare Runs

To compare runs, in the Fixed-Point Tool, right-click on the embedded run and select **Open SDI**.

On the upper axes, the Simulation Data Inspector plots the signal for the selected runs and the tolerance, if specified. On the lower axes, the Simulation Data Inspector plots the difference between those runs.



Histogram Plot of Signal

To view the histogram plot of a signal, select the signal in the **Results** spreadsheet. The **Result Details** pane includes a histogram plot that helps you visualize the dynamic range of a signal. It provides information about the:

- Total number of samples (N).
- Maximum number of bits to prevent overflow.
- Number of times each bit has represented the data (as a percentage of the total number of samples).
- Number of times that exact zero occurred (without the effect of quantization). This number does not include the number of zeroes that occurred due to rounding.

You can use this information to estimate the word size required to represent the signal.

RESULT DETAILS

fxpdemo_feedback/Controller/In1

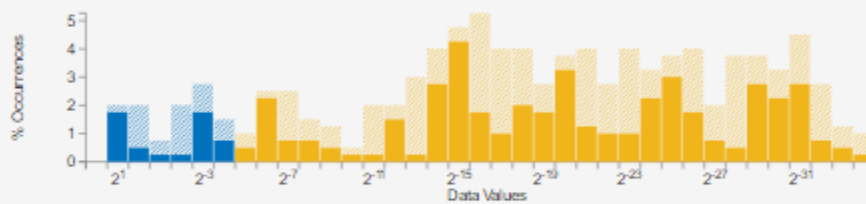
Proposed Data Type Summary

| Property | ProposedDT | SpecifiedDT |
|-----------|--------------|---------------|
| Data Type | fixdt(1,8,4) | Inherit: auto |
| Minimum | -8 | |
| Maximum | 7.9375 | |
| Precision | 0.0625 | |

Ranges used for proposal

| Property | Minimum | Maximum |
|-------------------|---------|--------------------|
| Shared Simulation | -2 | 3.9999999999711746 |

Visualization of Simulation Data



| | Potential Overflows | In-Range | Potential Underflows |
|-----------------|---------------------|----------|----------------------|
| Positive Values | 0 | 21 | 178 |
| Negative Values | 0 | 23 | 177 |

Number of times zero occurred: 0

Proposal Details

- There is a requirement for the data type of this result to match the data type of other results.
 - [Highlight Elements Sharing Same Data Type](#)

See Also

Related Examples

- “Propose Fraction Lengths Using Simulation Range Data” on page 42-45

Fixed-Point Instrumentation and Data Type Override

The conversion of a model from floating point to fixed point requires configuring fixed-point instrumentation and data type overrides. However, leaving these settings on after the conversion can lead to unexpected results.

The Fixed-Point Tool automatically enables fixed-point instrumentation when you click the **Collect Ranges** button in the tool. By default, the Fixed-Point Tool uses the current data type override set on the model. You can also choose to override data types with doubles, singles, or scaled doubles. When the simulation or derivation is complete, the tool automatically disables the instrumentation and removes the data type override, if data type override was selected in the tool. When you click the **Simulate with Embedded Types** button, the tool enables instrumentation during the simulation. Data type override settings on the model are not affected.

Control Instrumentation Settings

The fixed-point instrumentation mode controls which objects log minimum, maximum, and overflow data during simulation. Instrumentation is required to collect simulation ranges using the Fixed-Point Tool. These ranges are used to propose data types for the model. When you are not actively converting your model to fixed point, disable the fixed-point instrumentation to restore the maximum simulation speed to your model.

To enable instrumentation outside of the Fixed-Point Tool, at the command line set the `MinMaxOverflowLogging` parameter to `MinMaxAndOverflow` or `OverflowOnly`.

```
set_param('MyModel', 'MinMaxOverflowLogging', 'MinMaxAndOverflow')
```

Instrumentation requires a Fixed-Point Designer license. To disable instrumentation on a model, set the parameter to `ForceOff` or `UseLocalSettings`.

```
set_param('MyModel', 'MinMaxOverflowLogging', 'UseLocalSettings')
```

Control Data Type Override

Use data type override to simulate your model using double, single, or scaled double data types. If you do not have Fixed-Point Designer software, you can still configure data type override settings to simulate a model that specifies fixed-point data types. Using this setting, the software temporarily overrides data types with floating-point data types during simulation.

```
set_param('MyModel', 'DataTypeOverride', 'Double')
```

To observe the true behavior of your model, set the data type override parameter to `UseLocalSettings` or `Off`.

```
set_param('MyModel', 'DataTypeOverride', 'Off')
```

Instrumentation Settings and Data Type Override for a Model Reference Hierarchy

When you simulate a model that contains referenced models, the data type override and fixed-point instrumentation settings for the top-level model do not control the settings for the referenced models.

You must specify these settings separately for the referenced model. If the settings are inconsistent, for example, if you set the top-level model data type override setting to double and the referenced model to use local settings and the referenced model uses fixed-point data types, data type propagation issues might occur.

When you change the fixed-point instrumentation and data type override settings for any instance of a referenced model, the settings change on all instances of the model and on the referenced model itself.

Data Type Override Limitations

For these blocks, data type override is not supported.

- Stateflow Chart blocks that use MATLAB as the action language
- State Transition Table blocks that use MATLAB as the action language
- Truth Table blocks that use MATLAB as the action language
- Test Sequence block
- MATLAB Function block
- MATLAB Discrete-Event System block
- Requirements Table block

See Also

More About

- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9
- “Use Custom Data Type Override Settings for Range Collection” on page 44-9

Model Configuration Changes Made During Data Type Optimization

Data type optimization seeks to minimize an objective function while maintaining the original system behavior within a specified tolerance. You can optimize data types by using `fxpopt` at the command line, or by using the Optimized Fixed-Point Conversion workflow in the **Fixed-Point Tool**.

During the optimization process, the software changes the settings and model configuration parameters described below. You can restore these parameter settings after the optimization is complete in the Configuration Parameters dialog box or by using the `set_param` function. You can also use the `KeepOriginalModelParameters` option of `explore` to maintain the original values of model parameters.

| Reason for Parameter Change | Parameter | Default Value | Optimization Changes Value to |
|--|--|---------------|-------------------------------|
| Suppress Diagnostics | ParameterDowncastMsg on page 42-64 | 'error' | 'none' |
| | ParameterUnderflowMsg on page 42-65 | 'none' | 'none' |
| | FixptConstUnderflowMsg on page 42-65 | 'none' | 'none' |
| | ParameterPrecisionLossMsg on page 42-65 | 'none' | 'none' |
| | FixptConstPrecisionLossMsg on page 42-65 | 'none' | 'none' |
| | ParameterOverflowMsg on page 42-65 | 'error' | 'none' |
| | FixptConstOverflowMsg on page 42-65 | 'none' | 'none' |
| | IntegerOverflowMsg on page 42-65 | 'warning' | 'none' |
| | IntegerSaturationMsg on page 42-66 | 'warning' | 'none' |
| Logging with the Simulation Data Inspector | SignalLogging on page 42-66 | 'on' | 'on' |
| | ReturnWorkspaceOutputs on page 42-66 | 'on' | 'on' |
| | SaveFormat on page 42-66 | 'Dataset' | 'Dataset' |
| Reduce memory consumption of result | SaveTime on page 42-66 | 'on' | 'off' |
| | SaveOutput on page 42-66 | 'on' | 'off' |

| Reason for Parameter Change | Parameter | Default Value | Optimization Changes Value to |
|-----------------------------|-----------------------------------|--------------------|-------------------------------|
| Model validity | SignalRangeChecking on page 42-67 | 'none' | 'error' |
| Understand result | ShowPortDataTypes on page 42-67 | 'off' | 'on' |
| Accelerate optimization | SimulationMode on page 42-67 | 'normal' | 'accelerator' |
| Data type override | DataTypeOverride on page 42-67 | 'UseLocalSettings' | 'Off' |

You can use the `showContents` method of the `OptimizationSolution` object to print a summary of the changes made during data type optimization to the MATLAB Command Window. For example, after optimizing data types according to the example “Optimize Fixed-Point Data Types” use `showContents` to view the model parameter changes that were made during data type optimization:

```
solution = result.Solutions(1);
showContents(solution)
```

```
ModelName: 'ex_auto_gain_controller'
```

```
ModelParameters:
  Index          Name                               Value
-----
  1      SignalLogging                    'on'
  2      ReturnWorkspaceOutputs            'on'
  3      SaveFormat                        'Dataset'
  4      ShowPortDataTypes                  'on'
  5      SignalRangeChecking                'error'
  6      ParameterDowncastMsg              'none'
  7      ParameterUnderflowMsg             'none'
  8      ParameterPrecisionLossMsg          'none'
  9      ParameterOverflowMsg               'none'
 10      FixptConstPrecisionLossMsg         'none'
 11      FixptConstOverflowMsg              'none'
 12      FixptConstUnderflowMsg             'none'
 13      IntegerOverflowMsg                 'none'
 14      IntegerSaturationMsg                'none'
 15      SaveTime                           'off'
 16      SaveOutput                         'off'
 17      SimulationMode                      'accelerator'
 18      DataTypeOverride                   'off'
```

```
...
```

Detect downcast

Take no action when a parameter downcast occurs during simulation. Quantization effects, including downcasts, are expected during optimization. Optimized data types will meet all specified behavioral constraints. For more information, see “Detect downcast”.

Parameter: `ParameterDowncastMsg`

Value: 'none'

Default: 'error'

Detect underflow

Take no action when parameter quantization causes a non-zero value to underflow to zero during simulation. Take no action when a fixed-point constant underflow occurs during simulation. Quantization effects, including underflow, are expected during optimization. Optimized data types will meet all specified behavioral constraints. For more information, see “Detect underflow” and “Detect underflow”.

Parameter: ParameterUnderflowMsg

Value: 'none'

Default: 'none'

Parameter: FixptConstUnderflowMsg

Value: 'none'

Default: 'none'

Detect precision loss

Take no action when parameter precision loss or fixed-point constant precision loss occurs during simulation. Quantization effects, including precision loss, are expected during optimization. Optimized data types will meet all specified behavioral constraints. For more information, see “Detect precision loss” and “Detect precision loss”.

Parameter: ParameterPrecisionLossMsg

Value: 'none'

Default: 'warning'

Parameter: FixptConstPrecisionLossMsg

Value: 'none'

Default: 'none'

Detect overflow

Take no action if a parameter overflow or fixed-point constant overflow occurs during simulation. Quantization effects, including overflow, are expected during optimization. Optimized data types will meet all specified behavioral constraints. For more information, see “Detect overflow” and “Detect overflow”.

Parameter: ParameterOverflowMsg

Value: 'none'

Default: 'error'

Parameter: FixptConstOverflowMsg

Value: 'none'

Default: 'none'

Wrap on overflow

Take no action if the value of a signal overflows the signal data type and wraps around. Quantization effects, including overflow, are expected during optimization. Optimized data types will meet all specified behavioral constraints. For more information, see “Wrap on overflow”.

Parameter: IntegerOverflowMsg

Value: 'none'

Default: 'warning'

Saturate on overflow

Take no action if the value of a signal is too large to be represented by the signal data type, resulting in saturation. Quantization effects, including overflow, are expected during optimization. Optimized data types will meet all specified behavioral constraints. For more information, see “Saturate on overflow”.

Parameter: IntegerSaturationMsg

Value: 'none'

Default: 'warning'

Signal logging

Globally enable signal logging to the workspace for this model. Optimization requires this setting to log specified signals to the Simulation Data Inspector. For more information, see Signal logging.

Parameter: SignalLogging

Value: 'on'

Default: 'on'

Single simulation output

Enable the single-output format of the `sim` command to return the simulation result as a `Simulink.SimulationOutput` object. Optimization requires this setting to log signals with the Simulation Data Inspector. For more information, see Single simulation output.

Parameter: ReturnWorkspaceOutputs

Value: 'on'

Default: 'on'

Format

Store each logged state and output in a `Simulink.SimulationData.Dataset` object. Optimization requires this setting for signal logging with the Simulation Data Inspector. For more information, see “Log Data to Persistent Storage”.

Parameter: SaveFormat

Value: 'Dataset'

Default: 'Dataset'

Time

Do not export time data to the MATLAB workspace during simulation. Optimization requires this setting to avoid unnecessary memory consumption. For more information, see Time.

Parameter: SaveTime

Value: 'off'

Default: 'on'

Output

Do not export root output signal data to a specified MATLAB variable during simulation. Optimization requires this setting for signal logging with the Simulation Data Inspector. For more information, see Output.

Parameter: SaveOutput

Value: 'off'

Default: 'on'

Simulation range checking

Terminate the simulation when signals exceed specified minimum or maximum values. Optimization requires this setting to ensure that the simulation ranges of the model with optimized data types applied honors the specified design ranges. For more information, see “Simulation range checking”.

Parameter: SignalRangeChecking

Value: 'error'

Default: 'none'

Show port data types

Show data types of ports on the model block diagram. Optimization requires this setting to allow you to easily inspect the optimized data types applied to your model. For more information, see “Programmatic Model Editor Appearance Parameters”.

Parameter: ShowPortDataTypes

Value: 'on'

Default: 'off'

Simulation mode

Simulate the model in accelerator mode. Optimization uses accelerator mode to reduce the amount of time required to optimize data types on your model. For more information, see “What Is Acceleration?”.

Parameter: SimulationMode

Value: 'accelerator'

Default: 'normal'

Data type override

By default, optimization turns off any data type override set on your model so that the effect of optimized data types on model behavior is accurately represented during simulation. You can customize this behavior by using the advanced options of `fxpOptimizationOptions`. For more information, see “Fixed-Point Instrumentation and Data Type Override” on page 42-61.

Parameter: DataTypeOverride

Value: 'Off'

Default: 'UseLocalSettings'

See Also

`fxpopt` | “Specify Behavioral Constraints” on page 42-18

Related Examples

- “Optimize Fixed-Point Data Types for a System” on page 40-14
- “Optimize the Fixed-Point Data Types of a System Using the Fixed-Point Tool” on page 40-40

Range Analysis

- “How Range Analysis Works” on page 43-2
- “Derive Ranges at the Subsystem Level” on page 43-6
- “Derive Ranges Using Design Ranges” on page 43-8
- “Derive Ranges Using Block Initial Conditions” on page 43-10
- “Derive Ranges for Simulink.Parameter Objects” on page 43-12
- “Insufficient Design Range Information” on page 43-14
- “Troubleshoot Range Analysis of System Objects” on page 43-16
- “Provide More Design Range Information” on page 43-19
- “Fix Design Range Conflicts” on page 43-22
- “Intermediate Range Results” on page 43-24
- “Unsupported Simulink Software Features” on page 43-27
- “Simulink Blocks Supported for Range Analysis” on page 43-28

How Range Analysis Works

In this section...

- “Analyzing a Model with Range Analysis” on page 43-2
- “Automatic Stubbing” on page 43-4
- “Model Compatibility with Range Analysis” on page 43-4
- “How to Derive Ranges” on page 43-4

Analyzing a Model with Range Analysis

The model that you want to analyze **must** be compatible with range analysis. If your model is not compatible, either replace unsupported blocks or divide the model so that you can analyze the parts of the model that are compatible. For more information, see “Model Compatibility with Range Analysis” on page 43-4.

When you specify **Derived ranges** as the range collection mode, the Fixed-Point Designer software performs a static range analysis of your model to derive minimum and maximum range values for signals in the model. The software analyzes the model behavior and computes the values that can occur during simulation for each block Output. The range of these values is called a derived range.

The software statically analyzes the ranges of the individual computations in the model based on:

- Specified design ranges, known as design minimum and maximum values, for example, minimum and maximum values specified for:
 - Inport and Outport blocks
 - Block outputs
 - Input, output, and local data used in MATLAB Function and Stateflow Chart blocks
 - Simulink data objects (`Simulink.Signal` and `Simulink.Parameter` objects)
- Inputs
- The semantics of each calculation in the blocks

If the model contains objects that the analysis cannot support, where possible, the software uses automatic stubbing on page 43-4. For more information, see “Automatic Stubbing” on page 43-4.

The range analysis tries to narrow the derived range by using all the specified design ranges in the model. The more design range information you specify, the more likely the range analysis is to succeed. As the software performs the analysis, it derives new range information for the model. The software then attempts to use this new information, together with the specified ranges, to derive ranges for the remaining objects in the model.

For models that contain floating-point operations, range analysis might report a range that is slightly larger than expected. This difference is due to rounding errors. The software approximates floating-point numbers with infinite-precision rational numbers for analysis and then converts to floating point for reporting.

The following table summarizes how the analysis derives range information and provides links to examples.

| When... | How the Analysis Works | Examples |
|---|--|---|
| You specify design minimum and maximum data for a block output. | <p>The derived range at the block output is based on these specified values and on the following values for blocks connected to its inputs and outputs:</p> <ul style="list-style-type: none"> • Specified minimum and maximum values • Derived minimum and maximum values <p>Only block output signals participate in derived range analysis. If a block has additional data type controls, such as for the accumulator or intermediate results, ranges are not derived for these elements.</p> | “Derive Ranges Using Design Ranges” on page 43-8 |
| A parameter on a block has initial conditions and a design range. | The analysis takes both factors into account by taking the union of the design range and the initial conditions. | “Derive Ranges Using Block Initial Conditions” on page 43-10 |
| The model contains a parameter with a specified range and the parameter storage class is set to Auto. | The analysis does not take into account the range specified for the parameter. Instead, it uses the parameter value. | “Derive Ranges for Simulink.Parameter Objects” on page 43-12 |
| The model contains a parameter with a specified range and the parameter storage class is not set to Auto. | The analysis takes into account the range specified for the parameter and ignores the value. | “Derive Ranges for Simulink.Parameter Objects” on page 43-12 |
| The model contains insufficient design range information. | The analysis cannot determine derived ranges. Specify more design range information and rerun the analysis. | <p>“Troubleshoot Range Analysis of System Objects” on page 43-16</p> <p>The range analysis results might depend on the block sorted order, which determines the order in which the software analyzes the blocks. For more information, see “Control and Display Execution Order”.</p> |

| When... | How the Analysis Works | Examples |
|--|---|--|
| The model contains conflicting design range information. | The analysis cannot determine the derived minimum or derived maximum value for an object. The Fixed-Point Tool generates an error. To fix this error, examine the design ranges specified in the model to identify inconsistent design specifications. Modify them to make them consistent. | "Fix Design Range Conflicts" on page 43-22 |

Automatic Stubbing

What is Automatic Stubbing?

Automatic stubbing is when the software considers only the interface of the unsupported objects in a model, not their actual behavior. Automatic stubbing lets you analyze a model that contains objects that the Fixed-Point Designer software does not support. However, if any unsupported model element affects the derivation results, the analysis might achieve only partial results.

How Automatic Stubbing Works

With automatic stubbing, when the range analysis comes to an unsupported block, the software ignores ("stubs") that block. The analysis ignores the behavior of the block. As a result, the block output can take any value.

The software cannot "stub" all Simulink blocks, such as the Integrator block. See the blocks marked "not stubbable" in "Simulink Blocks Supported for Range Analysis" on page 43-28.

Model Compatibility with Range Analysis

To verify that your model is compatible with range analysis, see:

- "Unsupported Simulink Software Features" on page 43-27
- "Simulink Blocks Supported for Range Analysis" on page 43-28
- "Limitations of Support for Model Blocks" on page 43-35

How to Derive Ranges

- 1 Verify that your model is compatible with range analysis.
- 2 In Simulink, open your model and set it up for use with the Fixed-Point Tool. For more information, see Set Up the Model on page 42-13.
- 3 From the Simulink **Apps** tab, select **Fixed-Point Tool**.
- 4 In the Fixed-Point Tool, under **New**, select the **Iterative Fixed-Point Conversion** workflow.
- 5 Under **System Under Design (SUD)**, select the system or subsystem of interest.
- 6 Under **Range Collection Mode**, select **Derived ranges** as the method of range collection. This configures the model to collect ranges using idealized floating-point data types.

By default, the tool collects ranges using design information from the system under design. For more information, see “Derive Ranges at the Subsystem Level” on page 43-6.

- 7 Click **Prepare** to have the Fixed-Point Tool check the system under design for compatibility with the conversion process and report any issues found in the model.

The Fixed-Point Tool:

- Checks the model against fixed-point guidelines.
- Identifies unsupported blocks.
- Identifies blocks that need design range information.

- 8 Click the **Collect Ranges** button to run the analysis.

The analysis tries to derive range information for objects in the selected system under design. Your next steps depend on the analysis results.

| Analysis Results | Fixed-Point Tool Behavior | Next Steps | For More Information |
|---|--|---|---|
| Successfully derives range data for the model. | Displays the derived minimum and maximum values for the blocks in the selected system. | Review the derived ranges to determine if the results are suitable for proposing data types. If not, you must specify additional design information and rerun the analysis. | “Derive Ranges Using Design Ranges” on page 43-8 |
| Fails because the model contains blocks that the software does not support. | Generates an error and provides information about the unsupported blocks. | To fix the error, review the error message information and replace the unsupported blocks. | “Model Compatibility with Range Analysis” on page 43-4 |
| Cannot derive range data because the model contains conflicting design range information. | Generates an error. | To fix this error, examine the design ranges specified in the model to identify inconsistent design specifications. Modify the design ranges to make them consistent. | “Fix Design Range Conflicts” on page 43-22 |
| Cannot derive range data for an object because there is insufficient design range information specified on the model. | Highlights the results for the object. | Examine the model to determine which design range information is missing. | “Troubleshoot Range Analysis of System Objects” on page 43-16 |

Derive Ranges at the Subsystem Level

In this section...

“When to Derive Ranges at the Subsystem Level” on page 43-6

“Derive Ranges at the Subsystem Level” on page 43-6

You can derive range information for individual atomic subsystems and atomic charts. When you derive ranges at the model level, the software takes into account all information in the scope of the model. When you derive ranges at the subsystem level only, the software treats the subsystem as a standalone unit and the derived ranges are based on only the local design range information specified in the subsystem or chart. Therefore, when you derive ranges at the subsystem level, the analysis results might differ from the results of the analysis at the model level.

For example, consider a subsystem that has an input with a design minimum of -10 and a design maximum of 10 that is connected to an input signal with a constant value of 1 . When you derive ranges at the model level, the range analysis software uses the constant value 1 as the input. When you derive ranges at the subsystem level, the range analysis software does not take the constant value into account and instead uses $[-10..10]$ as the range.

When to Derive Ranges at the Subsystem Level

Derive ranges at the subsystem level to facilitate:

- System validation

It is a best practice to analyze individual subsystems in your model one at a time. This practice makes it easier to understand the atomic behavior of the subsystem. It also makes debugging easier by isolating the source of any issues.

- Calibration

The results from the analysis at subsystem level are based only on the settings specified within the subsystem. The proposed data types cover the full intended design range of the subsystem. Based on these results, you can determine whether you can reuse the subsystem in other parts of your model.

Derive Ranges at the Subsystem Level

The complete procedure for deriving ranges is described in “How to Derive Ranges” on page 43-4.

To derive ranges at the subsystem level, the key points to remember are:

- The subsystem or subchart must be atomic.

An *atomic subsystem* executes as a unit relative to the parent model. Atomic subsystem block execution does not interleave with parent block execution. You can extract atomic subsystems for use as standalone models.

- In the Fixed-Point Tool, under **System Under Design (SUD)**, select the subsystem of interest.
- Under **Range Collection Mode**, select **Derived ranges** as the method of range collection.

See Also

More About

- “How Range Analysis Works” on page 43-2

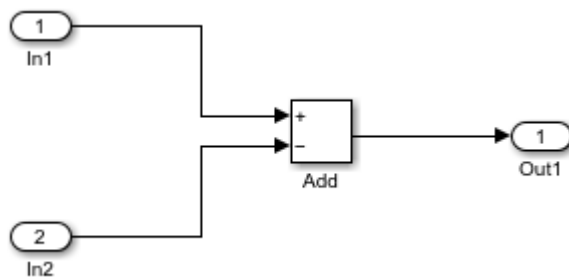
Derive Ranges Using Design Ranges

This example shows how the range analysis narrows the derived range for the Outport block. This range is based on the range derived for the Add block using the design ranges specified on the two Inport blocks and the design range specified for the Add block.

Open the Model and View Design Ranges

Open the model. At the MATLAB command line, enter:

```
open_system('ex_derived_min_max_1')
```



Update the diagram to display the specified design minimum and maximum values for each block.

- In1 design range is $[-50 \dots 100]$.
- In2 design range is $[-50 \dots 35]$.
- Add block design range is $[-125 \dots 55]$.

Derive Ranges

- 1 From the Simulink **Apps** tab, select **Fixed-Point Tool**.
- 2 In the Fixed-Point Tool, under **New** workflow, select **Iterative Fixed-Point Conversion**.
- 3 Under **System Under Design (SUD)**, select `ex_derived_min_max_1` as the system you want to convert.
- 4 Under **Range Collection Mode**, select **Derived ranges**.
- 5 Click the **Collect Ranges** button.

To calculate the derived range at the Add block input, the software uses the design minimum and maximum values specified for the Inport blocks, $[-50 \dots 100]$ and $[-50 \dots 35]$. The derived range at the Add block input is $[-85 \dots 150]$.

When the analysis is complete, the Fixed-Point Tool displays the derived and design minimum and maximum values for the blocks in the selected system in the spreadsheet.

| Results | | | | | | | |
|--------------|------------|----------------------|-----------|-----------|------------|------------|--|
| Name | CompiledDT | SpecifiedDT | DesignMin | DesignMax | DerivedMin | DerivedMax | |
| Add : Output | double | Inherit: Inherit ... | -125 | 55 | -85 | 55 | |
| In1 | double | Inherit: auto | -50 | 100 | -50 | 100 | |
| In2 | double | Inherit: auto | -50 | 35 | -50 | 35 | |
| Out1 | | Inherit: auto | | | -85 | 55 | |

- The derived range for the Add block output signal is narrowed to $[-85..55]$. This derived range is the intersection of the range derived from the block inputs, $[-85..150]$, and the design minimum and maximum values specified for the block output, $[-125..55]$.

Note The accumulator in the Add block does not participate in derived range analysis. Ranges are derived only for block output signals.

- The derived range for the Outport block Out1 is $[-85..55]$, the same as the Add block output.

Tip To display design ranges in your model, in the **Debug** tab, select **Information Overlays > Signal Data Ranges**.

See Also

Related Examples

- “Insufficient Design Range Information” on page 43-14
- “Troubleshoot Range Analysis of System Objects” on page 43-16
- “How Range Analysis Works” on page 43-2

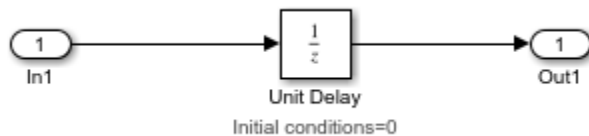
Derive Ranges Using Block Initial Conditions

This example shows how range analysis takes into account block initial conditions.

Open the Model

Open the model. At the MATLAB command line, enter:

```
open_system('ex_derived_min_max_2')
```



The model uses information overlays to display the specified design minimum and maximum values for the Inport block, and block annotations to display the initial conditions for the Unit Delay block.

- In1 design range is [5 . . 10].
- Unit Delay block initial condition is 0.

Derive Ranges

- 1 From the Simulink **Apps** tab, select **Fixed-Point Tool**.
- 2 In the Fixed-Point Tool, under **New** workflow, select **Iterative Fixed-Point Conversion**.
- 3 Under **System Under Design (SUD)**, select `ex_derived_min_max_2` as the system you want to convert.
- 4 Under **Range Collection Mode**, select **Derived ranges**.
- 5 Click the **Collect Ranges** button.

In the spreadsheet, the Fixed-Point Tool displays the derived minimum and maximum values for the blocks in the model.

The derived minimum and maximum range for the Output block, Out1, is [0 . . 10]. The range analysis derives this range by taking the union of the initial value, 0, on the Unit Delay block and the design range on the block, [5 . . 10].

- 6 Change the initial condition of the Unit Delay block to 7.
 - a In the model, double-click the Unit Delay block.
 - b In the **Block Parameters** dialog box, set **Initial condition** to 7, then click **OK**.
 - c In the Fixed-Point Tool, click the **Collect Ranges** button.

Because the analysis takes the union of the initial conditions, 7, and the design range, [5 . . 10], on the Unit Delay block, the derived range for the Output block is now [5 . . 10].

Tip To display design ranges in your model, in the **Debug** tab, select **Information Overlays > Signal Data Ranges**.

See Also

More About

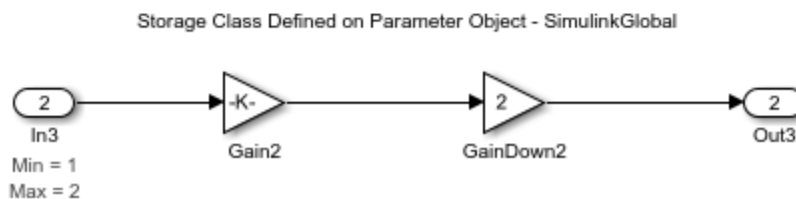
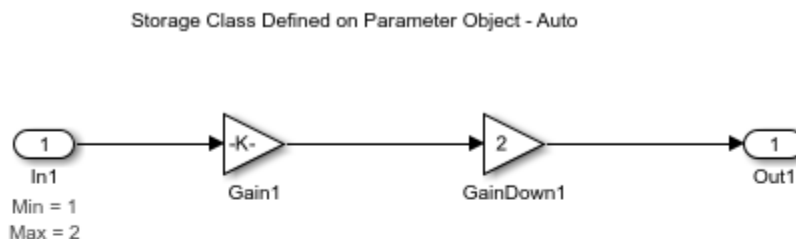
- “How Range Analysis Works” on page 43-2

Derive Ranges for Simulink.Parameter Objects

This example shows how the range analysis takes into account design range information for Simulink.Parameter objects unless the parameter storage class is set to Auto. If the parameter storage class is set to Auto, the analysis uses the value of the parameter.

Open the ex_derived_min_max_3 Model

```
open_system("ex_derived_min_max_3.slx")
```



The model displays the specified design minimum and maximum values for the Inport blocks. The design range for both Inport blocks is [1 . . . 2].

To display design ranges in your model, in the **Debug** tab, select **Information Overlays > Signal Data Ranges**.

Examine Gain Parameters

- 1 Double-click each Gain block and note the name of the Gain parameter on the Main tab.

| Gain Block | Gain Parameter |
|------------|----------------|
| Gain1 | paramObjOne |
| Gain2 | paramObjTwo |

- 2 In the **Modeling** tab, expand the **Design** gallery and select **Model Explorer**.
- 3 In **Model Explorer** window, select the base workspace and view information for each of the gain parameters used in the model.

| Gain Parameter | Type Information | Value | Storage Class |
|----------------|---------------------------|-------|---------------|
| paramObjOne | Simulink.Parameter object | 2 | Auto |
| paramObjTwo | Simulink.Parameter object | 2 | Model default |

Derive Ranges

- 1 From the Simulink **Apps** tab, select **Fixed-Point Tool**.
- 2 In the Fixed-Point Tool, under **New** workflow, select **Iterative Fixed-Point Conversion**.
- 3 In the Fixed-Point Tool, under **System Under Design (SUD)**, select `ex_derived_min_max_3` as the system you want to convert.
- 4 Under **Range Collection Mode**, select **Derived ranges**.
- 5 Click the **Collect Ranges** button.

When the analysis is finished, the Fixed-Point Tool displays the derived minimum and maximum values for the blocks in the model in the spreadsheet.

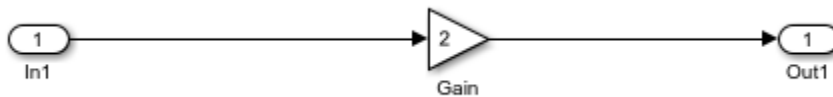
| Block | Derived Range | Reason |
|-------|---------------|---|
| Gain1 | [2..4] | The gain parameter, paramObjOne, specified on Gain block Gain1 is a Simulink.Parameter object that has its storage class specified as Auto. The range analysis uses the Value property of the Simulink.Parameter object, whose value is 2, and ignores the design range specified for these parameters. |
| Gain2 | [1..20] | The gain parameter, paramObjTwo, specified on Gain block Gain2 is a Simulink.Parameter object that has its storage class specified as Model default. The range analysis takes into account the design range, [1..10], specified for this parameter. |

Insufficient Design Range Information

This example shows that if the analysis cannot derive range information because there is insufficient design range information, you can fix the issue by providing additional input design minimum and maximum values.

Open the `ex_derived_min_max_4` Model

```
open_system("ex_derived_min_max_4")
```



The model displays the specified design minimum and maximum values for the blocks in the model.

- The Inport block In1 has a design minimum but no specified maximum value, as shown by the annotation, [-1. .].
- The Gain block has a design range of [-1.5. .1.5].
- The Outport block Out1 has no design range specified.

To display design ranges in your model, on the **Debug** tab, select **Information Overlays > Signal Data Ranges**.

Collect Ranges

- 1 From the Simulink **Apps** tab, select **Fixed-Point Tool**.
- 2 In the Fixed-Point Tool, under **New** workflow, select **Iterative Fixed-Point Conversion**.
- 3 In the Fixed-Point Tool, under **System Under Design (SUD)**, select `ex_derived_min_max_4` as the system you want to convert.
- 4 Under **Range Collection Mode**, select **Derived ranges**.
- 5 Click the **Collect Ranges** button.

The Fixed-Point Tool displays the derived minimum and maximum values for the blocks in the model. The range analysis is unable to derive a maximum value for the Inport block, In1. The tool highlights this result.

| Results | | | | | | | |
|---------|------|------------|----------------------|-----------|-----------|------------|------------|
| | Name | CompiledDT | SpecifiedDT | DesignMin | DesignMax | DerivedMin | DerivedMax |
| 📁 | Gain | double | Inherit: Inherit ... | -1.5 | 1.5 | -1.5 | 1.5 |
| 📁 | In1 | double | double | -1 | | -1 | Inf |
| 📁 | Out1 | | Inherit: auto | | | -1.5 | 1.5 |

Fix Insufficient Design Ranges

- 1 To fix the issue, specify a design maximum value for In1:
 - a In the model, double-click the Inport block, In1.
 - b In the block parameters dialog box, select the **Signal Attributes** tab.
 - c In this tab, set **Maximum** to 1 and click **OK**. To update the diagram, press (Ctrl + D).

The model displays the updated maximum value in the block annotation for In1, [-1..1].

- 2 Clear previously collected ranges and rerun the range analysis.
 - a In the Fixed-Point Tool, under **New** workflow, select **Range Collection**.

Changing workflows clears range data collected during the active workflow.

- b Switch back to the **Iterative Fixed-Point Conversion** workflow.
- c Select **Derived ranges** as the range collection mode.
- d Click the **Collect Ranges** button again to rerun the range analysis.

The range analysis can now derive ranges for the Inport and Gain blocks.

| Block | Derived Range | Reason |
|--------------|---------------|--|
| Inport In1 | [-1..1] | Uses specified design range on the block. |
| Gain | [-1.5..1.5] | The design range specified on the Gain block is [-1.5..1.5]. The derived range at the block input is [-1..1] (the derived range at the output of In1). Therefore, because the gain is 2, the derived range at the Gain block output is the intersection of the propagated range, [-2..2], and the design range, [-1.5..1.5]. |
| Outport Out1 | [-1.5..1.5] | Same as Gain block output because there is no locally specified design range on Outport block. |

See Also

Related Examples

- “Troubleshoot Range Analysis of System Objects” on page 43-16

Troubleshoot Range Analysis of System Objects

When deriving ranges for a model that uses a system object, the analysis fails if the model contains variables that can refer to multiple handle objects. The following example shows how to reconfigure the code so that the Fixed-Point Tool can derive ranges for the model.

In this example, range analysis of the first model `ex_HandleVariableRefersToMultipleObjects` produces an error because there is a variable in the code that can refer to different system objects depending on other conditions. The model `ex_HandleVariableRefersToSingleObject` is a rewrite of the first model with the same functionality, but the Fixed-Point Tool is able to derive ranges for the model.

Derive Ranges for Model with System Object

Open the first model, `ex_HandleVariableRefersToMultipleObjects`.

```
open_system("ex_HandleVariableRefersToMultipleObjects.slx")
```

The code inside the MATLAB Function block refers to the custom System Object `fAddConstant`.

```
function y = fcn(u, c)
    %#codegen

    persistent hSysObjAddTen
    persistent hSysObjAddNegTen
    persistent hSysObjForStep

    if isempty(hSysObjAddTen)
        hSysObjAddTen = fAddConstant(10);
    end

    if isempty(hSysObjAddNegTen)
        hSysObjAddNegTen = fAddConstant(-10);
    end

    if c > 0
        hSysObjForStep = hSysObjAddTen;
    else
        hSysObjForStep = hSysObjAddNegTen;
    end

    y = step(hSysObjForStep, u);
```

From the Simulink® **Apps** tab, select Fixed-Point Tool.

In the Fixed-Point Tool, select the **Iterative Fixed-Point Conversion** workflow.

Under **System Under Design (SUD)**, select `ex_HandleVariableRefersToMultipleObjects` as the system you want to convert.

Under **Range Collection Mode**, select **Derived ranges**.

Click the **Collect Ranges** button. The analysis fails because there is a handle variable in the code that can refer to different system objects depending on the value of `c`.

Rewrite Code to Enable Derived Range Analysis

You can rewrite the code inside the MATLAB Function block so that the Fixed-Point Tool is able to derive ranges for the System Object.

Close the Fixed-Point Tool and the `ex_HandleVariableRefersToMultipleObjects` model. Open the `ex_HandleVariableRefersToSingleObject` model.

```
open_system("ex_HandleVariableRefersToSingleObject.slx")
```

This model contains the rewritten code:

```
function y = fcn(u, c)
%#codegen

persistent hSysObjAddTen
persistent hSysObjAddNegTen

if isempty(hSysObjAddTen)
    hSysObjAddTen = fAddConstant(10);
end

if isempty(hSysObjAddNegTen)
    hSysObjAddNegTen = fAddConstant(-10);
end

if c > 0
    y = step(hSysObjAddTen, u);
else
    y = step(hSysObjAddNegTen, u);
end
```

From the Simulink® **Apps** tab, select Fixed-Point Tool.

In the Fixed-Point Tool, select the **Iterative Fixed-Point Conversion** workflow.

Under **System Under Design (SUD)**, select `ex_HandleVariableRefersToSingleObject` as the system you want to convert.

Under **Range Collection Mode**, select **Derived ranges**.

Click the **Collect Ranges** button. This time, the Fixed-Point Tool successfully derives ranges for the variables used in the model.

| Results | | | | | | | |
|----------------------------|------------|-------------------|-----------|-----------|------------|----------|----------|
| Name | CompiledDT | SpecifiedDT | DesignMin | DesignMax | DerivedMin | DerivedM | DerivedM |
| In1 | double | double | -10 | 10 | -10 | 10 | 10 |
| In2 | double | double | -10 | 10 | -10 | 10 | 10 |
| MATLAB Function.y | double | Inherit: Same ... | | | -20 | 20 | 20 |
| MATLAB Function/fAddCon... | double | | | | -10 | 10 | 10 |
| MATLAB Function/fAddCon... | double | | | | -10 | 10 | 10 |
| MATLAB Function/fAddCon... | double | | | | -20 | 20 | 20 |
| MATLAB Function/fcn : c | double | | | | -10 | 10 | 10 |
| MATLAB Function/fcn : u | double | | | | -10 | 10 | 10 |
| MATLAB Function/fcn : y | double | | | | -20 | 20 | 20 |
| Out1 | | Inherit: auto | | | -20 | 20 | 20 |

See Also

Fixed-Point Tool | “How Range Analysis Works” on page 43-2

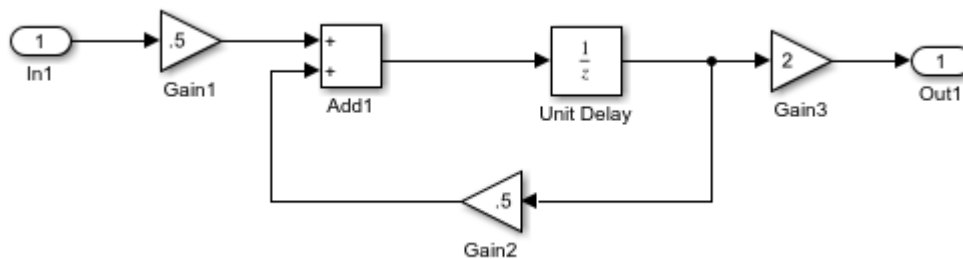
Provide More Design Range Information

This example shows that if the analysis cannot derive range information because there is insufficient design range information, you can fix the issue by providing additional design range information.

Open Model

Open the `ex_derived_min_max_5` model.

```
open_system("ex_derived_min_max_5.slx")
```



The model displays the specified design minimum and maximum values for the blocks in the model.

- The Inport block In1 has a design range of $[-10..20]$.
- The rest of the blocks in the model have no specified design range.

To display design ranges in your model, in the **Debug** tab, select **Information Overlays > Signal Data Ranges**.

Collect Ranges in the Fixed-Point Tool

- 1 From the Simulink **Apps** tab, select **Fixed-Point Tool**.
- 2 In the Fixed-Point Tool, under **New** workflow, select **Iterative Fixed-Point Conversion**.
- 3 In the Fixed-Point Tool, under **System Under Design (SUD)**, select `ex_derived_min_max_5` as the system you want to convert.
- 4 Under **Range Collection Mode**, select **Derived ranges**.
- 5 Click the **Collect Ranges** button.

When the analysis is complete, the Fixed-Point Tool displays the derived minimum and maximum values for the blocks in the model in the spreadsheet. Because the model contains a feedback loop, the analysis is unable to derive an output range for the Add block or for any of the blocks connected to this output. The Fixed-Point Tool highlights these results.

| Results | | | | | | | |
|---------|---------------|------------|----------------------|-----------|-----------|------------|------------|
| | Name ▲ | CompiledDT | SpecifiedDT | DesignMin | DesignMax | DerivedMin | DerivedMax |
| ☐ | Add1 : Output | double | double | | | -Inf | Inf |
| ☐ | Gain1 | double | Inherit: Inherit ... | | | -5 | 10 |
| ☐ | Gain2 | double | Inherit: Inherit ... | | | -Inf | Inf |
| ☐ | Gain3 | double | Inherit: Inherit ... | | | -Inf | Inf |
| ☐ | In1 | double | double | -10 | 20 | -10 | 20 |
| ☐ | Out1 | | Inherit: auto | | | -Inf | Inf |
| ☐ | Unit Delay | double | | | | -Inf | Inf |

- 6 To fix the issue, specify design minimum and maximum values inside the feedback loop. For this example, specify the range for the Gain2 block:
 - a In the model, double-click the Gain2 block.
 - b In the block parameters dialog box, select the **Signal Attributes** tab.
 - c In this tab, set **Output minimum** to -20 and **Output maximum** to 40 and click **OK**.
- 7 Clear previously collected ranges and rerun the range analysis.
 - a In the Fixed-Point Tool, under **New** workflow, select **Range Collection**.
Changing workflows clears range data collected during the active workflow.
 - b Switch back to the **Iterative Fixed-Point Conversion** workflow.
 - c Select **Derived ranges** as the range collection mode.
 - d Click the **Collect Ranges** button again to rerun the range analysis.

The range analysis uses the minimum and maximum values specified for Gain2 and In1 to derive ranges for all objects in the model.

Provide Additional Design Range Information

- 1 To fix the issue, specify design minimum and maximum values inside the feedback loop. For this example, specify the range for the Gain2 block:
 - a In the model, double-click the Gain2 block.
 - b In the block parameters dialog box, select the **Signal Attributes** tab.
 - c In this tab, set **Output minimum** to -20 and **Output maximum** to 40 and click **OK**.
- 2 Clear previously collected ranges and rerun the range analysis.
 - a In the Fixed-Point Tool, under **New** workflow, select **Range Collection**.
Changing workflows clears range data collected during the active workflow.
 - b Switch back to the **Iterative Fixed-Point Conversion** workflow.
 - c Select **Derived ranges** as the range collection mode.
 - d Click the **Collect Ranges** button again to rerun the range analysis.

The range analysis uses the minimum and maximum values specified for Gain2 and In1 to derive ranges for all objects in the model.

See Also

Related Examples

- “Insufficient Design Range Information” on page 43-14

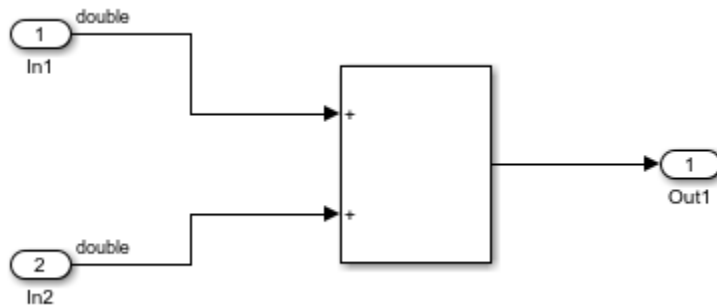
Fix Design Range Conflicts

This example shows how to fix design range conflicts. If you specify conflicting design minimum and maximum values in your model, the range analysis software reports an error. To fix this error, examine the design ranges specified in the model to identify inconsistent design specifications. Modify them to make them consistent. In this example, the output design range specified on the Output block conflicts with the input design ranges specified on the Inport blocks.

Open Model

Open the `ex_range_conflict` model.

```
open_system("ex_range_conflict.slx")
```



The model displays the specified design minimum and maximum values for the blocks in the model.

- The Inport blocks In1 and In2 have a design range of $[-1..1]$.
- The Output block Out1 has a design range of $[10..20]$.

To display design ranges in your model, in the **Debug** tab, select **Information Overlays > Signal Data Ranges**.

Collect Ranges in the Fixed-Point Tool

- 1 From the Simulink **Apps** tab, select **Fixed-Point Tool**.
- 2 In the Fixed-Point Tool, under **New** workflow, select **Iterative Fixed-Point Conversion**.
- 3 In the Fixed-Point Tool, under **System Under Design (SUD)**, select `ex_range_conflict` as the system you want to convert.
- 4 Under **Range Collection Mode**, select **Derived ranges**.
- 5 Click the **Collect Ranges** button.

The Fixed-Point Tool reports an error because the derived range for the Sum block, $[-2..2]$ is outside the specified design range for the Output block, $[10..20]$.

Fix Design Range Conflicts

- 1 To fix the conflict, change the design range on the Outport block to [-10..20] so that this range includes the derived range for the Sum block.
 - a In the model, double-click the Outport block.
 - b In the block parameters dialog box, click the **Signal Attributes** tab.
 - c In this tab, set **Minimum** to -10 and click **OK**.
- 2 Clear previously collected ranges and rerun the range analysis.
 - a In the Fixed-Point Tool, under **New** workflow, select **Range Collection**.
Changing workflows clears range data collected during the active workflow.
 - b Switch back to the **Iterative Fixed-Point Conversion** workflow.
 - c Select **Derived ranges** as the range collection mode.
 - d Click the **Collect Ranges** button again to rerun the range analysis.

The range analysis derives a minimum value of -2 and a maximum value of 2 for the Outport block.

See Also

More About

- “How Range Analysis Works” on page 43-2

Intermediate Range Results

In this section...

“Open Model” on page 43-24

“Collect Ranges in the Fixed-Point Tool” on page 43-24

“Propose Data Types” on page 43-25

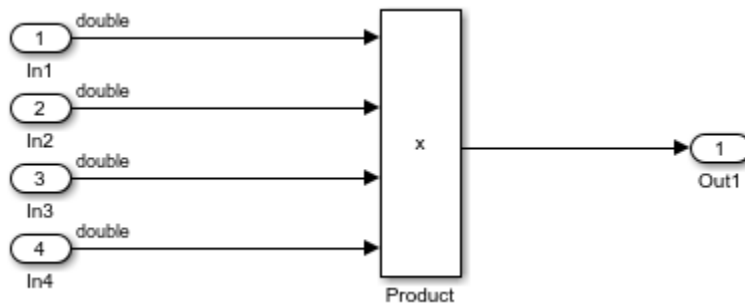
“Inspect Result Details” on page 43-25

This example shows how to interpret the **Intermediate Maximum** and **Intermediate Minimum** results in the **Result Details** tab.

Open Model

Open the `ex_intermediateRange` model.

```
open_system("ex_intermediateRange.slx")
```



Update the diagram (Ctrl+D). Notice the design range information for each of the input ports.

To display design ranges in your model, in the **Debug** tab, select **Information Overlays > Signal Data Ranges**.

Collect Ranges in the Fixed-Point Tool

- 1 Open the Fixed-Point Tool. From the Simulink **Apps** tab, select **Fixed-Point Tool**.
- 2 In the Fixed-Point Tool, under **New** workflow, select **Iterative Fixed-Point Conversion**.
- 3 In the Fixed-Point Tool, under **System Under Design (SUD)**, select `ex_intermediateRange` as the system you want to convert.
- 4 Under **Range Collection Mode**, select **Derived ranges**.
- 5 Click the **Collect Ranges** button.

The Fixed-Point Tool displays the derived minimum and maximum values for each object in the `ex_intermediateRange` model.

Propose Data Types

- 1 In the **Convert** section of the toolstrip, open the **Settings** menu.

In the **Default word length** field, enter 32

- 2 Click the **Propose Data Types** button .

The tool displays the proposed data types appear in the spreadsheet.

Inspect Result Details

- 1 Look at the proposed data type of the Product block. The Fixed-Point Tool proposed a data type with 32-bit word length and 12-bit fraction length. The derived maximum value of the Product block is 1, but the maximum representable value for the proposed data type is approximately 1,048,575.

To learn more about the data type proposal, select the product block in the spreadsheet. The **Result Details** pane populates with information about the result.

- 2 In the **Result Details** pane, in the **Ranges used for proposal** table, notice the row labeled **Intermediate**. After the first two inputs to the Product block are multiplied, the block has a maximum value of 1000000 before being multiplied by the next two inputs for a final maximum value of 1. The data type proposal for the Product block in this model is based on the intermediate minimum and maximum values. It is not based on the derived minimum and maximum values to prevent overflows at the intermediate stages of the block.

RESULT DETAILS

ex_intermediateRange/Product

Proposed Data Type Summary

| Property | ProposedDT | SpecifiedDT |
|-----------|--------------------|--------------------------|
| DataType | fixdt(0,32,12) | double |
| Minimum | 0 | -1.7976931348623157e+... |
| Maximum | 1048575.9997558594 | 1.7976931348623157e+308 |
| Precision | 0.000244140625 | 4.94065645841247e-324 |

Ranges used for proposal

| Property | Minimum | Maximum |
|----------------|---------|--------------------|
| Shared derived | 0 | 1.0000000000000002 |
| Derived | 0 | 1.0000000000000002 |
| Intermediate | 0 | 1000000 |

Proposal Details

- There is a requirement for the data type of this result to match the data type of other results.
 - [Highlight Elements Sharing Same Data Type](#)

See Also

More About

- “How Range Analysis Works” on page 43-2

Unsupported Simulink Software Features

Range analysis does not support the following Simulink software features. Avoid using these unsupported features.

| Not Supported | Description |
|---|--|
| Variable-step solvers | The software supports only fixed-step solvers. For more information, see “Fixed Step Solvers in Simulink”. |
| Callback functions | The software does not execute model callback functions during the analysis. The results that the analysis generates may behave inconsistently with the expected behavior. <ul style="list-style-type: none"> • If a model or any referenced model calls a callback function that changes any block parameters, model parameters, or workspace variables, the analysis does not reflect those changes. • Changing the storage class of base workspace variables on model callback functions or mask initializations is not supported. • Callback functions called prior to analysis, such as the <code>PreLoadFcn</code> or <code>PostLoadFcn</code> model callbacks, are fully supported. |
| Model callback functions | The software only supports model callback functions if the <code>InitFcn</code> callback of the model is empty. |
| Algebraic loops | The software does not support models that contain algebraic loops. For more information, see “Algebraic Loop Concepts”. |
| Masked subsystem initialization functions | The software does not support models whose masked subsystem initialization modifies any attribute of any workspace parameter. |
| Variable-size signals | The software does not support variable-size signals. A variable-size signal is a signal whose size (number of elements in a dimension), in addition to its values, can change during model execution. |
| Arrays of buses | The software does not support arrays of buses. For more information, see “Group Nonvirtual Buses in Arrays of Buses”. |
| Multiword fixed-point data types | The software does not support multiword fixed-point data types. |
| Nonfinite data | The software does not support nonfinite data (for example, <code>NaN</code> and <code>Inf</code>) and related operations. |
| Signals with nonzero sample time offset | The software does not support models with signals that have nonzero sample time offsets. |
| Models with no output ports | The software only supports models that have one or more output ports. |

Note The software does not report initial or intermediate values for Stateflow variables. Range analysis will only report the ranges at the output of the block.

Simulink Blocks Supported for Range Analysis

Overview of Simulink Block Support

The following tables summarize range analysis support for Simulink blocks. Each table lists all the blocks in each Simulink library and describes support information for that particular block. If the software does not support a given block, where possible, automatic stubbing considers the interface of the unsupported blocks, but not their behavior, during the analysis. However, if any of the unsupported blocks affect the simulation outcome, the analysis may achieve only partial results. If the analysis cannot use automatic stubbing for a block, the block is marked as “not stubbable”. For more information, see “Automatic Stubbing” on page 43-4.

Not all blocks that are supported for range analysis are supported for fixed-point conversion. To check if a block supports fixed-point data types, see “Blocks That Do Not Support Fixed-Point Data Types” on page 49-19.

Additional Math and Discrete Library

The software supports all blocks in the Additional Math and Discrete library.

Commonly Used Blocks Library

The Commonly Used Blocks library includes blocks from other libraries. Those blocks are listed under their respective libraries.

Continuous Library

| Block | Support Notes |
|---------------------------------|---------------|
| Derivative | Not supported |
| Integrator | Not supported |
| Integrator Limited | Not supported |
| PID Controller | Not supported |
| PID Controller (2DOF) | Not supported |
| Second-Order Integrator | Not supported |
| Second-Order Integrator Limited | Not supported |
| State-Space | Not supported |
| Transfer Fcn | Not supported |
| Transport Delay | Not supported |
| Variable Time Delay | Not supported |
| Variable Transport Delay | Not supported |
| Zero-Pole | Not supported |

Discontinuities Library

The software supports all blocks in the Discontinuities library.

Discrete Library

| Block | Support Notes |
|---------------------------------|--|
| Delay | Supported |
| Difference | Supported |
| Discrete Derivative | Supported |
| Discrete Filter | The software analyzes through the filter. It does not derive any range information for the filter. |
| Discrete FIR Filter | Supported |
| Discrete PID Controller | Supported |
| Discrete PID Controller (2 DOF) | Supported |
| Discrete State-Space | Not supported |
| Discrete-Time Integrator | Supported |
| Discrete Transfer Fcn | Supported |
| Discrete Zero-Pole | Not supported |
| Memory | Supported |
| Tapped Delay | Supported |
| Transfer Fcn First Order | Supported |
| Transfer Fcn Lead or Lag | Supported |
| Transfer Fcn Real Zero | Supported |
| Unit Delay | Supported |
| Zero-Order Hold | Supported |

Logic and Bit Operations Library

The software supports all blocks in the Logic and Bit Operations library.

Lookup Tables Library

| Block | Support Notes |
|-------------------------------|---|
| Cosine | Supported |
| Direct Lookup Table (n-D) | Supported |
| Interpolation Using Prelookup | Partially supported when: <ul style="list-style-type: none"> The Interpolation method parameter is Linear and the Number of table dimensions parameter is greater than 4. or <ul style="list-style-type: none"> The Interpolation method parameter is Linear and the Number of sub-table selection dimensions parameter is not 0. |
| 1-D Lookup Table | Partially supported when the Interpolation method or the Extrapolation method parameter is Cubic Spline . |

| Block | Support Notes |
|----------------------|--|
| 2-D Lookup Table | Not supported when the Interpolation method or the Extrapolation method parameter is Akima Spline. |
| n-D Lookup Table | Partially supported when: <ul style="list-style-type: none"> The Interpolation method or the Extrapolation method parameter is Cubic Spline. or <ul style="list-style-type: none"> The Interpolation method parameter is Linear and the Number of table dimensions parameter is greater than 5. Not supported when the Interpolation method or the Extrapolation method parameter is Akima Spline. |
| Lookup Table Dynamic | Supported |
| Prelookup | Supported |
| Sine | Supported |

Math Operations Library

| Block | Support Notes |
|----------------------------|---|
| Abs | Supported |
| Add | Supported |
| Algebraic Constraint | Supported |
| Assignment | Supported |
| Bias | Supported |
| Complex to Magnitude-Angle | Supported |
| Complex to Real-Imag | Supported |
| Divide | Supported |
| Dot Product | Supported |
| Find Nonzero Elements | Not supported |
| Gain | Supported |
| Magnitude-Angle to Complex | Supported |
| Math Function | Supported. Support for pow function is limited to integer exponents only. |
| Matrix Concatenate | Supported |
| MinMax | Supported |
| MinMax Running Resettable | Supported |
| Permute Dimensions | Supported |
| Polynomial | Supported |
| Product | Supported |
| Product of Elements | Supported |

| Block | Support Notes |
|---------------------------|--|
| Real-Imag to Complex | Supported |
| Reciprocal Sqrt | Partially supported |
| Reshape | Supported |
| Rounding Function | Supported |
| Sign | Supported |
| Signed Sqrt | Partially supported |
| Sine Wave Function | Partially supported |
| Slider Gain | Supported |
| Sqrt | Partially supported |
| Squeeze | Supported |
| Subtract | Supported |
| Sum | Supported |
| Sum of Elements | Supported |
| Trigonometric Function | Supported if Function is sin, cos, or sincos, and Approximation method is CORDIC. Partially supported otherwise. |
| Unary Minus | Supported |
| Vector Concatenate | Supported |
| Weighted Sample Time Math | Supported |

Model Verification Library

The software supports all blocks in the Model Verification library.

Model-Wide Utilities Library

| Block | Support Notes |
|-----------------------------|---------------|
| Block Support Table | Supported |
| DocBlock | Supported |
| Model Info | Supported |
| Timed-Based Linearization | Not supported |
| Trigger-Based Linearization | Not supported |

Ports & Subsystems Library

| Block | Support Notes |
|--|---------------|
| Subsystem, Atomic Subsystem, CodeReuse Subsystem | Supported |
| Configurable Subsystem | Supported |
| Enable | Supported |

| Block | Support Notes |
|---------------------------------|--|
| Enabled Subsystem | Range analysis does not consider the design minimum and maximum values specified for blocks connected to the output of the subsystem. |
| Enabled and Triggered Subsystem | <p>Not supported when the trigger control signal specifies a fixed-point data type.</p> <p>Range analysis does not consider the design minimum and maximum values specified for blocks connected to the output of the subsystem.</p> |
| For Each | <p>Supported with the following limitations:</p> <ul style="list-style-type: none"> • When For Each Subsystem contains another For Each Subsystem, not supported. • When For Each Subsystem contains one or more Simulink Design Verifier™ Test Condition, Test Objective, Proof Assumption, or Proof Objective blocks, not supported. |
| For Each Subsystem | <p>Supported with the following limitations:</p> <ul style="list-style-type: none"> • When For Each Subsystem contains another For Each Subsystem, not supported. • When For Each Subsystem contains one or more Simulink Design Verifier Test Condition, Test Objective, Proof Assumption, or Proof Objective blocks, not supported. |
| For Iterator Subsystem | Supported |
| Function-Call Feedback Latch | Supported |
| Function-Call Generator | Supported |
| Function-Call Split | Supported |
| Function-Call Subsystem | Range analysis does not consider the design minimum and maximum values specified for blocks connected to the output of the subsystem. |
| If | Supported |
| If Action Subsystem | Supported |
| Inport | — |
| Model | Supported except for the limitations described in “Limitations of Support for Model Blocks” on page 43-35. |
| Outport | Supported |
| Switch Case | Supported |
| Switch Case Action Subsystem | Supported |
| Trigger | Supported |

| Block | Support Notes |
|----------------------------------|---|
| Triggered Subsystem | Not supported when the trigger control signal specifies a fixed-point data type. Range analysis does not consider the design minimum and maximum values specified for blocks connected to the output of the subsystem. |
| Variant Subsystem, Variant Model | Supported |
| While Iterator Subsystem | Supported |

Signal Attributes Library

The software supports all blocks in the Signal Attributes library.

Signal Routing Library

| Block | Support Notes |
|------------------------|---|
| Bus Assignment | Supported |
| Bus Creator | Supported |
| Bus Selector | Supported |
| Data Store Memory | <ul style="list-style-type: none"> When the Data Store Memory variable is tunable, range analysis considers the design ranges specified on the block, and ignores local model writes. When the Data Store Memory variable is not tunable, or Auto, the analysis considers only local model writes. The derived range is the range of the last write to the variable. When the Data Store Memory variable is defined outside of the analyzed system, range analysis uses design ranges. |
| Data Store Read | Supported |
| Data Store Write | Supported |
| Demux | Supported |
| Environment Controller | Supported |
| From | Supported |
| Goto | Supported |
| Goto Tag Visibility | Supported |
| Index Vector | Supported |
| Manual Switch | The Manual Switch block is compatible with the software, but the analysis ignores this block in a model. |
| Merge | Supported |
| Multiport Switch | Supported |
| Mux | Supported |
| Selector | Supported |
| Switch | Supported |

| Block | Support Notes |
|--------------------|---------------|
| Vector Concatenate | Supported |

Sinks Library

| Block | Support Notes |
|-----------------|---------------|
| Display | Supported |
| Floating Scope | Supported |
| Outport (Out1) | Supported |
| Out Bus Element | Supported |
| Scope | Supported |
| Stop Simulation | Not supported |
| Terminator | Supported |
| To File | Supported |
| To Workspace | Supported |

Sources Library

| Block | Support Notes |
|---------------------------------|---|
| Band-Limited White Noise | Not supported |
| Chirp Signal | Partially supported |
| Clock | Supported |
| Constant | Supported unless Constant value is inf or nan (in which case, it is not supported). |
| Counter Free-Running | Supported |
| Counter Limited | Supported |
| Digital Clock | Supported |
| Enumerated Constant | Supported |
| From File | Partially supported. When MAT-file data is stored in MATLAB timeseries format, not supported. |
| From Workspace | Partially supported |
| Ground | Supported |
| Inport (In1) | Supported |
| In Bus Element | Supported if Simulink.Bus type is defined for the In Bus Element. |
| Pulse Generator | Supported |
| Ramp | Supported |
| Random Number | Not supported |
| Repeating Sequence | Partially supported |
| Repeating Sequence Interpolated | Partially supported |

| Block | Support Notes |
|--------------------------|---|
| Repeating Sequence Stair | Supported |
| Signal Editor | Not supported |
| Signal Generator | Partially supported if wave form is <code>sine</code> . Supported if wave form is <code>square</code> . Not supported if wave form is <code>random</code> . |
| Sine Wave | Partially supported |
| Step | Supported |
| Uniform Random Number | Not supported |

User-Defined Functions Library

| Block | Support Notes |
|-----------------------------|---|
| Interpreted MATLAB Function | Not supported |
| MATLAB Function | <p>The software uses the specified design minimum and maximum values and returned derived minimum and maximum values for instances of variables that correspond to input and output ports. It does not consider intermediate instances of these variables. For example, consider a MATLAB Function block that contains the following code:</p> <pre>function y = fcn(u,v) %#codegen y = 2*u; y = y + v;</pre> <p>Range analysis considers the design ranges specified for <code>u</code> and <code>v</code> for the instance of <code>y</code> in <code>y = y + v;</code> because this is the instance of <code>y</code> associated with the output of the block.</p> <p>The analysis does not consider design ranges for the instance of <code>y</code> in <code>y = 2*u;</code> because it is an intermediate instance.</p> |
| Level-2 MATLAB S-Function | Not supported |
| S-Function | Not supported |
| S-Function Builder | Not supported |
| Simulink Function | Simulink Functions with output arguments that are of complex type are not supported. |

Limitations of Support for Model Blocks

Range analysis supports the Model block with the following limitations. The software cannot analyze a model containing one or more Model blocks if:

- The referenced model is protected. Protected referenced models are encoded to obscure their contents. This allows third parties to use the referenced model without being able to view the intellectual property that makes up the model.

For more information, see “Reference Protected Models from Third Parties”.

- The parent model or any of the referenced models returns an error when you set the **Configuration Parameters > Diagnostics > Connectivity > Element name mismatch** parameter to error.

You can use the **Element name mismatch** diagnostic along with bus objects so that your model meets the bus element naming requirements imposed by some blocks.

- The Model block uses asynchronous function-call inputs.
- Any of the Model blocks in the model reference hierarchy creates an artificial algebraic loop. If this occurs, take the following steps:

- 1 On the **Diagnostics** pane of the Configuration Parameters dialog box, set the **Minimize algebraic loop** parameter to error so that Simulink reports an algebraic loop error.
- 2 On the **Model Referencing** Pane of the Configuration Parameters dialog box, select the **Minimize algebraic loop occurrences** parameter.

Simulink tries to eliminate the artificial algebraic loop during simulation.

- 3 Simulate the model.
- 4 Simulink will remove the algebraic loop if possible. If Simulink cannot eliminate the artificial algebraic loop, highlight the location of the algebraic loop by opening the **Modeling** tab and, in the **Compile** section, clicking **Update Model**.
- 5 Eliminate the artificial algebraic loop so that the software can analyze the model. Break the loop with Unit Delay blocks so that the execution order is predictable.

Note For more information, see “Algebraic Loop Concepts”.

- The parent model and the referenced model have mismatched data type override settings. The data type override setting of the parent model and its referenced models must be the same, unless the data type override setting of the parent model is `Use local settings`. You can configure data type override settings to simulate a model that specifies fixed-point data types. Using this setting, the software temporarily overrides data types with floating-point data types during simulation.

```
set_param('MyModel', 'DataTypeOverride', 'Double')
```

For more information, see `set_param`.

To observe the true behavior of your model, set the data type override parameter to `UseLocalSettings` or `Off`.

```
set_param('MyModel', 'DataTypeOverride', 'Off')
```

- The referenced model is a Model block with virtual buses at input ports, and the signals in the bus do not all have the same sample time at compilation. To make the model compatible with Simulink Design Verifier analysis, convert the virtual bus to a nonvirtual bus, or specify an explicit sample time for the port.
- When you run the analysis on Model block, then the code generated as a top model is not supported.

Range Collection Workflows

- “Use the Fixed-Point Tool to Explore Numerical Behavior” on page 44-2
- “Use Custom Data Type Override Settings for Range Collection” on page 44-9

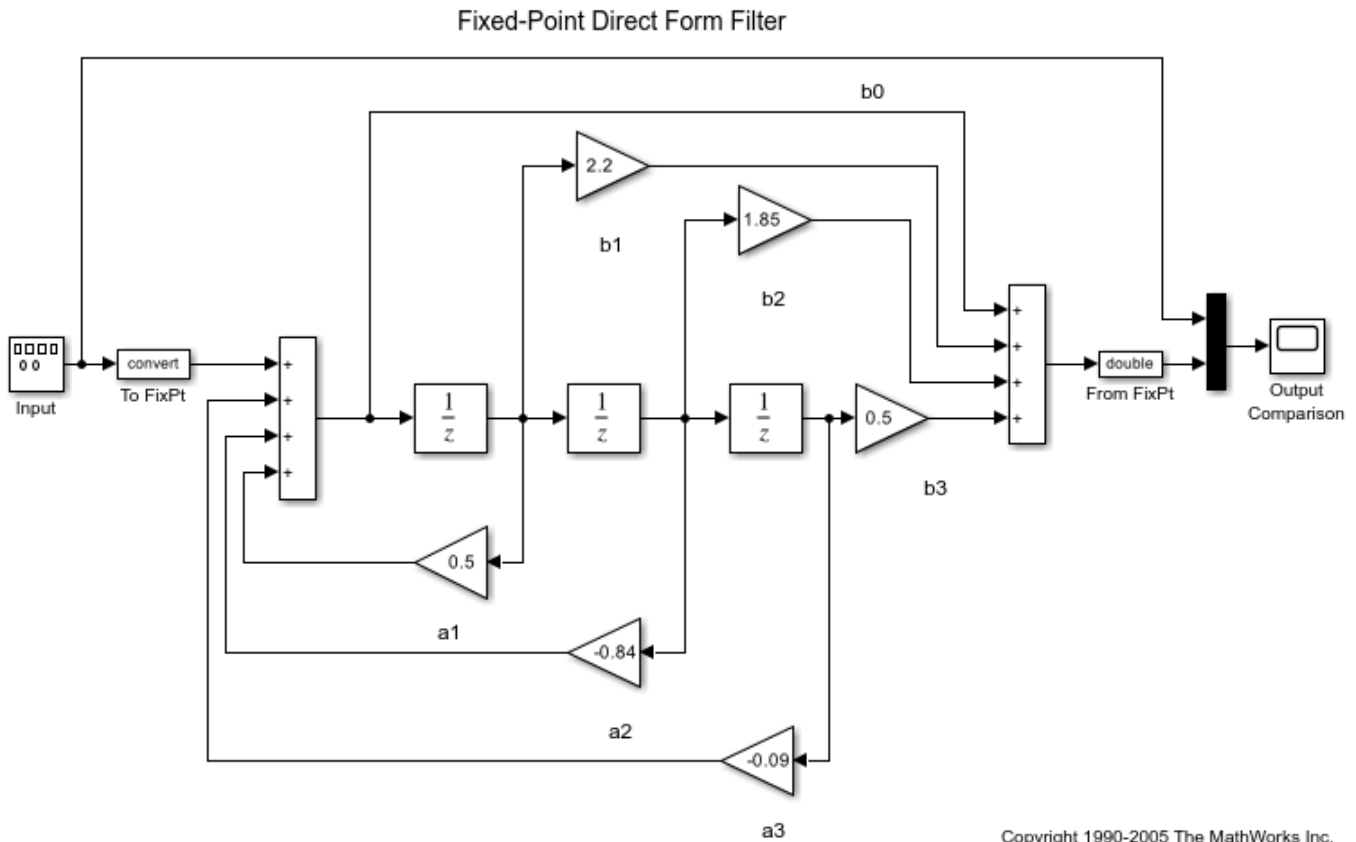
Use the Fixed-Point Tool to Explore Numerical Behavior

| In this section... |
|--|
| “Open the Fixed-Point Direct Form Filter Model” on page 44-2 |
| “Set Up the Model” on page 44-3 |
| “Open the Fixed-Point Tool and Collect Ranges” on page 44-4 |
| “Explore Fixed-Point Behavior of the Model” on page 44-6 |

This example shows how to use the Fixed-Point Tool to compare floating-point and fixed-point data types in your model. You can use the range collection functionality to explore and troubleshoot the numerical behavior of your model for different inputs.

Open the Fixed-Point Direct Form Filter Model

Open the `fxpdemo_direct_form2` model. This tutorial uses a fixed-point direct form filter implemented using fundamental building blocks such as Gain, Delay, and Sum. The model contains a Signal Generator block that supplies a square wave input to the filter.



Set Up the Model

In this tutorial, you explore the behavior of the filter for a range of signal inputs. To specify multiple simulation scenarios for range collection, define a `Simulink.SimulationInput` object in the base or model workspace. Define a `Simulink.SimulationInput` object, `simIn`, that specifies the amplitude of the square wave input for a range of values.

```
simIn(1:6) = Simulink.SimulationInput('fxpdemo_direct_form2');

simIn(1) = simIn(1).setBlockParameter('fxpdemo_direct_form2/Input',...
    'Amplitude','0.001');
simIn(2) = simIn(2).setBlockParameter('fxpdemo_direct_form2/Input',...
    'Amplitude','0.01');
simIn(3) = simIn(3).setBlockParameter('fxpdemo_direct_form2/Input',...
    'Amplitude','0.1');
simIn(4) = simIn(4).setBlockParameter('fxpdemo_direct_form2/Input',...
    'Amplitude','1');
simIn(5) = simIn(5).setBlockParameter('fxpdemo_direct_form2/Input',...
    'Amplitude','10');
simIn(6) = simIn(6).setBlockParameter('fxpdemo_direct_form2/Input',...
    'Amplitude','100');
```

To specify signal tolerances, enable signal logging at the output of the Sum1 block.

```
Simulink.sdi.markSignalForStreaming('fxpdemo_direct_form2/Sum1',1,'on');
```

Open the Fixed-Point Tool and Collect Ranges

- 1 In the **Apps** tab of the `fxpdemo_direct_form2` model, select **Fixed-Point Tool**.
- 2 In the Fixed-Point Tool, click **New > Range Collection**.
- 3 Under **System Under Design (SUD)**, select `fxpdemo_direct_form2`.
- 4 Under **Range Collection Mode**, select **Simulation Ranges** as the range collection method.
- 5 Under **Simulation Inputs**, select the `Simulink.SimulationInput` object, `simIn`, that you defined in the base workspace.
- 6 To specify tolerances for the system, under **Signal Tolerances**, specify tolerances for any signal in the model with signal logging enabled.

Set the relative tolerance (**Rel Tol**) of the signal that you logged to 15%.

▼ Signal Tolerances

Specify tolerances for signals in your model that have signal logging enabled. After converting your system to fixed point, the Workflow Browser displays whether the embedded run meets the specified signal tolerances.

Filter signal list: Refresh Signals

| Signal Name | Abs Tol | Rel Tol | Time Tol (seconds) |
|-------------|---------|---------|--------------------|
| Sum1:1 | | 0.15 | |

- 7 Under **Collect Ranges**, select **Double precision**.

Collect Ranges

Collect simulation ranges using data type override

Collect ranges using:

Use current settings
Use current data type override set on the model

Double precision
Override data types with doubles

Single precision
Override data types with singles

Scaled double precision
Override data types with scaled doubles

When you collect ranges via simulation, the Fixed-Point Tool will override the data types in your model with doubles and simulate the model with instrumentation to collect minimum and maximum values for each object in your model. You can also choose to override data types with singles or scaled doubles, or use the current data type override set on the model.

- 8 Click the **Collect Ranges** button.

Simulink simulates the `fxpdemo_direct_form2` model six times, once for each amplitude of the input square wave specified in the `Simulink.SimulationInput` object. The Fixed-Point Tool automatically enables fixed-point instrumentation and overrides the data types in your model with doubles to collect a floating-point baseline.

You can view the ranges of each simulation individually by selecting the simulation scenario in the **Workflow Browser**.

Selecting the `BaselineRun` node in the **Workflow Browser** shows the merged ranges from the six simulation scenarios.

The screenshot displays the Fixed-Point Tool interface. The **Results** table is shown below:

| Name | CompiledDT | SpecifiedDT | SimMin | SimMax |
|--------------------|------------|-----------------------------|---------------------|--------------------|
| From FixPt | double | double | -606 | 823.4057374709746 |
| Gain | double | fixdt(1,16,10) | -75 | 114.87944438716877 |
| Gain1 | double | fixdt(1,16,10) | -20.67829998969038 | 13.5 |
| Gain2 | double | fixdt(1,16,10) | -192.99746657044352 | 126 |
| Gain3 | double | fixdt(1,16,10) | -75 | 114.87944438716877 |
| Gain4 | double | fixdt(1,16,10) | -330 | 505.469553035427 |
| Gain5 | double | fixdt(1,16,10) | -277.5 | 425.0539442325245 |
| Sum : Accumulator | double | Inherit: Inherit via int... | -192.99746657044352 | 229.75888877433755 |
| Sum : Output | double | fixdt(1,16,10) | -150 | 229.75888877433755 |
| Sum1 : Accumulator | double | Inherit: Inherit via int... | -606 | 858.1419352251326 |
| Sum1 : Output | double | fixdt(1,16,10) | -606 | 823.4057374709746 |
| To FixPt | double | fixdt(1,16,10) | -100 | 100 |

The **Histograms of all results in the model** visualization shows the distribution of data points across various bins. The legend indicates the following counts:

| Values | Potential Overflows | In-Range | Potential Underflows |
|----------|---------------------|----------|----------------------|
| Positive | 102 | 491 | 1 |
| Negative | 101 | 509 | 2 |
| Zero | 0 | 0 | 0 |

- 9 Click **Settings**, then select Specified data types.
 10 Click **Simulate with Embedded Types**.

The Fixed-Point Tool simulates the model once for each simulation scenario, using the fixed-point data types specified in the model. Selecting the `EmbeddedRun` node in the **Workflow Browser** shows the merged results from the six simulation scenarios.

The screenshot displays the Range Collection tool interface. The top bar includes 'RANGE COLLECTION' and 'EXPLORE' tabs. Below the top bar are icons for 'New', 'Collect Ranges', 'Settings', and 'Simulate with Embedded Types'. The main area is divided into 'WORKFLOW BROWSER' on the left, 'Results' table in the center, and 'RESULT DETAILS' on the right. The 'Results' table lists simulation scenarios and their properties, with some rows highlighted in red to indicate overflows. The 'RESULT DETAILS' panel shows 'Needs Attention' with a red 'X' icon and a message: 'There are overflows associated with this result.' Below this, a table shows 'Property' and 'Specified Data Type' for 'fxpdemo_direct_form2/Sum1 : Output'. The 'Range Information' table shows 'Property', 'Minimum', and 'Maximum' for 'Simulation'. The bottom section shows 'Visualization of Simulation Data' with a histogram titled 'Histograms of all results in the model'. The histogram shows the distribution of data values, with a legend indicating 'Overflows', 'Representable', 'In-Range', and 'Underflows'. A secondary histogram shows 'Visualization of Simulation Data using fixdt(1,16,10)' with a legend for 'Values' (Positive, Negative, Zero) and 'In-Range' (550, 638, 18).

| Name | CompiledDT | SpecifiedDT | SimMin | SimMax |
|--------------------|----------------|-----------------------------|----------------|---------------|
| From FixPt | double | double | -31.7509765625 | 31.8525390625 |
| Gain | fixdt(1,16,10) | fixdt(1,16,10) | -15.384765625 | 14 |
| Gain1 | fixdt(1,16,10) | fixdt(1,16,10) | -2.5205078125 | 2.7685546875 |
| Gain2 | fixdt(1,16,10) | fixdt(1,16,10) | -23.5205078125 | 25.845703125 |
| Gain3 | fixdt(1,16,10) | fixdt(1,16,10) | -15.384765625 | 14 |
| Gain4 | fixdt(1,16,10) | fixdt(1,16,10) | -31.435030375 | 31.4009375 |
| Gain5 | fixdt(1,16,10) | fixdt(1,16,10) | -31.6904375 | 31.9912100375 |
| Sum : Accumulator | fixdt(1,32,10) | Inherit. Inherit via int... | -51.341796875 | 60.0615234375 |
| Sum : Output | fixdt(1,16,10) | fixdt(1,16,10) | -30.78953125 | 28 |
| Sum1 : Accumulator | fixdt(1,32,10) | Inherit. Inherit via int... | -74.0087890625 | 74.2822265625 |
| Sum1 : Output | fixdt(1,16,10) | fixdt(1,16,10) | -31.7509765625 | 31.8525390625 |
| To FixPt | fixdt(1,16,10) | fixdt(1,16,10) | 28 | 28 |

| Property | Specified Data Type |
|-----------|---------------------|
| Data Type | fixdt(1,16,10) |
| Minimum | -32 |
| Maximum | 31.9990234375 |
| Precision | 0.0009765625 |

| Property | Minimum | Maximum |
|------------|----------------|---------------|
| Simulation | -31.7509765625 | 31.8525390625 |

Visualization of Simulation Data using fixdt(1,16,10)
 ▲ Overflows (# of wraps) : 286

| Values | In-Range |
|----------|----------|
| Positive | 550 |
| Negative | 638 |
| Zero | 18 |

The **Workflow Browser** indicates that of the six simulation scenarios, only `EmbeddedRun_Scenario_4` met the tolerances specified. Results with overflows are highlighted in red.

Explore Fixed-Point Behavior of the Model

- 1 Select the **Explore** tab of the Fixed-Point Tool to investigate further. Under **Numerical Issues**, select **Overflow**, then click **Execution Order**.

The screenshot shows the Fixed-Point Tool interface. The top navigation bar includes 'RANGE COLLECTION' and 'EXPLORE'. The 'EXPLORE' section has filters for 'Data Types' (Fixed-Point, Double), 'Numerical Issues' (Overflow, Underflow), 'Number Space' (Integer, Fraction), and 'Signedness' (Signed, Unsigned). A 'Clear Filter' button is also present. The 'WORKFLOW BROWSER' on the left shows a hierarchy of runs, with 'EmbeddedRun' selected. The 'MODEL HIERARCHY' on the bottom left shows 'Simulink Root' and 'Data Objects'. The main 'Results' table is as follows:

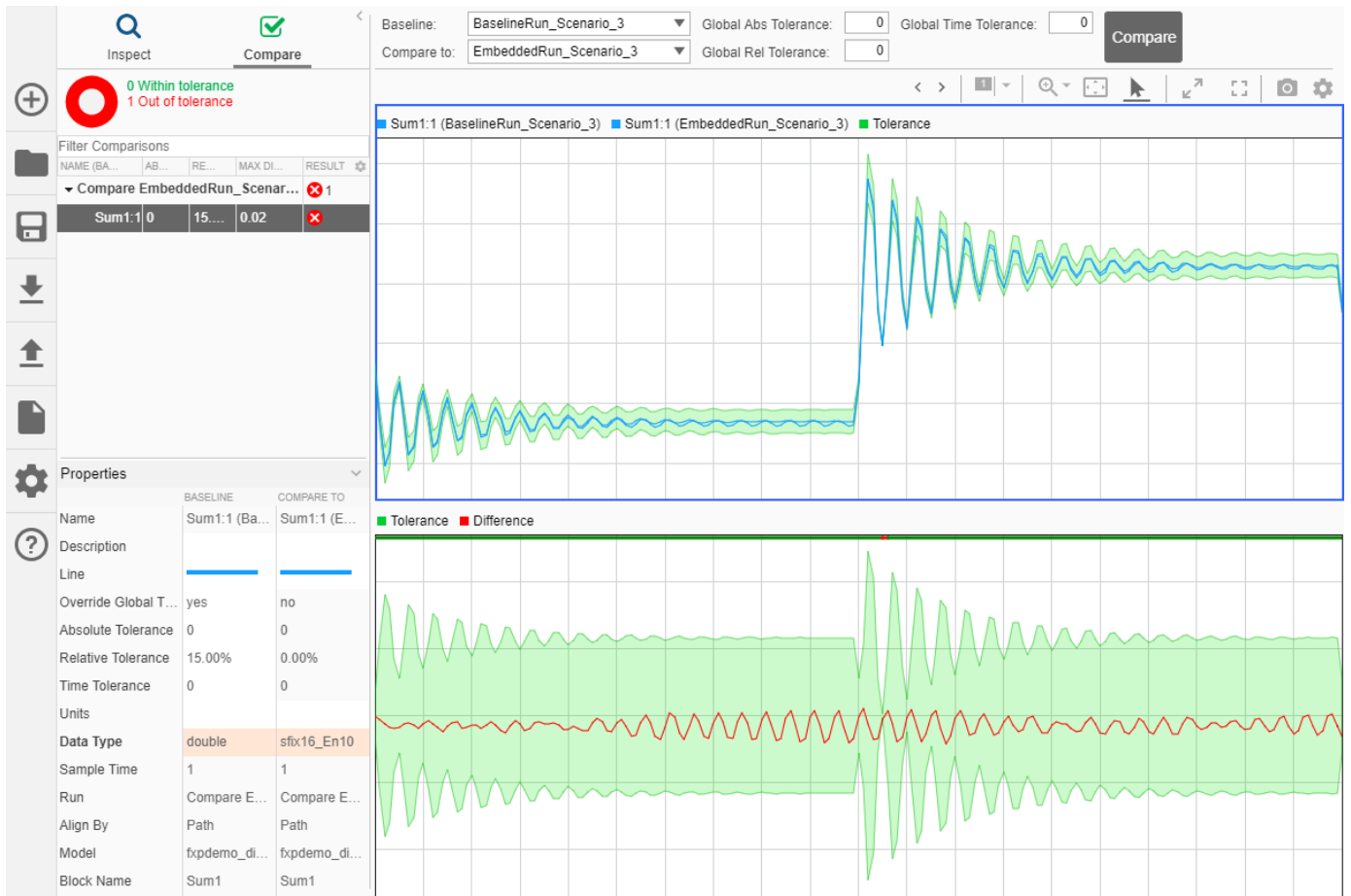
| Name | CompiledDT | SpecifiedDT | SimMin | SimMax |
|---------------|----------------|----------------|----------------|---------------|
| Gain4 | fixed(1,16,10) | fixed(1,16,10) | -31.435846875 | 31.4609375 |
| Gain6 | fixed(1,16,10) | fixed(1,16,10) | -31.8984375 | 31.9912109375 |
| Sum1 - Output | fixed(1,16,10) | fixed(1,16,10) | -30.76983125 | 28 |
| Sum1 - Output | fixed(1,16,10) | fixed(1,16,10) | -31.7509765625 | 31.8528390625 |
| To FixPt | fixed(1,16,10) | fixed(1,16,10) | 28 | 28 |

Below the table is a 'Visualization of Simulation Data' section titled 'Histograms of all results in the model'. It shows five histograms corresponding to the blocks in the table above. A legend indicates that red bars represent 'Overflows', white bars represent 'Representable', blue bars represent 'In-Range', and yellow bars represent 'Underflows'. The 'Gain4' histogram shows a significant portion of red bars, indicating overflows.

The Fixed-Point Tool displays only the EmbeddedRun results with overflows and sorts the list based on block execution order. In this example, the first overflow occurs in the Gain4 block.

You can double-click on any row in the **Results** spreadsheet to highlight the block in the model.

- 2 You can compare the fixed-point and floating-point behavior of the model for a specific simulation scenario using the Simulation Data Inspector. For example, the Fixed-Point Tool indicates that EmbeddedRun_Scenario_3 did not meet the specified tolerance. To compare this embedded run to the floating-point behavior for this simulation scenario, right-click on EmbeddedRun_Scenario_3 and select **Open SDI** to compare with BaselineRun_Scenario_3.



The Simulation Data Inspector plots the logged signal associated with the output of the Sum1 block for `BaselineRun_Scenario_3` and `EmbeddedRun_Scenario_3`, as well as their difference and the tolerance specified for this signal.

See Also

“Propose Data Types For Merged Simulation Ranges” on page 42-54 | “Control Views in the Fixed-Point Tool” on page 39-13 | “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9

Use Custom Data Type Override Settings for Range Collection

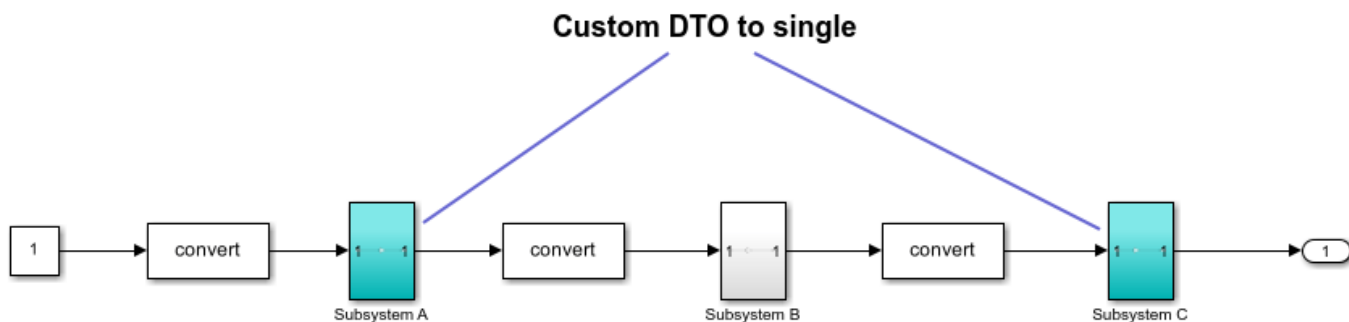
This example shows how to specify custom data type override settings to use during the range collection step in the Fixed-Point Tool.

By default, the Fixed-Point Tool honors the data types and any data type override specified on the model. You can use the Fixed-Point Tool to override data types in your model with doubles, singles, or scaled doubles. To specify custom data type overrides for elements within your model, use the `set_param` function.

Load a Simple Model

Open the `fxp_custom_dto` model. The model consists of three subsystems. Update the diagram (**Ctrl+D**) to display the data types currently set on the model.

```
open_system('fxp_custom_dto')
```



Copyright 2020 The MathWorks, Inc.

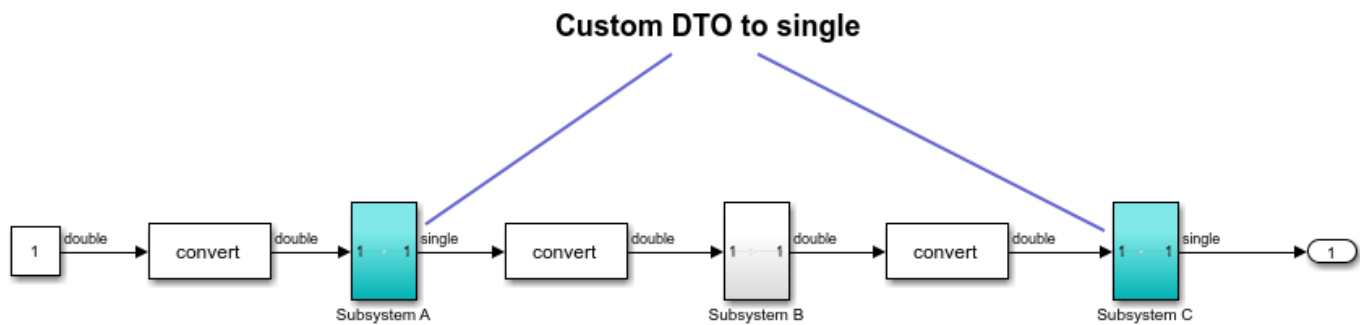
The original model uses the double-precision data type for each of the three subsystems within the model.

Override Data Types for a Subsystem

To override the data types in Subsystem A and Subsystem C with singles, use the `set_param` function:

```
set_param('fxp_custom_dto/Subsystem A', 'DataTypeOverride', 'Single');
set_param('fxp_custom_dto/Subsystem C', 'DataTypeOverride', 'Single');
```

Update the diagram and inspect the model to confirm that the data type override has been applied.



Copyright 2020 The MathWorks, Inc.

Collect Ranges Using the Fixed-Point Tool

In the **Apps** tab of the `fxp_custom_dto` model, select **Fixed-Point Tool**.

In the Fixed-Point Tool, select **New > Range Collection**. Under **System Under Design (SUD)**, select `fxp_custom_dto`. Under **Range Collection Mode**, select **Simulation Ranges**.

Under **Collect Ranges**, select **Use current settings**. Click **Collect Ranges**.

The Fixed-Point Tool collects ranges via simulation using the current data type override applied to your model. In this example, the data types of Subsystem A and Subsystem C are overrode with singles, and Subsystem B remains in double precision.

Verify Data Type Override Settings

To verify that the custom data type override settings specified using the `set_param` function were applied to the model during the range collection run, inspect the **Results** spreadsheet in the Fixed-Point Tool.

The compiled data type (**CompiledDT**) column for `BaselineRun` shows that Subsystem A and Subsystem C used the `single` data type, while the rest of the model was simulation using the `double` data type.

The screenshot displays the Range Collection tool interface. The 'Results' table lists simulation components and their data types. The 'Visualization of Simulation Data' section shows histograms for all results, with a legend indicating 'In-Range' (blue) and 'Underflows' (yellow) categories. The 'RESULT DETAILS' panel shows 'fxp_custom_dto/Subsystem A/Gain' with a 'Specified Data Type' of 'Inherit: Inherit via internal rule' and 'Range Information' showing a simulation range from 1 to 1.

| Name | CompiledDT | SpecifiedDT | SimMin | SimMax |
|-----------------------|------------|---------------------------|--------|--------|
| Data Type Conversion | double | Inherit: Inherit via b... | 1 | 1 |
| Data Type Conversion1 | double | Inherit: Inherit via b... | 1 | 1 |
| Data Type Conversion2 | double | Inherit: Inherit via b... | 1 | 1 |
| Subsystem A/Gain | single | Inherit: Inherit via i... | 1 | 1 |
| Subsystem B/Gain1 | double | Inherit: Inherit via i... | 1 | 1 |
| Subsystem C/Gain1 | single | Inherit: Inherit via i... | 1 | 1 |

| Property | Specified Data Type |
|-----------|------------------------------------|
| Data Type | Inherit: Inherit via internal rule |
| Minimum | |
| Maximum | |
| Precision | |

| Property | Minimum | Maximum |
|------------|---------|---------|
| Simulation | 1 | 1 |

| Values | Potential Overflows | In-Range | Potential Underflows |
|----------|---------------------|----------|----------------------|
| Positive | 0 | 1 | 0 |
| Negative | 0 | 0 | 0 |
| Zero | 0 | 0 | 0 |

Data Type Override for a Model Reference Hierarchy

When you simulate a model that contains referenced models, the data type override settings for the top-level model do not control the settings for the referenced models. You must specify these settings separately for referenced models, and if the settings must be consistent. For example, if you set the top-level model data type override setting to double and the referenced model to use local settings, and the referenced model uses fixed-point data types, then data type propagation issues might occur.

When you change the data type override settings for any instance of a referenced model, the settings change on all instances of the model and on the referenced model itself.

See Also

“Fixed-Point Instrumentation and Data Type Override” on page 42-61 | “Convert a Referenced Model to Fixed Point” on page 39-7

Working with the MATLAB Function Block

- “Convert MATLAB Function Block to Fixed Point” on page 45-2
- “Replace Functions in a MATLAB Function Block with a Lookup Table” on page 45-9
- “Best Practices for Working with the MATLAB Function Block in Automated Fixed-Point Conversion Workflows” on page 45-12
- “Control Data Types and Generate Code with MATLAB Function Block” on page 45-13
- “Specify Fixed-Point Math Properties in MATLAB Function Block” on page 45-19
- “Generate Fixed-Point FIR Code Using MATLAB Function Block” on page 45-26

Convert MATLAB Function Block to Fixed Point

This example shows how to use the Fixed-Point Tool to convert a model containing a MATLAB Function block to fixed point.

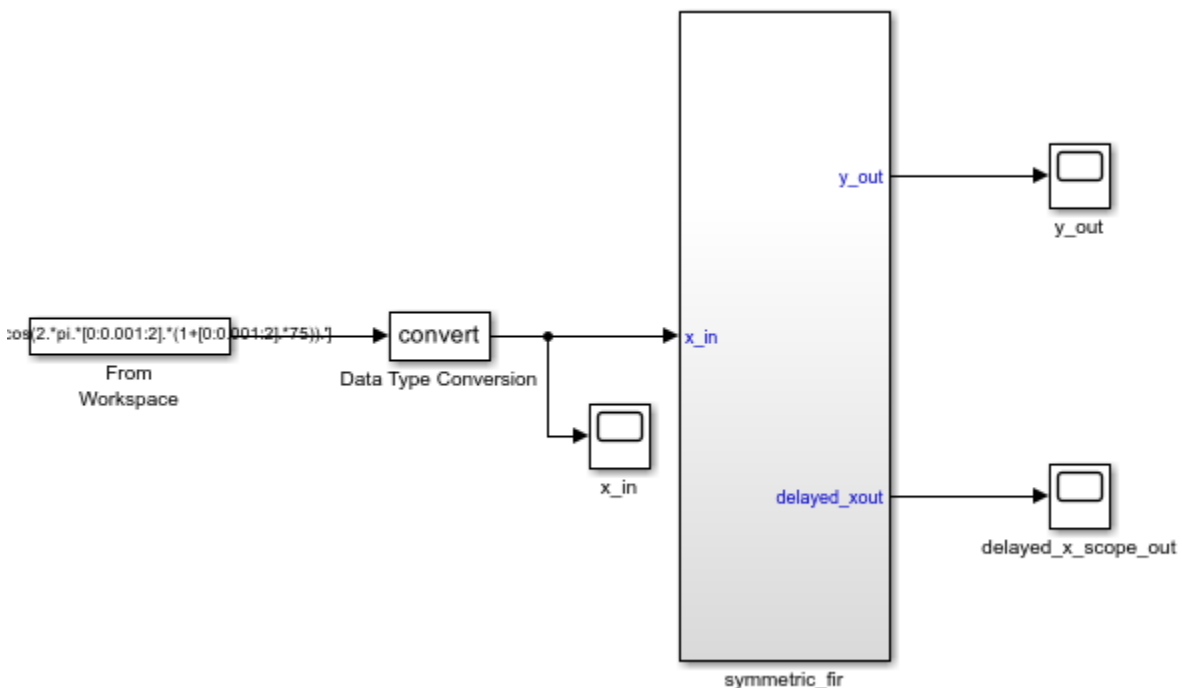
Best Practices for Working with the MATLAB Function Block in the Fixed-Point Tool

- Do not edit the fixed-point variant of your MATLAB Function block algorithm. Use the code view to edit the floating-point variant of your MATLAB code and re-propose and apply data types.
- For a successful conversion, only use modeling constructs supported for automated fixed-point conversion. For a list of the supported modeling constructs, see “MATLAB Language Features Supported for Automated Fixed-Point Conversion” on page 7-35.
- While collecting range information, do not edit the MATLAB code in the MATLAB Function block. Editing the code will cause problems if you try to merge results.
- During the fixed-point conversion process using the Fixed-Point Tool, do not use the “Save as” option to save the MATLAB Function block with a different name. If you do, you might lose existing results for the original block.

Open the Model

Open the `ex_symm_fir` model.

```
open_system("ex_symm_fir.slx")
```



The `ex_symm_fir` model uses a symmetric FIR filter. Simulate the model and inspect the model output. Inspect the symmetric FIR filter algorithm by double-clicking the MATLAB Function block.

Prepare for Fixed-Point Conversion

- 1 To open the Fixed-Point Tool, in the **Apps** tab, expand the **Apps** gallery and select **Fixed-Point Tool**.
- 2 In the Fixed-Point Tool, expand the **New** button arrow and select **Iterative Fixed-Point Conversion**.
- 3 Under **System Under Design (SUD)**, select the `symmetric_fir` subsystem, which contains the MATLAB Function block, as the system to convert.
- 4 Under **Range Collection Mode**, select **Simulation ranges** as the method of range collection. This configures the model to collect ranges using idealized floating-point data types.
- 5 In the **Prepare** section of the toolstrip, click **Prepare**.

Collect Range Information

Collect idealized ranges to use for data type proposal. Expand the **Collect Ranges** button arrow and select **Double precision**. Click **Collect Ranges** to start the simulation.

The Fixed-Point Tool stores the simulation data in a run titled **BaselineRun**. Examine the range information of the MATLAB variables in the spreadsheet.

Propose Data Types

Configure the proposal settings and propose fixed-point data types for the model.

- 1 In the **Convert** section of the toolstrip you can configure the data type proposal settings for the MATLAB Function block variables.

In this example, use the default proposal settings.

- 2 Click **Propose Data Types**.

The data type proposals appear in the **ProposedDT** column of the spreadsheet.

Note The **SpecifiedDT** column is always blank for MATLAB Function block variables.

The screenshot displays the MATLAB Function Block tool interface. The main window is titled "ITERATIVE FIXED-POINT CONVERSION" and "EXPLORE". The interface includes a toolbar with buttons for "New", "Prepare", "Collect Ranges", "MATLAB Functions", "Propose Data Types", "Apply Data Types", "Simulate with Embedded Types", "Run to compare in SDI", "Compare Results", and "Restore Original Model".

The "RESULTS" table shows the following data:

| Name | CompiledDT | SpecifiedDT | ProposedDT | Accept | SimMin | SimMax |
|--------------------------------|------------|-------------|------------------|--------|-----------------|-----------------|
| MATLAB Function/sfir : a1 | double | | numerictype(1... | ✓ | -1.999289346... | 1.9977514492... |
| MATLAB Function/sfir : a2 | double | | numerictype(1... | ✓ | -1.992314672... | 1.9998072404... |
| MATLAB Function/sfir : a3 | double | | numerictype(1... | ✓ | -1.986123539... | 1.9997334214... |
| MATLAB Function/sfir : a4 | double | | numerictype(1... | ✓ | -1.998136979... | 1.9999771889... |
| MATLAB Function/sfir : a5 | double | | numerictype(1... | ✓ | -0.434971836... | 0.4223737819... |
| MATLAB Function/sfir : a6 | double | | numerictype(1... | ✓ | -1.214431497... | 1.2178942850... |
| MATLAB Function/sfir : dela... | double | | numerictype(1... | ✓ | -0.999987366... | 1 |
| MATLAB Function/sfir : h_in1 | double | | numerictype(1... | ✓ | -0.1339 | -0.1339 |
| MATLAB Function/sfir : h_in2 | double | | numerictype(1... | ✓ | -0.0838 | -0.0838 |
| MATLAB Function/sfir : h_in3 | double | | numerictype(0... | ✓ | 0.2026 | 0.2026 |
| MATLAB Function/sfir : h_in4 | double | | numerictype(0... | ✓ | 0.4064 | 0.4064 |

The "Visualization of Simulation Data" section shows a histogram of all results in the model. The y-axis is labeled "Histogram Bins" and ranges from 2^{-40} to 2^2 . The x-axis is labeled "Data Values" and ranges from 2^{-40} to 2^2 . The histogram shows a distribution of values, with a significant peak at 2^{-40} . A legend indicates that the histogram is color-coded by status: Overflows (red), Representable (grey), In-Range (blue), and Underflows (yellow).

The "Proposed Data Type Summary" table shows the following data:

| Property | Proposed Data Type | Compil |
|-----------|----------------------|------------|
| Data Type | numerictype(1,16,13) | double |
| Minimum | -4 | -1.7976931 |
| Maximum | 3.9998779296875 | 1.7976931 |
| Precision | 0.0001220703125 | 4.9406564 |

The "Ranges used for proposal" table shows the following data:

| Property | Minimum | Maximum |
|------------|-----------------|-----------------|
| Simulation | -1.999289346... | 1.9977514492... |

The "Visualization of Simulation Data using numerictype(1,16,13)" section shows a histogram of the data values. The y-axis is labeled "% Occurrences" and ranges from 0 to 30. The x-axis is labeled "Data Values" and ranges from 2^{-40} to 2^2 . The histogram shows a distribution of values, with a significant peak at 2^{-40} .

The "Values" table shows the following data:

| Values | Potential Overflows | In-Range | Potential Underflows |
|----------|---------------------|----------|----------------------|
| Positive | 0 | 1022 | 0 |
| Negative | 0 | 978 | 0 |
| Zero | 0 | 1 | 0 |

Inspect Code Using the Code View

To launch the code view, click the **MATLAB Functions** button.

Using the code view you can:

- View detailed variable and expression information.
- Adjust proposal settings, such as `fimath` settings.
- Edit proposed data types.
- Manage function replacements.

For examples showing how to replace MATLAB functions with a lookup table, see “Replace Functions in a MATLAB Function Block with a Lookup Table” on page 45-9.

- Edit your code.
- Propose fixed-point data types.
- Apply proposed data types to your code.

To view the current proposal settings, click **Settings**. Here you can edit the `fimath` properties for the function. For this example, the default `fimath` properties are sufficient.

System Under Design: symmetric_fir

```

21 %#codegen
22 function [y_out, delayed_xout] = sfir(x_in, h_in1, h_in2, h_in3, h_in4)
23 % Symmetric FIR Filter
24
25 % declare and initialize the delay registers
26 persistent ud1 ud2 ud3 ud4 ud5 ud6 ud7 ud8;
27 if isempty(ud1)
28     ud1 = 0; ud2 = 0; ud3 = 0; ud4 = 0; ud5 = 0; ud6 = 0; ud7 = 0; ud8 = 0;
29 end
30
31 % access the previous value of states/registers
32 a1 = ud1 + ud8; a2 = ud2 + ud7;
33 a3 = ud3 + ud6; a4 = ud4 + ud5;
34
35 % multiplier chain
36 m1 = h_in1 * a1; m2 = h_in2 * a2;
37 m3 = h_in3 * a3; m4 = h_in4 * a4;
38
39 % adder chain
40 a5 = m1 + m2; a6 = m3 + m4;
41

```

Variables | Function Replacements

Show data for run: Ranges(Double) Only show runs with data [Go to converted code](#)

| Variable | Type | Sim Min | Sim Max | Proposed Type |
|---------------|--------|---------|---------|------------------------|
| Input | | | | |
| x_in | double | -1 | 1 | numerictype(1, 16, 14) |
| h_in1 | double | -0.13 | -0.13 | numerictype(1, 16, 17) |
| h_in2 | double | -0.08 | -0.08 | numerictype(1, 16, 18) |
| h_in3 | double | 0.2 | 0.2 | numerictype(0, 16, 18) |
| h_in4 | double | 0.41 | 0.41 | numerictype(0, 16, 17) |
| Output | | | | |
| y_out | double | -1.22 | 1.22 | numerictype(1, 16, 14) |
| delayed_xout | double | -1 | 1 | numerictype(1, 16, 14) |

Apply Proposed Data Types

When you have finished examining the proposed types, editing proposal settings, and implementing any function replacements, apply the proposed data types to the model. You can apply the data types either from the code view, or from the Fixed-Point Tool.

In the code view window, click **Apply**. The left pane displays both the original floating-point MATLAB Function block, as well as a newly generated fixed-point variant MATLAB Function block.

Right-click on the MATLAB Function block node in the left pane. Select **Go to Block** to navigate to the MATLAB Function block in the model.

The screenshot shows the 'Fixed Point Tool - MATLAB Function Block Converter' window for a subsystem named 'symmetric_fir'. The interface is divided into three main sections:

- System Under Design:** A tree view on the left shows the hierarchy: MATLAB Function Blocks > symmetric_fir > MATLAB Function > MATLAB Function > sfir > MATLAB Function_FixPt [fi].
- Code Editor:** The main area displays the MATLAB code for the 'sfir' function, converted to fixed-point. The code includes comments and function definitions for delay registers, multiplier chains, and adder chains. Line numbers 19 through 46 are visible.
- Function Replacements:** A table at the bottom allows for selecting data types for function blocks. The 'Show data for run:' dropdown is set to 'Ranges(Double)'. There is an unchecked checkbox for 'Only show runs with data' and a 'Go to original code' button.

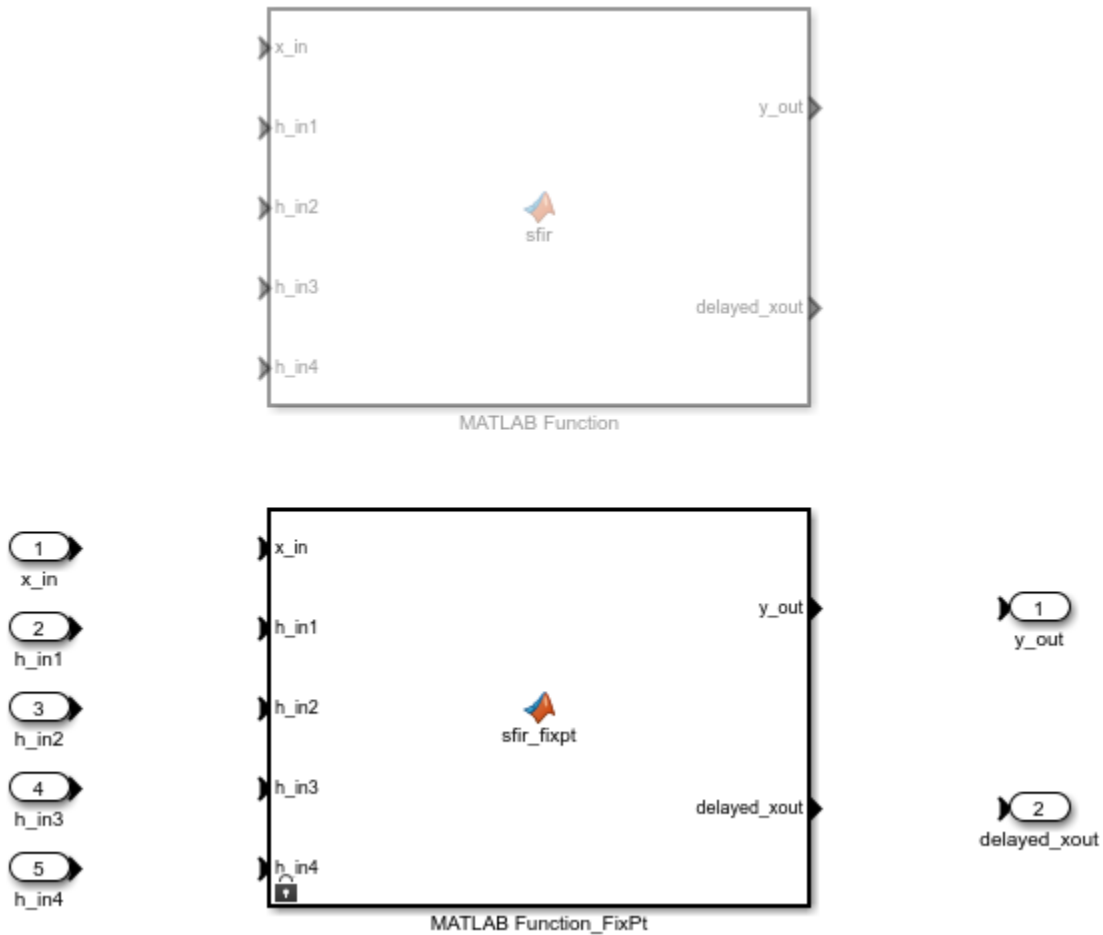
The MATLAB code in the editor is as follows:

```

19 % Copyright 2011-2015 The MathWorks, Inc.
20 %#codegen
21 function [y_out, delayed_xout] = sfir_fixpt(x_in, h_in1, h_in2, h_in3, h_in4)
22 % Symmetric FIR Filter
23
24 % declare and initialize the delay registers
25 fm = get_fimath();
26 h_in1 = fi(h_in1, 1, 16, 17, fm);
27 h_in2 = fi(h_in2, 1, 16, 18, fm);
28 h_in3 = fi(h_in3, 0, 16, 18, fm);
29 h_in4 = fi(h_in4, 0, 16, 17, fm);
30 x_in = fi(x_in, 1, 16, 14, fm);
31
32 persistent ud1 ud2 ud3 ud4 ud5 ud6 ud7 ud8;
33 if isempty(ud1)
34     ud1 = fi(0, 1, 16, 14, fm); ud2 = fi(0, 1, 16, 14, fm); ud3 = fi(0, 1, 16, 14, fm); ud4 =
35 end
36
37 % access the previous value of states/registers
38 a1 = fi(ud1 + ud8, 1, 16, 13, fm); a2 = fi(ud2 + ud7, 1, 16, 13, fm);
39 a3 = fi(ud3 + ud6, 1, 16, 13, fm); a4 = fi(ud4 + ud5, 1, 16, 13, fm);
40
41 % multiplier chain
42 m1 = fi(h_in1 * a1, 1, 16, 16, fm); m2 = fi(h_in2 * a2, 1, 16, 17, fm);
43 m3 = fi(h_in3 * a3, 1, 16, 16, fm); m4 = fi(h_in4 * a4, 1, 16, 15, fm);
44
45 % adder chain
46 a5 = fi(m1 + m2, 1, 16, 16, fm); a6 = fi(m3 + m4, 1, 16, 14, fm);

```

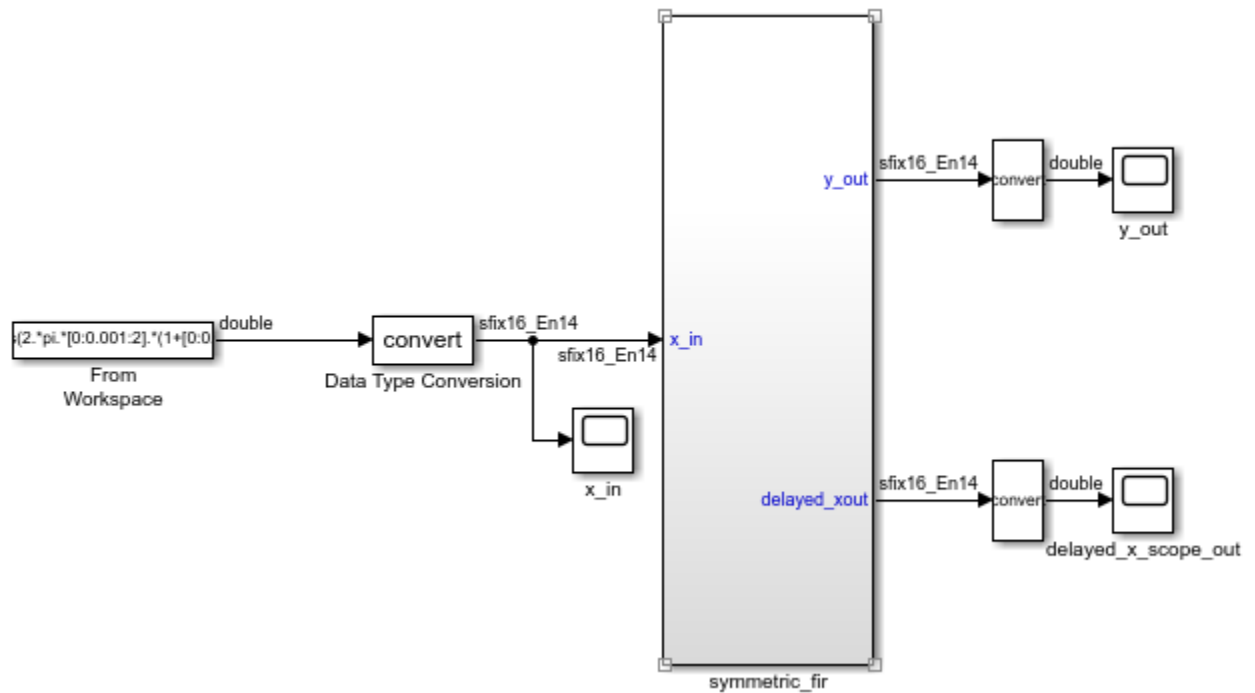
A variant subsystem is now in the place of the MATLAB Function block. The variant subsystem contains both floating-point and fixed-point versions of the MATLAB Function block. The active version is automatically controlled by the Fixed-Point Tool based on the data type override settings of the model. Data Type Override is not currently active on the model, so the fixed-point version is active.



Verify Results

Return to the Fixed-Point Tool to verify the results of the conversion.

In the **Verify** section of the toolstrip, click the **Simulate with Embedded Types** button to simulate the model using the newly applied fixed-point data types. The model simulates with the fixed-point variant as the active variant.



See Also

Related Examples

- “Replace Functions in a MATLAB Function Block with a Lookup Table” on page 45-9
- “Best Practices for Working with the MATLAB Function Block in Automated Fixed-Point Conversion Workflows” on page 45-12
- “MATLAB Language Features Supported for Automated Fixed-Point Conversion” on page 7-35

Replace Functions in a MATLAB Function Block with a Lookup Table

In this section...

“Open the Model” on page 45-9

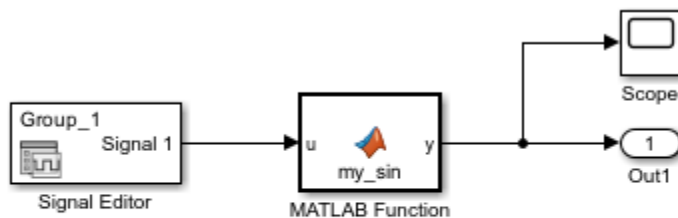
“Replace Sine Function with Lookup Table Approximation” on page 45-9

This example shows how to replace a function that is used inside a MATLAB Function block with a more efficient implementation. Use the Code View to replace the built-in `sin` function with a lookup table approximation.

Open the Model

Open the `ex_mySin` model.

```
open_system("ex_mySin.slx")
```



Copyright 2023 The MathWorks, Inc.

This model contains a MATLAB Function block which computes the sine of the input.

```
function y = my_sin(u)
%#codegen
y = sin(u);
```

Replace Sine Function with Lookup Table Approximation

- 1 To open the Fixed-Point Tool, in the **Apps** tab, expand the **Apps** gallery and select **Fixed-Point Tool**
- 2 In the Fixed-Point Tool, expand the **New** button arrow and select **Iterative Fixed-Point Conversion**.
- 3 Under **System Under Design (SUD)**, select the model `ex_mySin` as the system to convert.
- 4 Under **Range Collection Mode**, select **Simulation ranges** as the method of range collection. This configures the model to collect ranges using idealized floating-point data types.
- 5 In the **Prepare** section of the toolstrip, click **Prepare**.
- 6 Expand the **Collect Ranges** button arrow and select **Double precision**. Click **Collect Ranges** to start the simulation.

The Fixed-Point Tool stores the simulation data in a run titled `BaselineRun`. Examine the range information of the MATLAB variables in the spreadsheet.

- 7 To launch the code view, in the **Convert** section of the toolstrip, click **MATLAB Functions**.
- 8 Select the **Function Replacements** tab.
- 9 Enter the name of the function you want to replace. For this example, enter `sin`. Select **Lookup Table**, and then click **+**.

The fixed-point conversion process infers the ranges for the function and then uses an interpolated lookup table to replace the function. By default, the lookup table uses linear interpolation, 1000 points, and the minimum and maximum values detected by running the test file.

- 10 Click **Propose** to get data type proposals for the variables.
- 11 Click **Apply** to apply the data type proposals and generate a fixed-point lookup table.

System Under Design: `ex_mySin` SETTINGS ▾ PROPOSE APPLY ?

```

13
14 % calculate replacement_sin via lookup table between extents x = fi([-0.01,4.94]),
15 % interpolation degree = 1, number of points = 1000
16 function y = replacement_sin( x )
17     persistent LUT
18     if ( isempty(LUT) )
19         T = numerictype( true, 16, 14);
20         LUT = fi([-0.00999983333416666, -0.00504502364358801, -9.00900899682239e-05, ...
21             0.00486484567550096, 0.00981966200157351, 0.0147742372399358, ...
22             0.0197284497481933, 0.024682177892857, 0.0296353000523299, ...
23             0.0345876946198926, 0.0395392400066895, 0.0444898146447135, ...
24             0.0494392969897909, 0.0543875655245655, 0.0593344987614818, ...
25             0.0642799752457683, 0.0692238735584187, 0.0741660723191732, ...
26             0.0791064501894988, 0.084044885875568, 0.0889812581312369, ...
27             0.0939154457610221, 0.0988473276230759, 0.103776782632161, ...
28             0.108703689762622, 0.113627928051361, 0.118549376600799, ...
29             0.123467914581854, 0.128383421236901, 0.133295775882739, ...
30             0.138204857913553, 0.143110546803876, 0.148012722111549, ...
31             0.152911263480676, 0.157806050644578, 0.16269696342875, ...
32             0.167583881753809, 0.172466685638439, 0.177345255202343, ...
33             0.182219470669183, 0.187089212369518, 0.191954360743747, ...
34             0.196814796345041, 0.201670399842279, 0.206521052022973, ...
35             0.211366633796198, 0.216207026195516, 0.221042110381896, ...
36             0.225871767646631, 0.230695879414254, 0.235514327245448, ...
37             0.240326992839954, 0.245133758039476, 0.249934504830582, ...

```

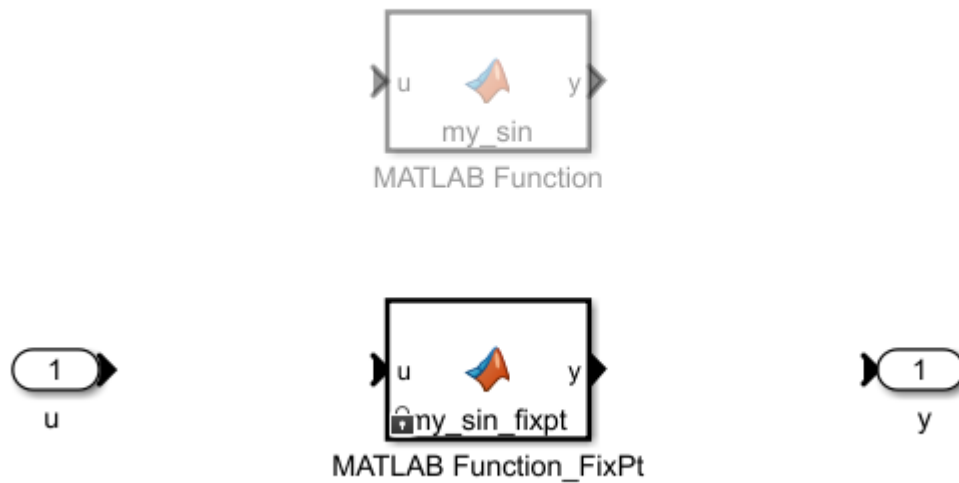
Variables **Function Replacements**

Enter a function to replace: Lookup Table ▾ + -

| Function or Operator | Replacement | Interpolation Method | Design Min | Design Max | Number of Points |
|----------------------|--------------|----------------------|------------|------------|------------------|
| sin | Lookup Table | Linear | Auto | Auto | 1000 |

If the behavior of the generated fixed-point code does not match the behavior of the original code closely enough, modify the interpolation method or number of points used in the lookup table and then regenerate the fixed-point code.

- 12** Return to the Fixed-Point Tool. In the **Verify** section of the toolstrip, click the **Simulate with Embedded Types** button to simulate the model using the newly applied fixed-point data types. The model simulates with the fixed-point variant as the active variant.



See Also

Related Examples

- "Convert MATLAB Function Block to Fixed Point" on page 45-2
- "Best Practices for Working with the MATLAB Function Block in the Fixed-Point Tool" on page 45-2
- "MATLAB Language Features Supported for Automated Fixed-Point Conversion" on page 7-35

Best Practices for Working with the MATLAB Function Block in Automated Fixed-Point Conversion Workflows

Follow these best practices when using MATLAB Function blocks with automated fixed-point conversion workflows, including the Fixed-Point Tool and `fxpopt`.

- Do not edit the fixed-point variant of your MATLAB Function block algorithm. Use the code view to edit the floating-point variant of your MATLAB code and reconvert to fixed-point data types.
- For a successful conversion, only use modeling constructs supported for automated fixed-point conversion. For a list of the supported modeling constructs, see “MATLAB Language Features Supported for Automated Fixed-Point Conversion” on page 7-35.
- While collecting range information, do not edit the MATLAB code in the MATLAB Function block. Editing the code will cause problems if you try to merge results.
- During the fixed-point conversion process using the Fixed-Point Tool, do not use the "Save As" option to save the MATLAB Function block with a different name. If you do, you might lose existing results for the original block.

Unsupported MATLAB Function Block Features

These MATLAB Function block features are not supported for automated fixed-point conversion workflows.

- Some MATLAB language features are not supported for automated fixed-point conversion. See “MATLAB Language Features Supported for Automated Fixed-Point Conversion” on page 7-35.
- Variable and struct names must be less than 63 characters long.

See Also

“MATLAB Language Features Supported for Automated Fixed-Point Conversion” on page 7-35

Related Examples

- “Convert MATLAB Function Block to Fixed Point” on page 45-2
- “Replace Functions in a MATLAB Function Block with a Lookup Table” on page 45-9

Control Data Types and Generate Code with MATLAB Function Block

In this section...

“Data Type Override with MATLAB Function Block” on page 45-13

“Fixed-Point Data Types with MATLAB Function Block” on page 45-14

“Share Models Containing Fixed-Point MATLAB Function Blocks” on page 45-17

You can use the MATLAB Function block to compose a MATLAB language function in a Simulink model that generates embeddable code. When you simulate the model or generate code for a target environment, a function in a MATLAB Function block generates efficient C/C++ code. This code meets the strict memory and data type requirements of embedded target environments.

Data Type Override with MATLAB Function Block

When you use the MATLAB Function block in a Simulink model that specifies data type override, the block determines the data type override equivalents of the input signal and parameter types. The block then uses these equivalent values to run the simulation. This table shows how the MATLAB Function block determines the data type override equivalent using the data type of the input signal or parameter and the data type override settings in the Simulink model. For more information about data type override, see “Fixed-Point Instrumentation and Data Type Override” on page 42-61.

| Input Signal or Parameter Type | Data Type Override Setting | Data Type Override Applies To Setting | Override Data Type |
|--------------------------------|----------------------------|---------------------------------------|--------------------|
| Inherited single | Double | All numeric types or Floating-point | Built-in double |
| | Single | All numeric types or Floating-point | Built-in single |
| | Scaled double | All numeric types or Floating-point | fi scaled double |
| Specified single | Double | All numeric types or Floating-point | Built-in double |
| | Single | All numeric types or Floating-point | Built-in single |
| | Scaled double | All numeric types or Floating-point | fi scaled double |
| Inherited double | Double | All numeric types or Floating-point | Built-in double |
| | Single | All numeric types or Floating-point | Built-in single |
| | Scaled double | All numeric types or Floating-point | fi scaled double |
| Specified double | Double | All numeric types or Floating-point | Built-in double |

| Input Signal or Parameter Type | Data Type Override Setting | Data Type Override Applies To Setting | Override Data Type |
|--------------------------------|----------------------------|---------------------------------------|--------------------|
| | Single | All numeric types or Floating-point | Built-in single |
| | Scaled double | All numeric types or Floating-point | fi scaled double |
| Inherited Fixed | Double | All numeric types or Fixed-point | fi double |
| | Single | All numeric types or Fixed-point | fi single |
| | Scaled double | All numeric types or Fixed-point | fi scaled double |
| Specified Fixed | Double | All numeric types or Fixed-point | fi double |
| | Single | All numeric types or Fixed-point | fi single |
| | Scaled double | All numeric types or Fixed-point | fi scaled double |

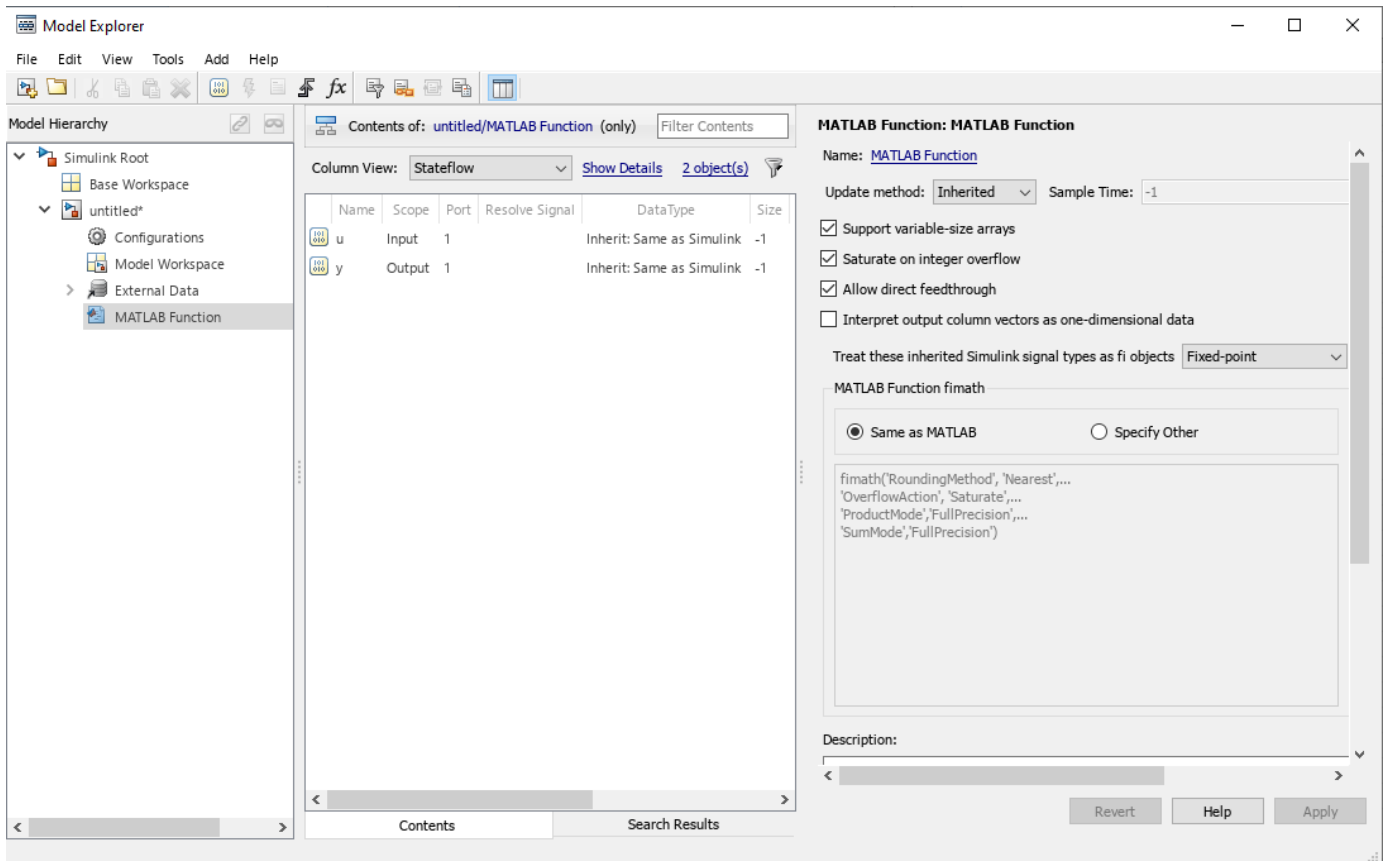
Fixed-Point Data Types with MATLAB Function Block

Code generation from MATLAB supports a significant number of Fixed-Point Designer functions. For information about which Fixed-Point Designer functions are supported, see “Functions Supported for Code Acceleration or C Code Generation” on page 12-4. To simulate models using fixed-point data types in Simulink, you must have a Fixed-Point Designer license.

Specify Fixed-Point Parameters in the Model Explorer

You can use the Model Explorer to specify parameters for a MATLAB Function block in a fixed-point model. For more information, see “Specify MATLAB Function Block Properties”.

- 1 Create a new model. Place a MATLAB Function block in the model.
- 2 Open the Model Explorer. On the **Modeling** tab, in the **Design** section, click **Model Explorer**.
- 3 Expand the **untitled*** node in the **Model Hierarchy** pane. Then, select the **MATLAB Function** node.



These parameters apply to MATLAB Function blocks in models that use fixed-point and integer data types:

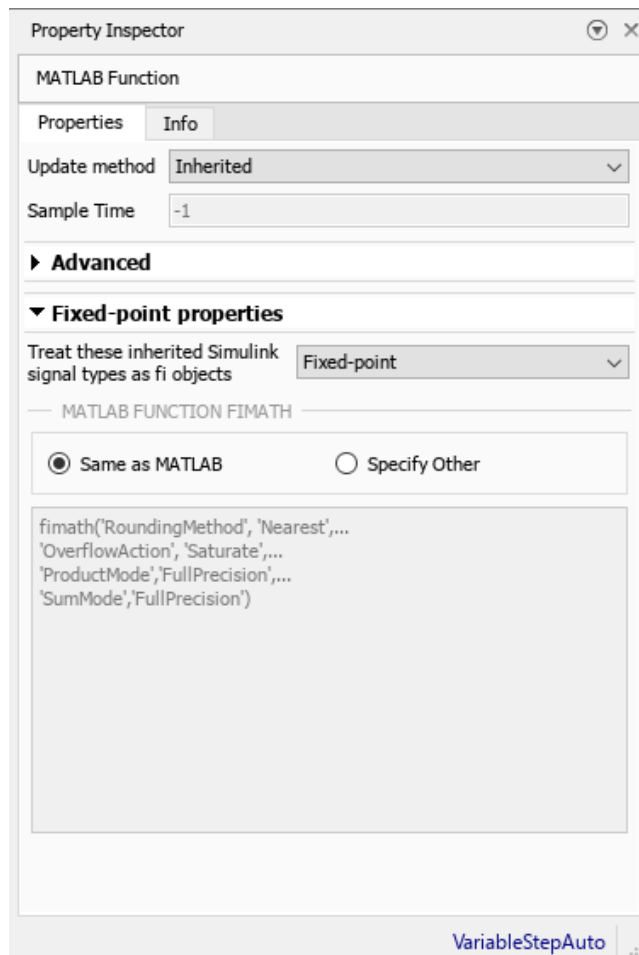
- **Treat these inherited Simulink signal types as fi objects** — Choose whether to treat inherited fixed-point and integer signals as `fi` objects.
 - **Fixed-point** — The MATLAB Function block treats all fixed-point inputs as Fixed-Point Designer `fi` objects.
 - **Fixed-point & Integer** — The MATLAB Function block treats all fixed-point and integer inputs as Fixed-Point Designer `fi` objects.
- **MATLAB Function fimath** — Specify the `fimath` properties for the block to associate with these objects:
 - All fixed-point and integer input signals to the MATLAB Function block that you choose to treat as `fi` objects.
 - All `fi` and `fimath` objects constructed in the MATLAB Function block.

Select one of these options:

- **Same as MATLAB** — The block uses the same `fimath` properties as the current default `fimath`. The edit box displays the current default `fimath` in read-only form.
- **Specify Other** — Specify your own `fimath` object in the edit box.

Use fimath Objects in MATLAB Function Blocks

Open the **Property Inspector** pane for the MATLAB Function block. On the **Modeling** tab, in the **Design** section, select **Property Inspector**. In the **Property Inspector** pane, expand **Fixed-point properties**.



The **MATLAB Function block fimath** parameter enables you to specify one set of fimath object properties for the MATLAB Function block. The block associates the fimath properties you specify with the following objects:

- All fixed-point and integer input signals to the MATLAB Function block that you choose to treat as fi objects.
- All fi and fimath objects constructed in the MATLAB Function block.

You can select one of the following options:

- **Same as MATLAB** — The block uses the same fimath properties as the current default fimath. The edit box displays the current default fimath in read-only form.
- **Specify Other** — Specify your own fimath object in the edit box. You can do this in two ways:
 - Construct the fimath object inside the edit box.

- Construct the `fimath` object in the MATLAB or model workspace and then enter its variable name in the edit box.

Note If you use this option and plan to share your model with others, make sure you define the variable in the model workspace.

For an example showing the **MATLAB Function `fimath`** options work, see “Specify Fixed-Point Math Properties in MATLAB Function Block” on page 45-19.

The Fixed-Point Designer `isfimathlocal` function supports code generation for MATLAB.

Share Models Containing Fixed-Point MATLAB Function Blocks

To share a fixed-point model containing a MATLAB Function block, you must first move any variables you define in the MATLAB workspace, including `fimath` objects, to the model workspace. For example:

- 1 Create a new model. Place a MATLAB Function block in the model.
- 2 Define a `fimath` object in the MATLAB workspace.

```
F = fimath('RoundingMethod','Floor','OverflowAction','Wrap',...
          'ProductMode','KeepLSB','ProductWordLength',32,...
          'SumMode','KeepLSB','SumWordLength',32)
```

F =

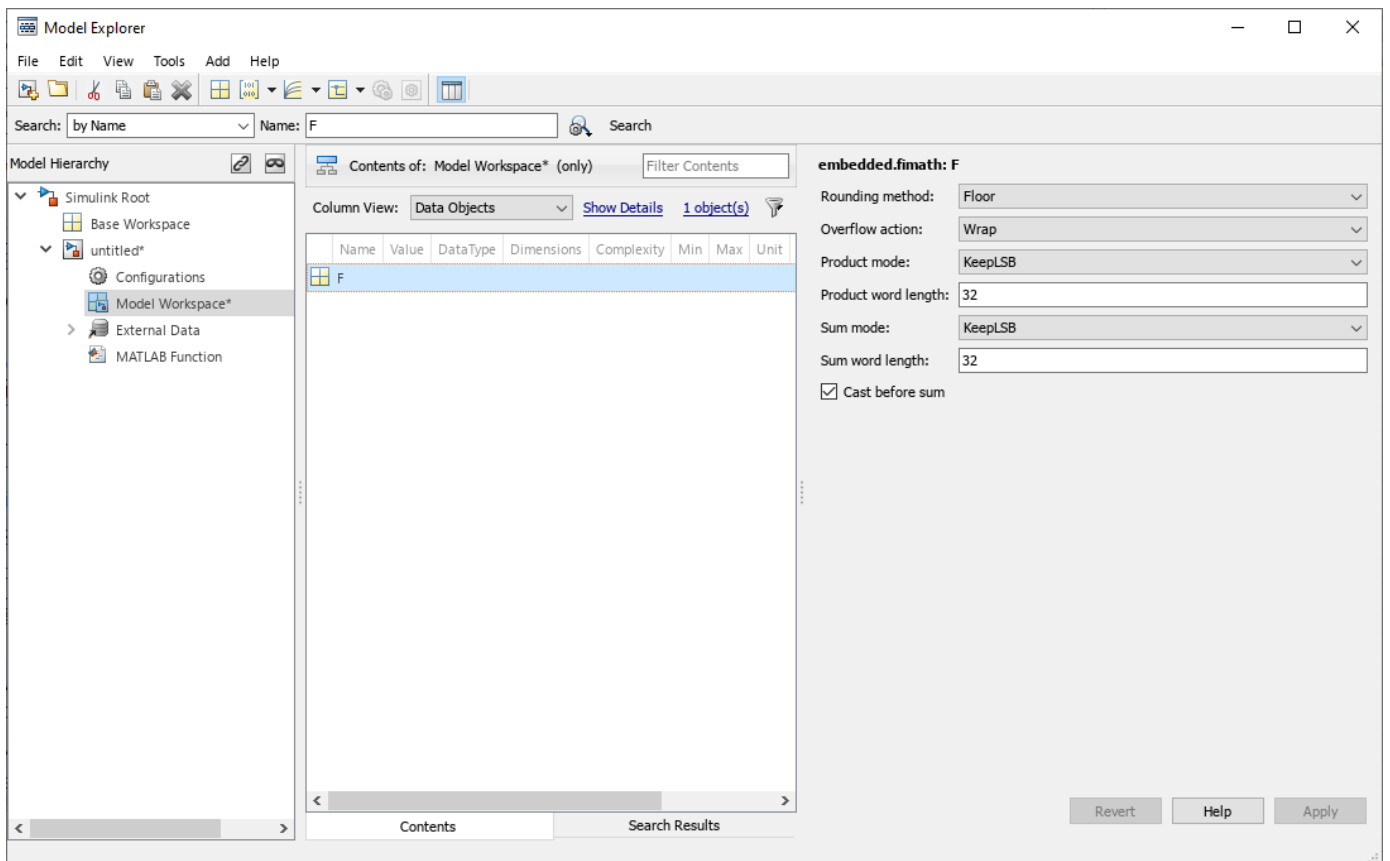
```
    RoundingMethod: Floor
    OverflowAction: Wrap
           ProductMode: KeepLSB
ProductWordLength: 32
           SumMode: KeepLSB
    SumWordLength: 32
    CastBeforeSum: true
```

You can use this `fimath` for any Simulink fixed-point signal entering the MATLAB Function block as an input.

- 3 Open the Model Explorer.
- 4 Expand the **untitled*** node in the **Model Hierarchy** pane of the Model Explorer. Select the **MATLAB Function** node.
- 5 For **MATLAB Function block `fimath`**, select **Specify other**. In the edit box, enter the variable `F`. Click **Apply** to save your changes.

You have now defined the `fimath` properties to be associated with all Simulink fixed-point input signals and all `fi` and `fimath` objects constructed within the block.

- 6 In the **Model Hierarchy** pane, select **Base Workspace**. You can see the variable `F` that you have defined in the MATLAB workspace listed in the **Contents** pane. If you send this model to another user, that user must first define that same variable in the MATLAB workspace to get the same results.
- 7 Cut the variable `F` from the base workspace and paste it into the model workspace listed under the node for your model, in this case, **untitled***. The Model Explorer now appears as shown.



You can share your model with another user. Because you included the required variables in the workspace of the model itself, another user can run the model and get the correct results. Receiving and running the model does not require any extra steps.

See Also

MATLAB Function | “Code Generation Workflow” (MATLAB Coder)

Related Examples

- “Implement MATLAB Functions in Simulink with MATLAB Function Blocks”
- “Specify Fixed-Point Math Properties in MATLAB Function Block” on page 45-19

Specify Fixed-Point Math Properties in MATLAB Function Block

This example shows how to specify fixed-point math properties in a MATLAB® Function block. You can specify `fimath` properties to add fixed-point numbers without bit growth. The same methods also apply to subtraction and multiplication.

Define Fixed-Point Variables with Attached `fimath`

Define fixed-point variables A and B with attached `fimath`.

```
clearvars

F = fimath(...
    'RoundMode', 'Fix', ...
    'OverflowMode', 'Wrap', ...
    'ProductMode', 'SpecifyPrecision', ...
    'ProductWordLength', 32, ...
    'ProductFractionLength', 16, ...
    'SumMode', 'SpecifyPrecision', ...
    'SumWordLength', 32, ...
    'SumFractionLength', 16, ...
    'CastBeforeSum', true);

A = fi(1, true, 32, 16, F);
B = fi(1, true, 32, 16, F);
```

Define a structure T containing prototypes of the variables A and B.

```
T.A = cast(0, 'like', A);
T.B = cast(0, 'like', B);
```

The structure T functions as a types table. The *values* of the fields of T are not important. You will use the *data types* of the fields of T later in this example to specify fixed-point types that carry the `fimath` along with them.

Add the fixed-point variables A and B. In MATLAB, the `fimath` attached to variables A and B specify that the sum is to be computed as 32-bit word length and 16-bit fraction length.

```
Y = A + B
```

```
Y =
    2

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 32
    FractionLength: 16

    RoundingMethod: Zero
    OverflowAction: Wrap
    ProductMode: SpecifyPrecision
    ProductWordLength: 32
    ProductFractionLength: 16
    SumMode: SpecifyPrecision
    SumWordLength: 32
    SumFractionLength: 16
    CastBeforeSum: true
```

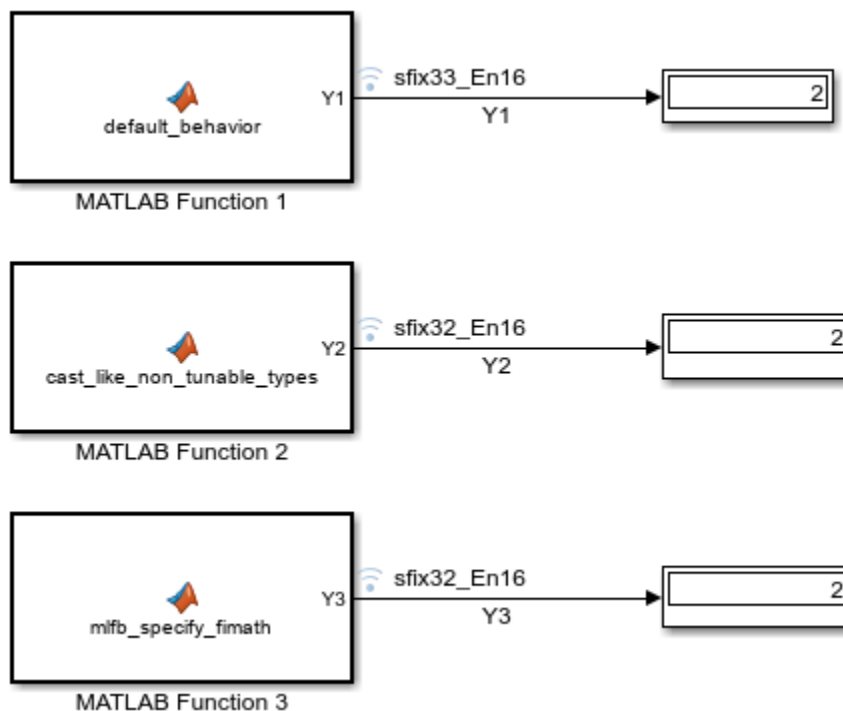
The `fimath` propagates to the variable `Y`.

Default MATLAB Function Block Behavior

In general, the Simulink software does not propagate `fimath` on fixed-point `fi` objects. This rule applies even if the `fi` objects have attached `fimath` in a Constant block or are passed in as a MATLAB Function block parameter. However, any attached `fimath` defined *inside* a MATLAB Function block are respected. One exception to the rule about parameters is described in the next section.

This function is defined in the block named MATLAB Function 1 in the `mParameterFIMath` model. If you execute this function in MATLAB, it returns the same 32-bit data type as the `Y = A + B` example.

```
function Y1 = default_behavior(A,B)
    Y1 = A + B;
end
```



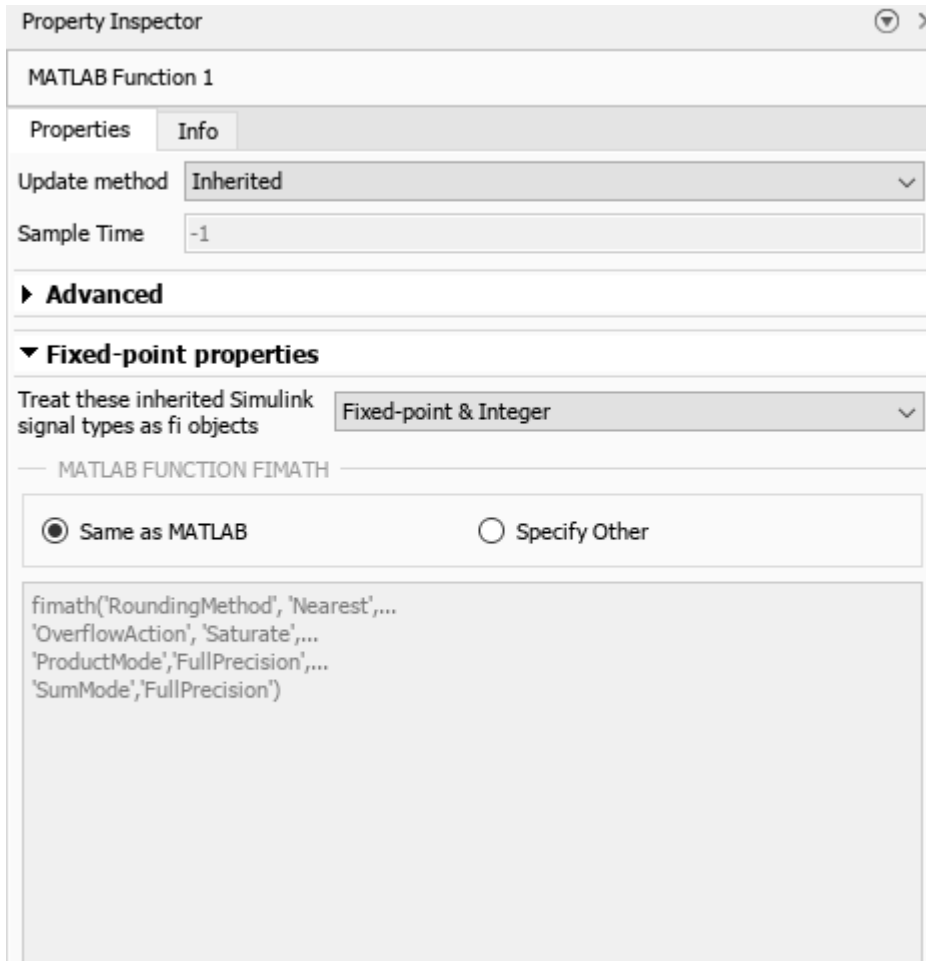
Simulate the `mParameterFIMath` model to see how the MATLAB Function block executes this code.

```
model = 'mParameterFIMath';
open_system(model);
sim(model);
Y1 = logouts.get('Y1').Values.Data
```

```
Y1 =
    2
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 33
    FractionLength: 16
```

The MATLAB Function 1 block returns Y1 with a 33-bit word length instead of the 32-bit word length returned in MATLAB. To see why, open the **Property Inspector** pane for the MATLAB Function 1 block.



The MATLAB Function 1 block returns a 33-bit word length output because its **MATLAB Function fimath** setting has been set to **Same as MATLAB**. The input parameters A and B are stripped of their attached `fimath` and instead use the default `fimath` settings from MATLAB. The default `fimath` in MATLAB does full-precision sums. Therefore, the sum of two 32-bit variables return a 33-bit result.

To see the default `fimath` settings in MATLAB, first reset any global `fimath` settings, then enter `fimath` in the MATLAB Command Window.

```
resetglobalfimath
fimath

ans =
    RoundingMethod: Nearest
    OverflowAction: Saturate
           ProductMode: FullPrecision
           SumMode: FullPrecision
```

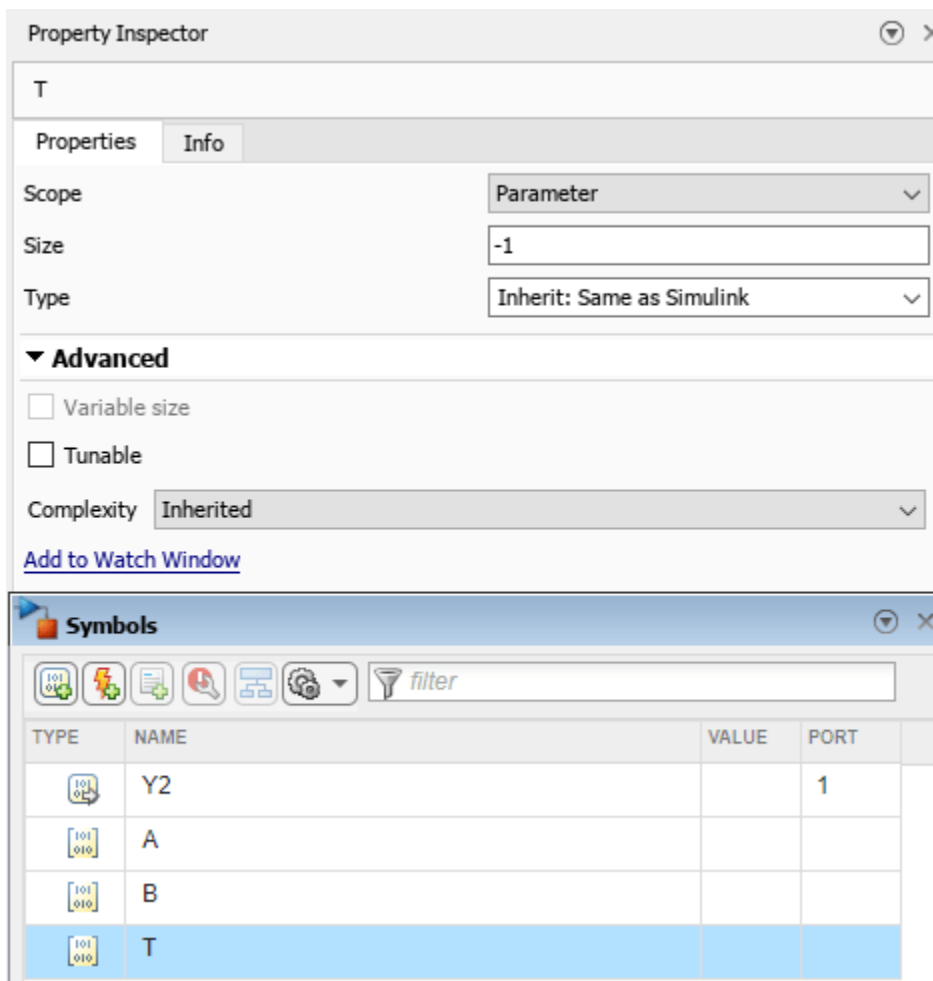
You can also change the value of `globalfimath`, but doing so changes the fixed-point math behavior globally. This method is not recommended.

Pass `fimath` into MATLAB Function Block from Parameter

In general, the MATLAB Function block strips `fimath` from fixed-point inputs. The exception to this rule is if an input parameter is a non-tunable structure, such as this one.

```
T.A = cast(0, 'like', A);
T.B = cast(0, 'like', B);
```

In the `mParameterFIMath` model, the MATLAB Function 2 block has the structure `T` defined as a non-tunable input parameter.



The fields `T.A` and `T.B` carry the data type and `fimath` of `A` and `B`. If you cast the inputs `A` and `B` like `T.A` and `T.B`, respectively, you recover the `fimath` that was defined in MATLAB.

```
function Y2 = cast_like_non_tunable_types(A, B, T)
    A = cast(A, 'like', T.A);
    B = cast(B, 'like', T.B);
    Y2 = A + B;
end
```

The output of the MATLAB Function 2 block has the desired 32-bit word length.

```
Y2 = logouts.get('Y2').Values.Data
Y2 =
    2
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 32
    FractionLength: 16
```

This workflow provides a robust way of defining fixed-point data types in MATLAB Function blocks because it allows for the definition of different `fimath` and data type for each different variable. This method has these advantages:

- The algorithm and type specification can be separate, with the types controlled in a dictionary separate from the block. Changing the type does not change the algorithm or block.
- MATLAB and MATLAB Function blocks in Simulink have identical behavior.
- Each parameter can have its own `fimath` and data type.
- Each parameter can change to be types other than fixed-point. For instance, `T.A = single(0); T.B = single(0);` would change all types in this example to single without having to use global data-type override settings.

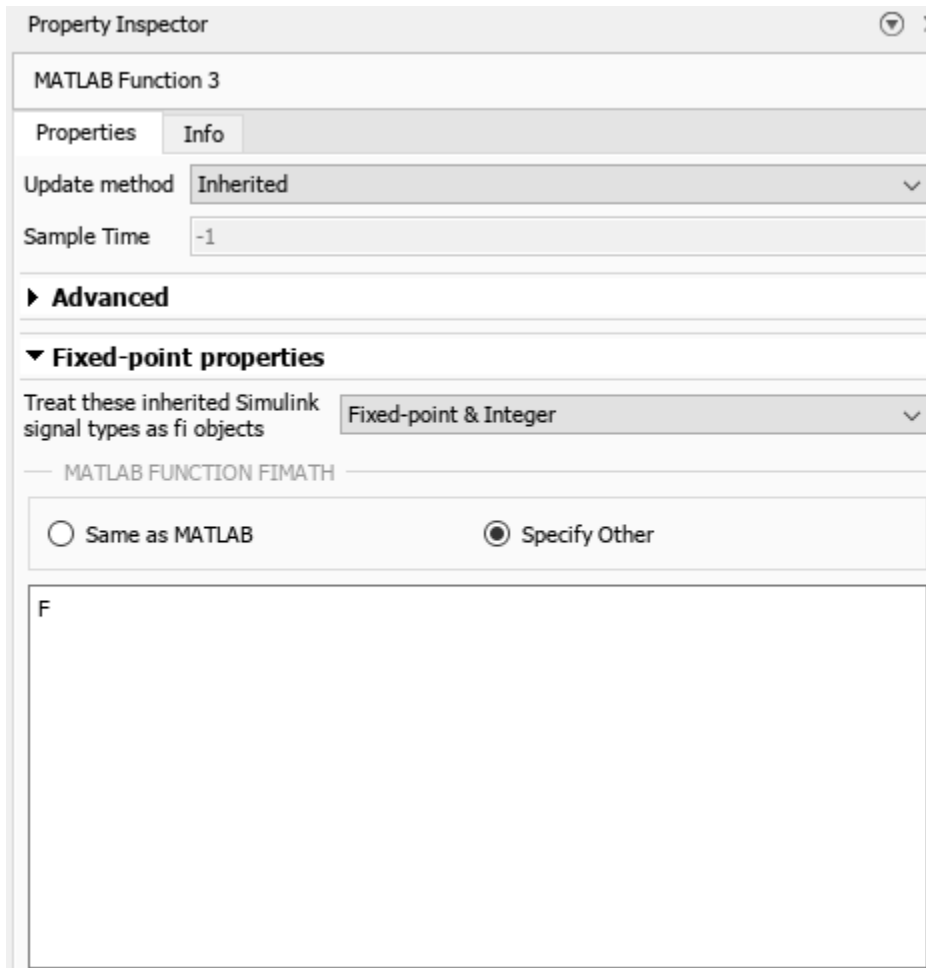
Set `fimath` for All Input Parameters

An alternative way to define `fimath` in a MATLAB Function block is to declare the `fimath` in the **Fixed-point properties** of the MATLAB Function block.

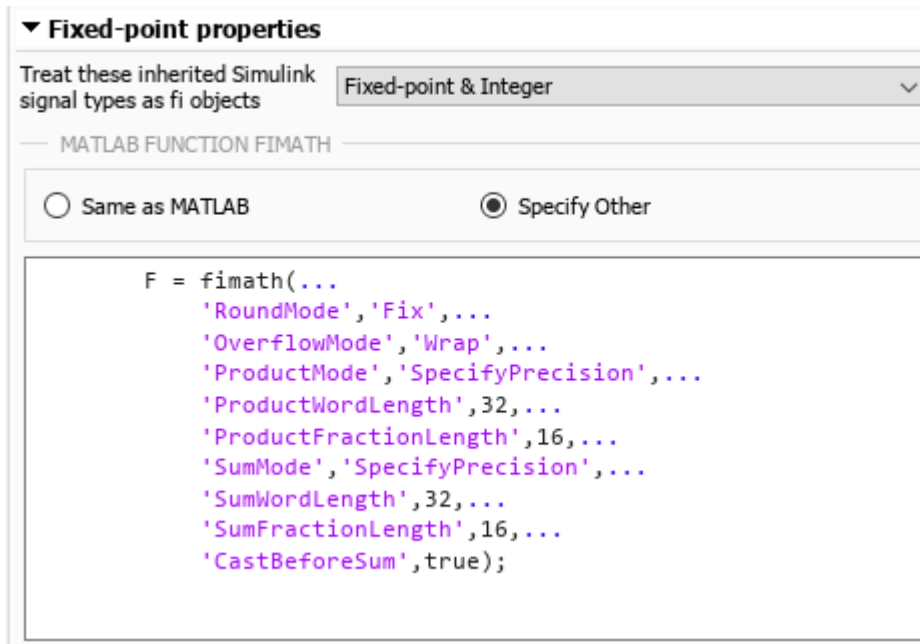
In the `mParameterFIMath` model, the MATLAB Function 3 block contains this code.

```
function Y3 = mlfb_specify_fimath(A, B)
    Y3 = A + B;
end
```

In the **Property Inspector** pane, under **Fixed-point properties > MATLAB Function `fimath`**, the option **Specify Other** is selected. The `fimath` is defined as the variable `F` from above.



Alternatively, you can write the `fimath` definition directly in the **MATLAB Function `fimath`** box.



Confirm that the output of this block has the desired 32-bit word length.

```
Y3 = logout.get('Y3').Values.Data
```

```
Y3 =  
    2
```

```
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 32  
    FractionLength: 16
```

This method has these limitations:

- All MATLAB Function block input parameters get the same `fimath`.
- Each MATLAB Function block must be modified to specify the `fimath`.

See Also

MATLAB Function | `fimath`

Related Examples

- “Control Data Types and Generate Code with MATLAB Function Block” on page 45-13
- “Implement MATLAB Functions in Simulink with MATLAB Function Blocks”

Generate Fixed-Point FIR Code Using MATLAB Function Block

In this section...

“Program the MATLAB Function Block” on page 45-26

“Prepare the Inputs” on page 45-26

“Create the Model” on page 45-27

“Define the fimath Object Using the Model Explorer” on page 45-28

“Run the Simulation” on page 45-28

Program the MATLAB Function Block

The following example shows how to create a fixed-point, lowpass, direct form FIR filter in Simulink. To create the FIR filter, you use Fixed-Point Designer software and the MATLAB Function block. In this example, you perform the following tasks in the sequence shown:

- 1 Place a MATLAB Function block in a new model. You can find the block in the Simulink User-Defined Functions library.
- 2 Save your model as `cgen_fi`.
- 3 Double-click the MATLAB Function block in your model to open the **MATLAB Function Block Editor**. Type or copy and paste the following MATLAB code, including comments, into the Editor:

```
function [yout,zf] = dffirdemo(b, x, zi) %#codegen
%codegen_fi doc model example
%Initialize the output signal yout and the final conditions zf
Ty = numericity(1,12,8);
yout = fi(zeros(size(x)), 'numericity', Ty);
zf = zi;

% FIR filter code
for k=1:length(x);
    % Update the states: z = [x(k);z(1:end-1)]
    zf(:) = [x(k);zf(1:end-1)];
    % Form the output: y(k) = b*z
    yout(k) = b*zf;
end

% Plot the outputs only in simulation.
% This does not generate C code.
figure;
subplot(211);plot(x); title('Noisy Signal');grid;
subplot(212);plot(yout); title('Filtered Signal');grid;
```

Prepare the Inputs

Define the filter coefficients b , noise x , and initial conditions zi by typing the following code at the MATLAB command line:

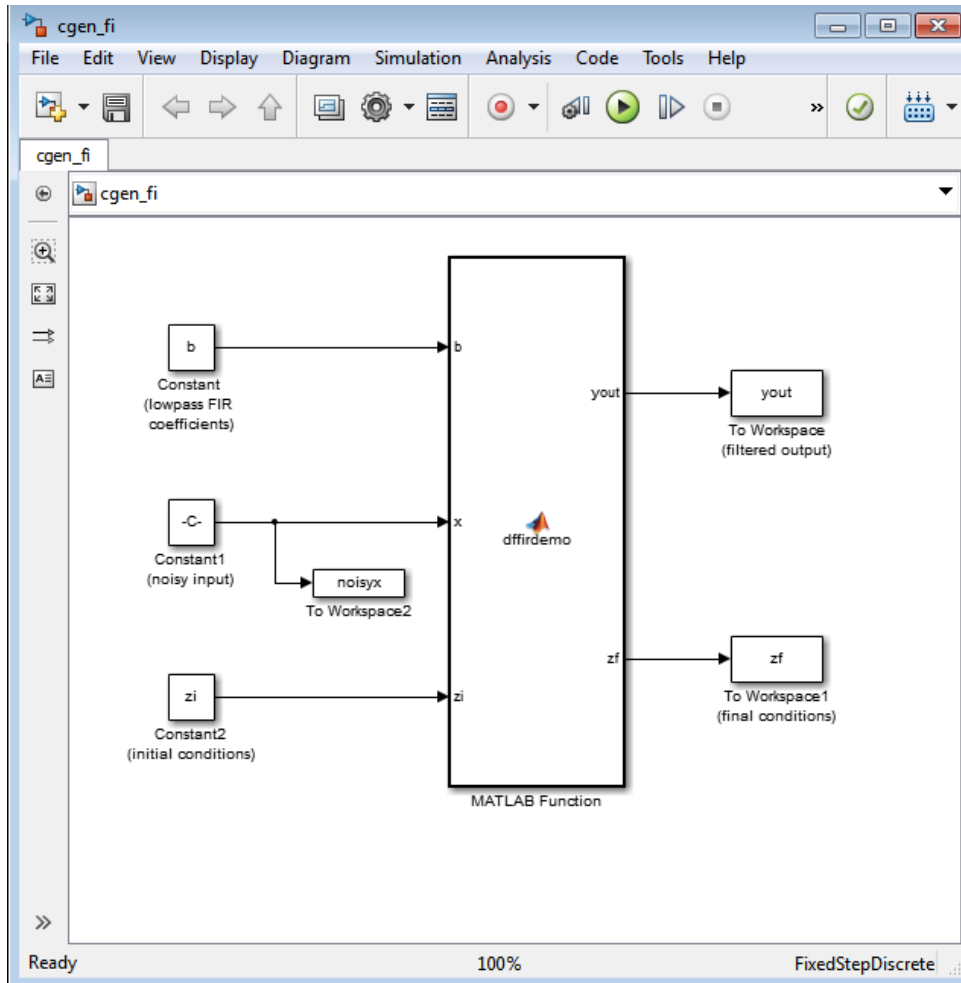
```
b=fidemo.fi_fir_coefficients;
load mtlb
x = mtlb;
n = length(x);
noise = sin(2*pi*2140*(0:n-1)'./Fs);
```



```
x = x + noise;
zi = zeros(length(b),1);
```

Create the Model

- 1 Add blocks to your model to create the following system.



- 2 Set the block parameters in the model to these “Fixed-Point FIR Code Example Parameter Values” on page 12-52.
- 3 On the **Modeling** tab, click **Model Settings**. Set the following configuration parameters.

| Parameter | Value |
|-----------|---------------------------------|
| Stop time | 0 |
| Type | Fixed-step |
| Solver | discrete (no continuous states) |

Click **Apply** to save your changes.

Define the fimath Object Using the Model Explorer

- 1 Open the Model Explorer for the model.
- 2 Click the **cgen_fi** > **MATLAB Function** node in the **Model Hierarchy** pane. The dialog box for the MATLAB Function block appears in the **Dialog** pane of the Model Explorer.
- 3 Select **Specify other** for the **MATLAB Function block fimath** parameter on the MATLAB Function block dialog box. You can then create the following **fimath** object in the edit box:

```
fimath('RoundingMethod','Floor','OverflowAction','Wrap',...
      'ProductMode','KeepLSB','ProductWordLength',32,...
      'SumMode','KeepLSB','SumWordLength',32)
```

The **fimath** object you define here is associated with fixed-point inputs to the MATLAB Function block as well as the **fi** object you construct within the block.

By selecting **Specify other** for the **MATLAB Function block fimath**, you ensure that your model always uses the **fimath** properties you specified.

Run the Simulation

- 1 Run the simulation by selecting your model and typing **Ctrl+T**. While the simulation is running, information outputs to the MATLAB command line. You can look at the plots of the noisy signal and the filtered signal.
- 2 Next, build embeddable C code for your model by selecting the model and typing **Ctrl+B**. While the code is building, information outputs to the MATLAB command line. A folder called `coder_fi_grt_rtw` is created in your current working folder.
- 3 Navigate to `coder_fi_grt_rtw > cgen_fi.c`. In this file, you can see the code generated from your model. Search for the following comment in your code:

```
/* codegen_fi doc model example */
```

This search brings you to the beginning of the section of the code that your MATLAB Function block generated.

Working with Data Objects in the Fixed-Point Workflow

- “Bus Objects in the Fixed-Point Workflow” on page 46-2
- “Autoscaling Data Objects Using the Fixed-Point Tool” on page 46-4

Bus Objects in the Fixed-Point Workflow

In this section...

“How Data Type Proposals Are Determined for Bus Objects” on page 46-2

“Bus Naming Conventions with Data Type Override” on page 46-3

“Limitations of Bus Objects in the Fixed-Point Workflow” on page 46-3

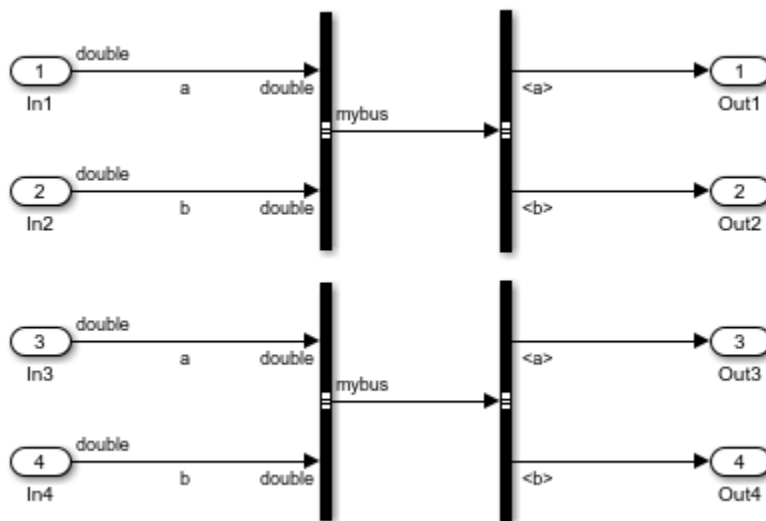
How Data Type Proposals Are Determined for Bus Objects

This example shows how the software determines the data types for elements of bus objects using the `ex_bus_range` model.

The data type proposal for a bus object is found by taking the union of the ranges of all sources driving the same bus element, and then proposing a data type for this range. The Fixed-Point Tool does not log minimum and maximum ranges for elements of a bus signal.

Open the `ex_bus_range` model and range data.

```
load bus_range.mat
open_system('ex_bus_range')
```



Each of the four input ports in this model have specified design ranges. The In2 and In4 input ports must share the same data type because they drive the same element of the `mybus` bus object.

The **Fixed-Point Tool** proposes a data type based on the union of these two ranges. After proposing data types for the model, select the **Data Objects** node in the **Model Hierarchy** pane. In the **Result Details** pane for the `mybus : b` element of the bus object, notice the row labeled **Shared Design** in

the **Ranges used for proposal** table. The proposed data type is based on this range, which is the union of the design ranges of the In2 and In4 blocks.

Bus Naming Conventions with Data Type Override

When you use data type override on a model that contains buses, the Fixed-Point Tool generates a new bus which uses the overridden data type. To indicate that a model is using an overridden bus, the tool adds a prefix to the name of the original bus object. While a model is in an overridden state, a bus object named myBus is renamed based on the following pattern.

| DTO Mode | DTO Applies To | | |
|---------------|-------------------|-----------------|-----------------|
| | All Numeric Types | Floating Point | Fixed Point |
| Scaled Double | dtoScl_myBus | dtoSclFlt_myBus | dtoSclFxp_myBus |
| Double | dtoDbl_myBus | dtoDblFlt_myBus | dtoDblFxp_myBus |
| Single | dtoSgl_myBus | dtoSglFlt_myBus | dtoSglFxp_myBus |

Note You cannot see bus objects with an overridden data type within the Type Editor because they are not stored in the base workspace.

Limitations of Bus Objects in the Fixed-Point Workflow

An update diagram error can occur if any of the following conditions occur.

- Your model is in accelerator mode and has a bus object with an overridden data type at the output port.

To perform data type override, run the model in normal mode.

- The data types in your model are overridden and the model contains Stateflow charts that use MATLAB as the action language.
- Your model contains tunable MATLAB structures assigned to a bus signal (such as Unit Delay blocks with a structure as the initial condition, Stateflow data, and MATLAB structures from the workspace).

To use the Fixed-Point Tool, change the structure to a non tunable structure. To avoid unnecessary quantization effects, specify the structure fields as doubles. For more information on using a structure as an initial condition with bus objects, see “Data Type Mismatch and Structure Initial Conditions” on page 49-28.

- Your model contains a structure parameter specified through the mask of an atomic subsystem.

To use the Fixed-Point Tool, make the system non-atomic.

Autoscaling Data Objects Using the Fixed-Point Tool

The **Fixed-Point Tool** generates a data type proposal for data objects based on ranges collected through simulation, derived range analysis, and design ranges specified on model objects. The **Fixed-Point Tool** also takes into consideration any data type constraints imposed by the model objects.

These types of data objects are supported for conversion using the **Fixed-Point Tool**.

- `Simulink.Parameter`
- `Simulink.Bus`
- `Simulink.NumericType`
- `Simulink.AliasType`
- `Simulink.Signal`
- `Simulink.LookupTable`
- `Simulink.Breakpoint`

The following sections describe how the tool collects the ranges and analyzes constraints.

Collecting Ranges for Data Objects

The objects in your model that use the same data object to specify its type must all share the same data type. The Fixed-Point Tool collects the ranges for all objects in your model. Objects that must share the same data type are placed in a data type group. The Fixed-Point Tool generates a data type proposal for the group based on the union of the ranges of all model objects in the group.

Collecting Ranges for Parameter Objects

Whenever possible, it is a best practice to specify design range information on the parameter object. When the data type of the parameter object is set to `auto`, the Fixed-Point Tool follows the same rules as when proposing for inherited data types. The Fixed-Point Tool determines the ranges to use for the data type proposal for a parameter object by taking the union of the parameter value, the parameter design ranges, and the design ranges of client blocks.

Data Type Constraints in Data Objects

Some objects in a shared data type group may contain constraints on the data types they can accept. For example, some blocks can accept only signed data types.

Autoscaling Parameter Objects

The Fixed-Point Tool is not able to detect when a parameter object must be integer only, such as when using a parameter object as a variable for dimensions, variant control, or a Boolean value. In these cases, you must clear the **Accept** box in the Fixed-Point Tool proposal stage before applying data types to your model.

Autoscaling Breakpoint Objects

Breakpoint data must always be strictly monotonically increasing. Although a breakpoint data set may be strictly monotonic in double format, due to saturation and quantization, it might not be after conversion to a fixed-point data type. The Fixed-Point Tool accounts for this behavior and proposes a

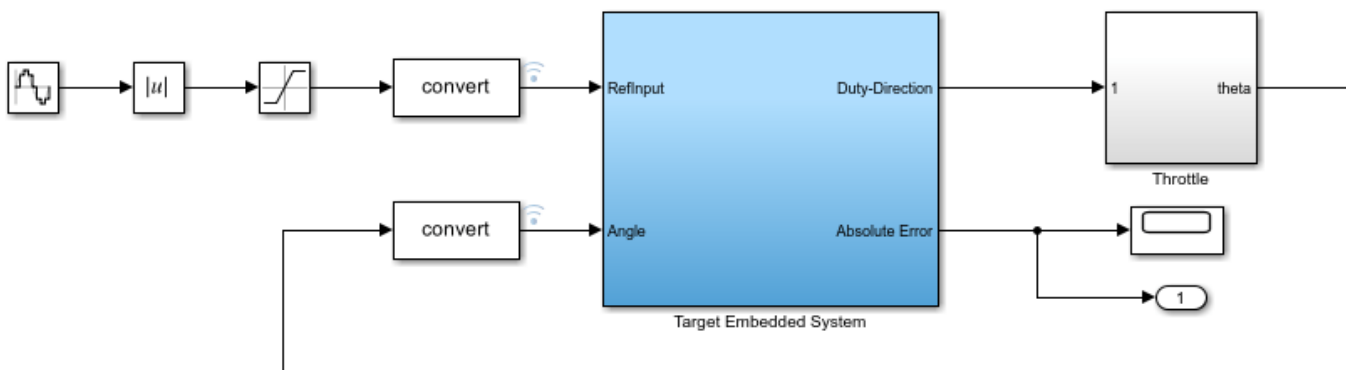
data type large enough to satisfy the monotonicity constraint after conversion. In some cases, the data type is very large in order to satisfy the constraint. In this case, consider editing your breakpoint data such that it can be represented efficiently in fixed point.

Autoscale a Model Using Data Objects for Data Type Definitions

The following model uses several different types of data objects, including `Simulink.Bus`, `Simulink.NumericType`, `Simulink.LookupTable`, and `Simulink.Breakpoint` objects for data type definition. Use the Fixed-Point Tool to convert the floating-point model, including the data objects used in the model, to fixed point.

Open the `ex_data_objects` Model

```
open_system('ex_data_objects')
```




Use the Fixed-Point Tool to Autoscale the Model


- 1 From the Simulink **Apps** tab, select **Fixed-Point Tool**.
- 2 In the Fixed-Point Tool, under **New** workflow, select **Iterative Fixed-Point Conversion**.
- 3 In the Fixed-Point Tool, under **System Under Design (SUD)**, select **Target Embedded System** as the system you want to convert.
- 4 Under **Range Collection Mode**, select **Simulation ranges**.
- 5 Click the **Prepare** button. The Fixed-Point Tool checks the system under design for compatibility with the conversion process and reports any issues found in the model.

When model objects within the system under design share a data type with objects outside of the system under design, data type propagation issues can occur after conversion to fixed point. For this reason, during the preparation stage of the conversion, the Fixed-Point Tool inserts Data Type Conversion blocks at the outputs of the system under design.

In this example, the tool is not able to automatically insert a Data Type Conversion block at the `ex_data_objects/Throttle` port because the port uses a bus signal. You can ignore this warning in this case because there are already Data Type Conversion blocks isolating this port inside the `Throttle` subsystem.

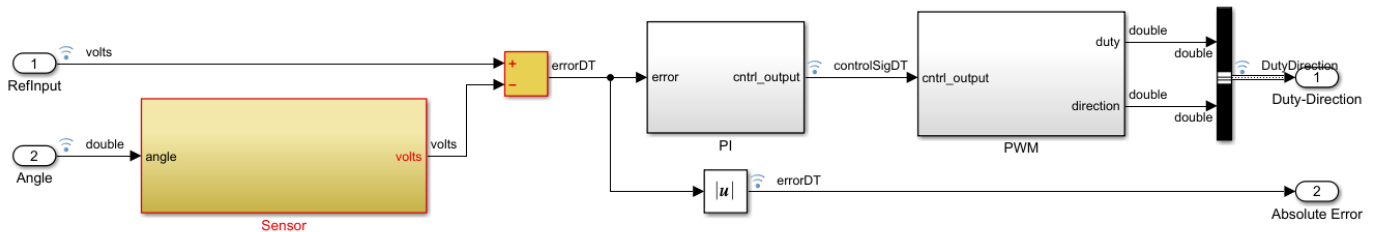
- 6 Expand the **Collect Ranges** button arrow and select **Double Precision**. Click the **Collect**

Ranges button  to start the simulation. The Fixed-Point Tool stores collected range information in a run titled `BaselineRun`.


- 7 In the **Convert** section, click the **Propose Data Types** button .

The Fixed-Point Tool detects data objects in the model and proposes a data type that satisfies the constraints of the data object. You can view all data objects used in a model by selecting **Data Objects** in the **Model Hierarchy** pane.

- 8 To learn more about a particular result, select the data object in the **Results** spreadsheet. The **Result Details** pane provides more details about the proposal, and gives a link to highlight all blocks in your model using a particular data object.



The tool displays the proposed data types for all results in the **ProposedDT** column of the **Results** spreadsheet.


- 9 To view the data type group that a result belongs to, add the **DTGroup** column to the spreadsheet. Click the add column button . Select **DTGroup** in the menu.

| Name | CompiledDT | SpecifiedDT | ProposedDT | Accept | Sim | |
|------------------|------------|----------------------|----------------|-------------------------------------|-----------------|-----------------|
| Abs1 | double | Inherit: Same ... | fixdt(1,16,11) | <input checked="" type="checkbox"/> | 0 | |
| Absolute Error | | Inherit: auto | n/a | | | |
| Add : Accumul... | double | double | fixdt(1,16,11) | <input checked="" type="checkbox"/> | -2 | |
| Add : Output | double | double | fixdt(1,16,11) | <input checked="" type="checkbox"/> | -2 | |
| Angle | | Inherit: auto | fixdt(0,16,15) | <input checked="" type="checkbox"/> | | |
| Duty-Direction | | Inherit: auto | n/a | | | |
| PI/Add : Accu... | double | Inherit: Inherit ... | n/a | | -0 | |
| PI/Add : Output | double | Inherit: Inherit ... | fixdt(1,16,18) | <input checked="" type="checkbox"/> | -0 | |
| PI/Add1 : Acc... | double | double | fixdt(1,16,14) | <input checked="" type="checkbox"/> | -0 | |
| PI/Add1 : Output | double | double | fixdt(1,16,14) | <input checked="" type="checkbox"/> | -0 | |
| PI/Gain | double | double | fixdt(1,16,15) | <input checked="" type="checkbox"/> | -0.666636995... | 0.0958976327... |

- CompiledDT
- SpecifiedDT
- ProposedDT
- Accept
- SimMin
- SimMax
- DesignMin
- DesignMax
- DerivedMin
- DerivedMax
- ProposedMin
- ProposedMax
- DTGroup

To sort by the **DTGroup** column, click the column header. You can now see results that must share the same data type next to each other.

10

Click the **Apply Data Types** button  to write the proposed data types to the model.

The Fixed-Point Tool applies the data type proposals to the data objects at their definition. In this example, the data objects are defined in the base workspace. View the details of a particular data object by entering the name of the data object at the MATLAB command line.

```
errorDT
```

```
NumericType with properties:
```

```
    DataTypeMode: 'Fixed-point: binary point scaling'  
    Signedness: 'Signed'  
    WordLength: 16  
    FractionLength: 11  
    IsAlias: 1  
    DataScope: 'Auto'  
    HeaderFile: ''  
    Description: ''
```

See Also

More About

- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9

Command Line Interface for the Fixed-Point Tool

- “The Command-Line Interface for the Fixed-Point Tool” on page 47-2
- “Convert a Model to Fixed Point Using the Command Line” on page 47-4

The Command-Line Interface for the Fixed-Point Tool

The methods of the `DataTypeWorkflow.Converter` class allow you to collect simulation and derived data, propose and apply data types to the model, and analyze results. The class performs the same fixed-point conversion tasks as the Fixed-Point Tool. The following table summarizes the steps in the workflow and lists the appropriate classes and methods to use at each step.

| Step in Workflow | Primary Objects and Object Functions for Step in Workflow |
|---|---|
| Set up model | <ul style="list-style-type: none"> • <code>DataTypeWorkflow.Converter</code> |
| Prepare the model for fixed-point conversion | <ul style="list-style-type: none"> • <code>applySettingsFromShortcut</code> • <code>applySettingsFromRun</code> |
| Gather range information | <ul style="list-style-type: none"> • <code>deriveMinMax</code> • <code>simulateSystem</code> |
| Propose data types | <ul style="list-style-type: none"> • <code>DataTypeWorkflow.ProposalSettings</code> • <code>addTolerance</code> • <code>clearTolerances</code> • <code>showTolerances</code> • <code>proposeDataTypes</code> • <code>proposalIssues</code> |
| Apply proposed data types | <ul style="list-style-type: none"> • <code>applyDataTypes</code> |
| Verify new fixed-point settings and analyze results | <ul style="list-style-type: none"> • <code>DataTypeWorkflow.Result</code> • <code>DataTypeWorkflow.VerificationResult</code> • <code>results</code> • <code>saturationOverflows</code> • <code>wrapOverflows</code> • <code>verify</code> • <code>explore</code> |

Note You should not use the Fixed-Point Tool and the command-line interface in the same conversion session.

To decide which workflow is right for you, consult the following table:

| Capability | Fixed-Point Tool | Command-Line Interface |
|---------------------------------------|------------------|------------------------|
| Populate runs to dataset | Supported | Supported |
| Delete result from dataset | Supported | Supported |
| Edit proposed data types | Supported | Not supported |
| Selectively apply data type proposals | Supported | Not supported |
| Run multiple simulations | Supported | Supported |
| Script workflow | Not supported | Supported |

See Also

Related Examples

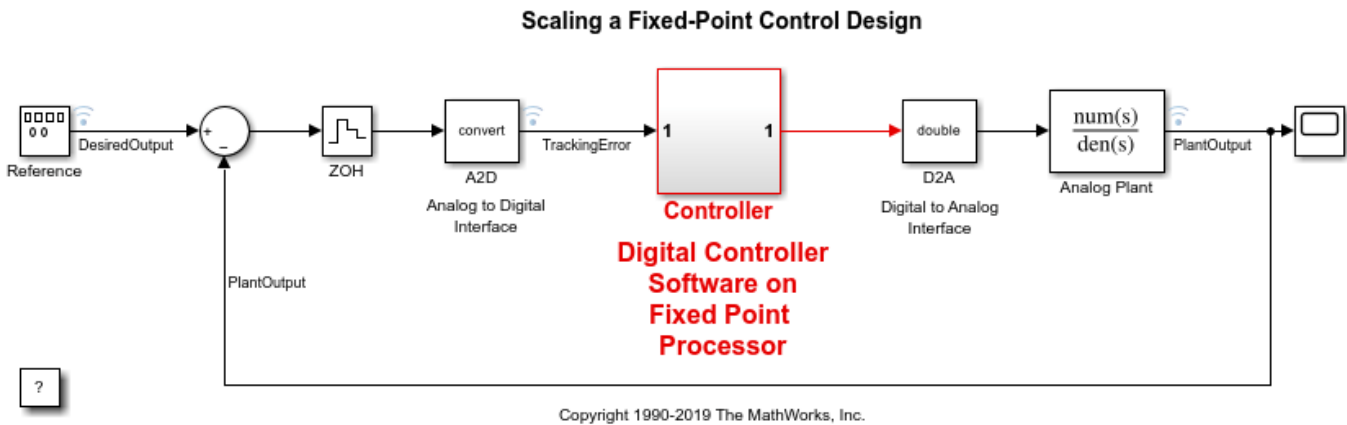
- “Convert a Model to Fixed Point Using the Command Line” on page 47-4

Convert a Model to Fixed Point Using the Command Line

This example shows how to refine the data types of a model using the command line.

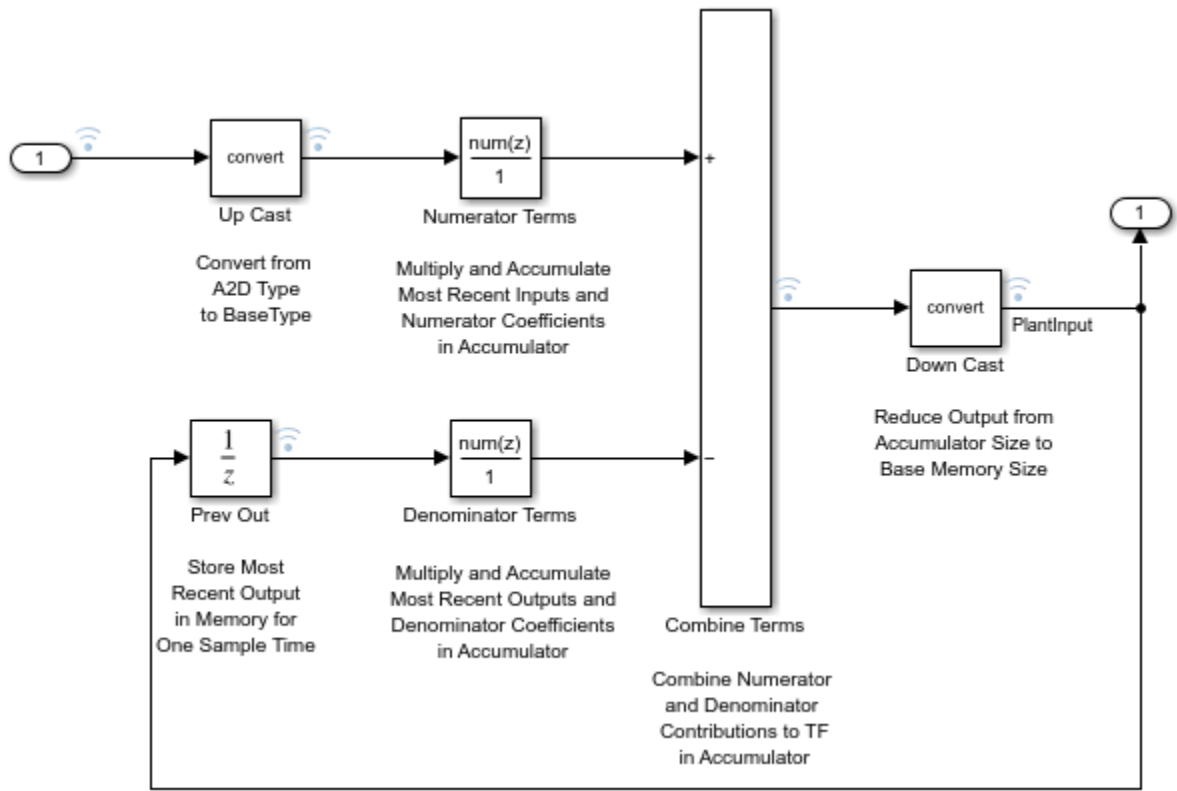
Open the `fxpdemo_feedback` model.

```
model = 'fxpdemo_feedback';
open_system(model);
```



The Controller subsystem uses fixed-point data types.

```
sud = 'fxpdemo_feedback/Controller';
open_system(sud)
```



Create a `DataTypeWorkflow.Converter` object to refine the data types of the Controller subsystem of the `fxpdemo_feedback` model.

```
converter = DataTypeWorkflow.Converter(sud);
```

Simulate the model and store the results in a run titled `InitialRun`.

```
converter.CurrentRunName = 'InitialRun';
converter.simulateSystem();
```

Determine any overflows occurred during the run.

```
 saturations = converter.saturationOverflows('InitialRun')
```

```
 saturations =
```

Result with properties:

```

    ResultName: 'fxpdemo_feedback/Controller/Up Cast'
  SpecifiedDataType: 'fixdt(1,16,14)'
   CompiledDataType: 'fixdt(1,16,14)'
   ProposedDataType: ''
         Wraps: []
   Saturations: 23
   WholeNumber: 0
         SimMin: -2
         SimMax: 1.9999
```

```
DerivedMin: []
DerivedMax: []
  RunName: 'InitialRun'
  Comments: {'An output data type cannot be specified on this result. The output type
DesignMin: []
DesignMax: []
```

```
wraps = converter.wrapOverflows('InitialRun')
```

```
wraps =
    []
```

A saturation occurs in the Up Cast block of the Controller subsystem during the simulation. There are no wrapping overflows. Refine the data types of the model so that there are no saturations.

Configure the model for conversion using a shortcut. Find the shortcuts that are available for the system by accessing the `ShortcutsForSelectedSystem` property of the converter object.

```
shortcuts = converter.ShortcutsForSelectedSystem
```

```
shortcuts =
    6x1 cell array
    {'Range collection using double override'      }
    {'Range collection with specified data types'  }
    {'Range collection using single override'     }
    {'Disable range collection'                  }
    {'Remove overrides and disable range collection'}
    {'Range collection using scaled double override'}
```

To collect idealized ranges for the system, using the 'Range collection using double override' shortcut, override the system with double-precision data types and enable instrumentation.

```
converter.applySettingsFromShortcut(shortcuts{1});
```

This shortcut also updates the current run name property of the converter object.

```
baselineRun = converter.CurrentRunName
```

```
baselineRun =
    'Ranges(Double)'
```

Simulate the model again to gather the idealized range information. These results are stored in the run `baselineRun`.

```
converter.simulateSystem();
```


Create a `ProposalSettings` object to control the data type proposal settings and specify tolerances for signals in the model.

```
propSettings = DataTypeWorkflow.ProposalSettings;
```

Specify a relative tolerance of 20% for the output signal of the `PlantOutput` signal in the model.

```
addTolerance(propSettings, 'fxpdemo_feedback/Analog Plant', 1, 'RelTol', 2e-1);
```

You can view all tolerances specified for a system using the `showTolerances` method.

```
showTolerances(propSettings)
```

| Path | Port_Index | Tolerance_Type | Tolerance_Value |
|-----------------------------------|------------|----------------|-----------------|
| {'fxpdemo_feedback/Analog Plant'} | 1 | {'RelTol'} | 0.2 |

Propose data types for the system using the proposal settings specified in `propSettings`, and the ranges stored in the `baselineRun` run.

```
converter.proposeDataTypes(baselineRun, propSettings)
```

Apply data types proposed for the `baselineRun` run to the model.

```
converter.applyDataTypes(baselineRun)
```

Verify that the behavior of the model using the new data types meets the tolerances specified on the proposal settings object, `propSettings`. The `verify` method removes the data type override and simulates the model using the updated fixed-point data types. It returns a `DataTypeWorkflow.VerificationResult` object.

```
result = verify(converter, baselineRun, 'FixedRun')
```

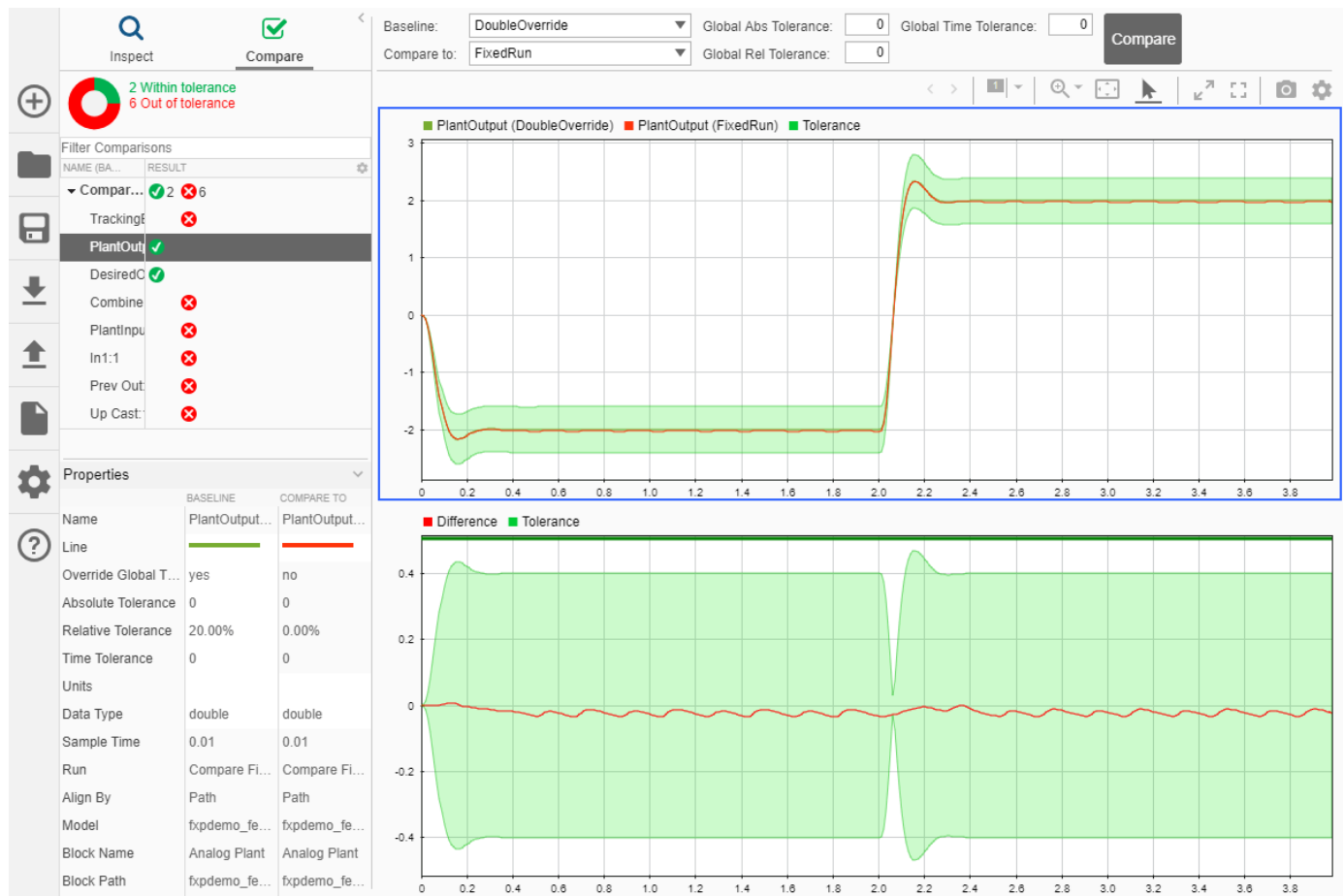
```
result =
```

```
VerificationResult with properties:
```

```
ScenarioResults: [0x0 DataTypeWorkflow.VerificationResult]
RunName: 'FixedRun'
BaselineRunName: 'Ranges(Double)'
Status: 'Pass'
MaxDifference: 0.0351
```

Using the `explore` method of the `DataTypeWorkflow.VerificationResult` object, launch the Simulation Data Inspector and examine the signals for which you specified a tolerance.

```
explore(result)
```



See Also

More About

- “The Command-Line Interface for the Fixed-Point Tool” on page 47-2

Code Generation

- “Fixed-Point Code Generation Support” on page 48-4
- “Accelerating Fixed-Point Models” on page 48-6
- “Using External Mode or Rapid Simulation Target” on page 48-7
- “Net Slope Computation” on page 48-8
- “Control the Generation of Fixed-Point Utility Functions” on page 48-16
- “Optimize Generated Code with the Model Advisor” on page 48-21
- “Lookup Table Optimization” on page 48-27
- “Selecting Data Types for Basic Operations” on page 48-29
- “Use of Shifts by C Code Generation Products” on page 48-31
- “Use Hardware-Efficient Algorithm to Solve Systems of Complex-Valued Linear Equations” on page 48-34
- “Implement Hardware-Efficient QR Decomposition Using CORDIC in a Systolic Array” on page 48-44
- “Perform QR Factorization Using CORDIC” on page 48-52
- “Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using QR Decomposition” on page 48-79
- “Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition” on page 48-82
- “Implement Hardware-Efficient Real Burst Matrix Solve Using QR Decomposition” on page 48-85
- “Implement Hardware-Efficient Real Burst Matrix Solve Using Q-less QR Decomposition” on page 48-88
- “Implement Hardware-Efficient Real Partial-Systolic QR Decomposition” on page 48-90
- “Implement Hardware-Efficient Real Partial-Systolic Q-less QR Decomposition” on page 48-93
- “Implement Hardware-Efficient Real Burst QR Decomposition” on page 48-96
- “Implement Hardware-Efficient Real Burst Q-less QR Decomposition” on page 48-99
- “Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using QR Decomposition” on page 48-102
- “Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition” on page 48-105
- “Implement Hardware-Efficient Complex Burst Matrix Solve Using QR Decomposition” on page 48-108
- “Implement Hardware-Efficient Complex Burst Matrix Solve Using Q-less QR Decomposition” on page 48-111
- “Implement Hardware-Efficient Complex Partial-Systolic QR Decomposition” on page 48-114
- “Implement Hardware-Efficient Complex Partial-Systolic Q-less QR Decomposition” on page 48-118
- “Implement Hardware-Efficient Complex Burst QR Decomposition” on page 48-121

- “Implement Hardware-Efficient Complex Burst Q-less QR Decomposition” on page 48-124
- “Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using QR Decomposition with Diagonal Loading” on page 48-127
- “Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using QR Decomposition with Diagonal Loading” on page 48-131
- “Determine Fixed-Point Types for QR Decomposition” on page 48-135
- “Determine Fixed-Point Types for Q-less QR Decomposition” on page 48-138
- “Algorithms to Determine Fixed-Point Types for Complex Q-less QR Matrix Solve $A'AX=B$ ” on page 48-140
- “Determine Fixed-Point Types for Complex Q-less QR Matrix Solve $A'AX=B$ ” on page 48-150
- “Algorithms to Determine Fixed-Point Types for Complex Least-Squares Matrix Solve $AX=B$ ” on page 48-154
- “Determine Fixed-Point Types for Complex Least-Squares Matrix Solve $AX=B$ ” on page 48-165
- “Algorithms to Determine Fixed-Point Types for Real Q-less QR Matrix Solve $A'AX=B$ ” on page 48-169
- “Determine Fixed-Point Types for Real Q-less QR Matrix Solve $A'AX=B$ ” on page 48-179
- “Algorithms to Determine Fixed-Point Types for Real Least-Squares Matrix Solve $AX=B$ ” on page 48-183
- “Determine Fixed-Point Types for Real Least-Squares Matrix Solve $AX=B$ ” on page 48-194
- “Compute Forgetting Factor Required for Streaming Input Data” on page 48-198
- “Estimate Standard Deviation of Quantization Noise of Complex-Valued Signal” on page 48-200
- “Estimate Standard Deviation of Quantization Noise of Real-Valued Signal” on page 48-202
- “Implement Hardware-Efficient Real Partial-Systolic Q-less QR with Forgetting Factor” on page 48-204
- “Implement Hardware-Efficient Complex Partial-Systolic Q-less QR with Forgetting Factor” on page 48-209
- “Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Forgetting Factor” on page 48-214
- “Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Forgetting Factor” on page 48-220
- “Implement Hardware-Efficient Complex Burst Matrix Solve Using QR Decomposition with Tikhonov Regularization” on page 48-226
- “Implement Hardware-Efficient Complex Burst Matrix Solve Using Q-less QR Decomposition with Tikhonov Regularization” on page 48-229
- “Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using QR Decomposition with Tikhonov Regularization” on page 48-234
- “Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Tikhonov Regularization” on page 48-237
- “Implement Hardware-Efficient Real Burst Matrix Solve Using QR Decomposition with Tikhonov Regularization” on page 48-242
- “Implement Hardware-Efficient Real Burst Matrix Solve Using Q-less QR Decomposition with Tikhonov Regularization” on page 48-245
- “Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using QR Decomposition with Tikhonov Regularization” on page 48-250

- “Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Tikhonov Regularization” on page 48-253
- “Determine Fixed-Point Types for Complex Least-Squares Matrix Solve with Tikhonov Regularization” on page 48-258
- “Determine Fixed-Point Types for Complex Q-less QR Matrix Solve with Tikhonov Regularization” on page 48-262
- “Determine Fixed-Point Types for Real Least-Squares Matrix Solve with Tikhonov Regularization” on page 48-266
- “Determine Fixed-Point Types for Real Q-less QR Matrix Solve with Tikhonov Regularization” on page 48-270
- “Implement Hardware-Efficient Real Burst Q-less QR with Forgetting Factor” on page 48-274
- “Implement Hardware-Efficient Complex Burst Q-less QR with Forgetting Factor” on page 48-279
- “Implement Hardware-Efficient Real Burst Matrix Solve Using Q-less QR Decomposition with Forgetting Factor” on page 48-285
- “Implement Hardware-Efficient Complex Burst Matrix Solve Using Q-less QR Decomposition with Forgetting Factor” on page 48-291
- “Implement Hardware-Efficient Real Burst Asynchronous Matrix Solve Using Q-less QR Decomposition” on page 48-297
- “Implement Hardware-Efficient Complex Burst Asynchronous Matrix Solve Using Q-less QR Decomposition” on page 48-300
- “How to Use Square Jacobi SVD HDL Optimized Block” on page 48-303
- “Implement HDL Optimized SVD in Feedforward Fashion Without Backpressure” on page 48-307
- “Implement HDL Optimized SVD with Backpressure Signal and HDL FIFO Block” on page 48-312

Fixed-Point Code Generation Support

In this section...

“Introduction” on page 48-4

“Languages” on page 48-4

“Data Types” on page 48-4

“Rounding Modes” on page 48-4

“Overflow Handling” on page 48-4

“Blocks” on page 48-4

“Scaling” on page 48-5

Introduction

All fixed-point blocks support code generation, except particular simulation features. The sections that follow describe the code generation support that the Fixed-Point Designer software provides. You must have a Simulink Coder license to generate C code or a HDL Coder license to generate HDL code.

Languages

C code generation is supported with the use of Simulink Coder. HDL code generation is supported with the use of HDL Coder.

Data Types

Fixed-point code generation supports all integer and fixed-point data types that are supported by simulation. Word sizes of up to 128 bits are supported in simulation. See “Supported Data Types” on page 34-13.

Rounding Modes

All rounding modes — Ceiling, Convergent, Floor, Nearest, Round, Simplest, and Zero — are supported.

Overflow Handling

- Saturation and wrapping are supported.
- Wrapping generates the most efficient code.
- Currently, you cannot choose to exclude saturation code automatically when hardware saturation is available. Select wrapping in order for the Simulink Coder product to exclude saturation code.

Blocks

All blocks generate code for all operations with a few exceptions. The Lookup Table Dynamic block generates code for all lookup methods except Interpolation-Extrapolation.

The Simulink Block Data Type Support table summarizes characteristics of blocks in the Simulink block library, including whether they support fixed-point data types and any limitations that apply for C code generation. To view the table, enter the following command at the MATLAB command line:

```
showblockdatatypetable
```

For information on block support for HDL code generation, see “Display Blocks for HDL Code Generation in Library Browser” (HDL Coder). You can also use the HDL Workflow Advisor to check your model for blocks not supported for HDL code generation.

Scaling

Any binary-point-only scaling and [Slope Bias] scaling that is supported in simulation is supported, bit-true, in code generation.

See Also

More About

- “Optimize Generated Code with the Model Advisor” on page 48-21

Accelerating Fixed-Point Models

If the model meets the code generation restrictions, you can use Simulink acceleration modes with your fixed-point model. The acceleration modes can drastically increase the speed of some fixed-point models. This is especially true for models that execute a very large number of time steps. The time overhead to generate code for a fixed-point model is generally larger than the time overhead to set up a model for simulation. As the number of time steps increases, the relative importance of this overhead decreases.

Note Rapid Accelerator mode does not support models with bus objects or 33+ bit fixed-point data types as parameters.

Every Simulink model is configured to have a start time and a stop time in the Configuration Parameters dialog box. Simulink simulations are usually configured for non-real-time execution, which means that the Simulink software tries to simulate the behavior from the specified start time to the stop time as quickly as possible. The time it takes to complete a simulation consists of two parts: overhead time and core simulation time, which is spent calculating changes from one time step to the next. For any model, the time it takes to simulate if the stop time is the same as the start time can be regarded as the overhead time. If the stop time is increased, the simulation takes longer. This additional time represents the core simulation time. Using an acceleration mode to simulate a model has an initially larger overhead time that is spent generating and compiling code. For any model, if the simulation stop time is sufficiently close to the start time, then Normal mode simulation is faster than an acceleration mode. But an acceleration mode can eliminate the overhead of code generation for subsequent simulations if structural changes to the model have not occurred.

In Normal mode, the Simulink software runs general code that can handle various situations. In an acceleration mode, code is generated that is tailored to the current usage. For fixed-point use, the tailored code is much leaner than the simulation code and executes much faster. The tailored code allows an acceleration mode to be much faster in the core simulation time. For any model, when the stop time is close to the start time, overhead dominates the overall simulation time. As the stop time is increased, there is a point at which the core simulation time dominates overall simulation time. Normal mode has less overhead compared to an acceleration mode when fresh code generation is necessary. Acceleration modes are faster in the core simulation portion. For any model, there is a stop time for which Normal mode and acceleration mode with fresh code generation have the same overall simulation time. If the stop time is decreased, then Normal mode is faster. If the stop time is increased, then an acceleration mode has an increasing speed advantage. Eventually, the acceleration mode speed advantage is drastic.

Normal mode generally uses more tailored code for floating-point calculations compared to fixed-point calculations. Normal mode is therefore generally much faster for floating-point models than for similar fixed-point models. For acceleration modes, the situation often reverses and fixed point becomes significantly faster than floating point. As noted above, the fixed-point code goes from being general to highly tailored and efficient. Depending on the hardware, the integer-based fixed-point code can gain speed advantages over similar floating-point code. Many processors can do integer calculations much faster than similar floating-point operations. In addition, if the data bus is narrow, there can also be speed advantages to moving around 1-, 2-, or 4-byte integer signals compared to 4- or 8-byte floating-point signals.

Using External Mode or Rapid Simulation Target

In this section...

“Introduction” on page 48-7

“External Mode” on page 48-7

“Rapid Simulation Target” on page 48-7

Introduction

If you are using the Simulink Coder external mode or rapid simulation (RSim) target, there are situations where you can get unexpected errors when tuning block parameters. These errors can arise when you specify the `Best precision` scaling option for blocks that support constant scaling for best precision.

The sections that follow provide further details about the errors you might encounter. To avoid these errors, specify a scaling value instead of using the `Best precision` scaling option.

External Mode

If you change a parameter such that the binary point moves during an external mode simulation or during graphical editing, and you reconnect to the target, a checksum error occurs and you must rebuild the code. When you use `Best Precision` scaling, the binary point is automatically placed based on the value of a parameter. Each power of two roughly marks the boundary where a parameter value maps to a different binary point. For example, a parameter value of 1-2 maps to a particular binary point position. If you change the parameter to a value of 2-4, the binary point moves one place to the right, while if you change the parameter to a value of 0.5-1, it moves one place to the left.

For example, suppose that a block has a parameter value of -2. You then build the code and connect in external mode. While connected, you change the parameter to -4. If the simulation is stopped and then restarted, this parameter change causes a binary point change. In external mode, the binary point is kept fixed. If you keep the parameter value of -4 and disconnect from the target, then when you reconnect, a checksum error occurs and you must rebuild the code.

Rapid Simulation Target

If a parameter change is great enough, and you are using the best precision mode for constant scaling, then you cannot use the RSim target.

If you change a block parameter by a sufficient amount (approximately a factor of two), the best precision mode changes the location of the binary point. Any change in the binary point location requires the code to be rebuilt because the model checksum is changed. This means that if best precision parameters are changed over a great enough range, you cannot use the rapid simulation target and a checksum error message occurs when you initialize the RSim executable.

Net Slope Computation

In this section...

“Handle Net Slope Computation” on page 48-8

“Use Division to Handle Net Slope Computation” on page 48-9

“Improve Numerical Accuracy of Simulation Results with Rational Approximations to Handle Net Slope” on page 48-9

“Improve Efficiency of Generated Code with Rational Approximations to Handle Net Slope” on page 48-12

“Use Integer Division to Handle Net Slope Computation” on page 48-15

Handle Net Slope Computation

The Fixed-Point Designer software provides an optimization parameter, **Use division for fixed-point net slope computation**, that controls how the software handles net slope computation. To learn how to enable this optimization, see “Use Integer Division to Handle Net Slope Computation” on page 48-15.

When a change of fixed-point slope is not a power of two, net slope computation is necessary. Normally, net slope computation is implemented using an integer multiplication followed by shifts. Under certain conditions, net slope computation can be implemented using integer division or a rational approximation of the net slope. One of the conditions is that the net slope can be accurately represented as a rational fraction or as the reciprocal of an integer. Under this condition, the division implementation gives more accurate numerical behavior. Depending on your compiler and embedded hardware, a division implementation might be more desirable than the multiplication and shifts implementation. The generated code for the rational approximation and/or integer division implementation might require less ROM or improve model execution time.

When to Use Division for Fixed-Point Net Slope Computation

This optimization works if:

- The net slope can be approximated with a fraction or is the reciprocal of an integer.
- Division is more efficient than multiplication followed by shifts on the target hardware.

Note The Fixed-Point Designer software is not aware of the target hardware. Before selecting this option, verify that division is more efficient than multiplication followed by shifts on your target hardware.

When Not to Use Division to Handle Net Slope Computation

This optimization does not work if:

- The software cannot perform the division using the production target `long` data type and therefore must use multiword operations.

Using multiword division does not produce code suitable for embedded targets. Therefore, do not use division to handle net slope computation in models that use multiword operations. If your

model contains blocks that use multiword operations, change the word length of these blocks to avoid these operations.

- Net slope is a power of 2 or a rational approximation of the net slope contains division by a power of 2.

Binary-point-only scaling, where the net slope is a power of 2, involves moving the binary point within the fixed-point word. This scaling mode already minimizes the number of processor arithmetic operations.

Use Division to Handle Net Slope Computation

To enable this optimization:

- 1 In the **Configuration Parameters** dialog box, on the **Math and Data Types > Data Types** pane, set **Use division for fixed-point net slope computation** to **On**, or **Use division for reciprocals of integers only**

For more information, see “Use division for fixed-point net slope computation”.

- 2 On the **Hardware Implementation > Device details** pane, set the **Signed integer division rounds to** configuration parameter to **Floor** or **Zero**, as appropriate for your target hardware. The optimization does not occur if the **Signed integer division rounds to** parameter is **Undefined**.

Note Set this parameter to a value that is appropriate for the target hardware. Failure to do so might result in division operations that comply with the definition on the **Hardware Implementation** pane, but are inappropriate for the target hardware.

- 3 Set the **Integer rounding mode** of the blocks that require net slope computation (for example, **Product**, **Gain**, and **Data Type Conversion**) to **Simplest** or match the rounding mode of your target hardware.

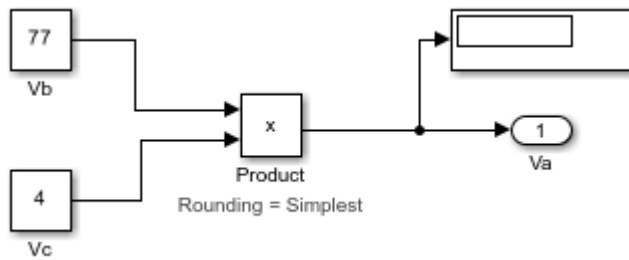
Note You can use the Model Advisor to alert you if you have not configured your model correctly for this optimization. Open the Model Advisor and run the **Identify questionable fixed-point operations** check. For more information, see “Identify blocks that will invoke net slope computation” on page 48-23 .

Improve Numerical Accuracy of Simulation Results with Rational Approximations to Handle Net Slope

This example illustrates how setting the **Math and Data Types > Use division for fixed-point net slope computation** parameter to **On** improves numerical accuracy.

Open `ex_net_slope1` Model

`ex_net_slope1`



Explore the Model

$$S_a Q_a = S_b Q_b \cdot S_c Q_c$$

or

$$Q_a = \frac{S_b S_c}{S_a} \cdot Q_b Q_c$$

where the net slope is:

$$\frac{S_b S_c}{S_a}$$

The net slope for the Product block is 7/11. Because the net slope can be represented as a fractional value consisting of small integers, you can use the **On** setting of the **Use division for fixed-point net slope computation** optimization parameter if your model and hardware configuration are suitable. For more information, see “When to Use Division for Fixed-Point Net Slope Computation” on page 48-8.

For the Product block in this model,

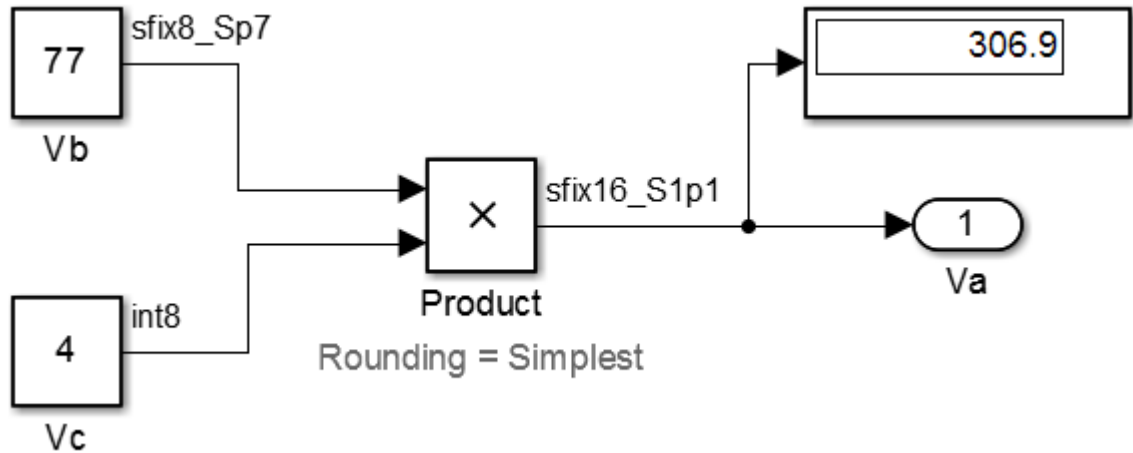
$$V_a = V_b \times V_c$$

These values are represented by the general [Slope Bias] encoding scheme described in “Scaling” on page 35-5: $V_i = S_i Q_i + B_i$.

Because there is no bias for the inputs or outputs:

Set Up Model and Run Simulation

- 1 For the Constant block Vb, set the **Output data type** to `fixdt(1, 8, 0.7, 0)`. For the Constant block Vc, set the **Output data type** to `fixdt(1, 8, 0)`.
- 2 For the Product block, set the **Output data type** to `fixdt(1, 16, 1.1, 0)`. Set the **Integer rounding mode** to **Simplest**.
- 3 In the **Configuration Parameters** dialog box, set the **Hardware Implementation > Device details > Signed integer division rounds to** configuration parameter to **Zero**.
- 4 Set the **Math and Data Types > Use division for fixed-point net slope computation** to **Off**.
- 5 In your Simulink model window, in the **Simulation** tab, click **Run**.

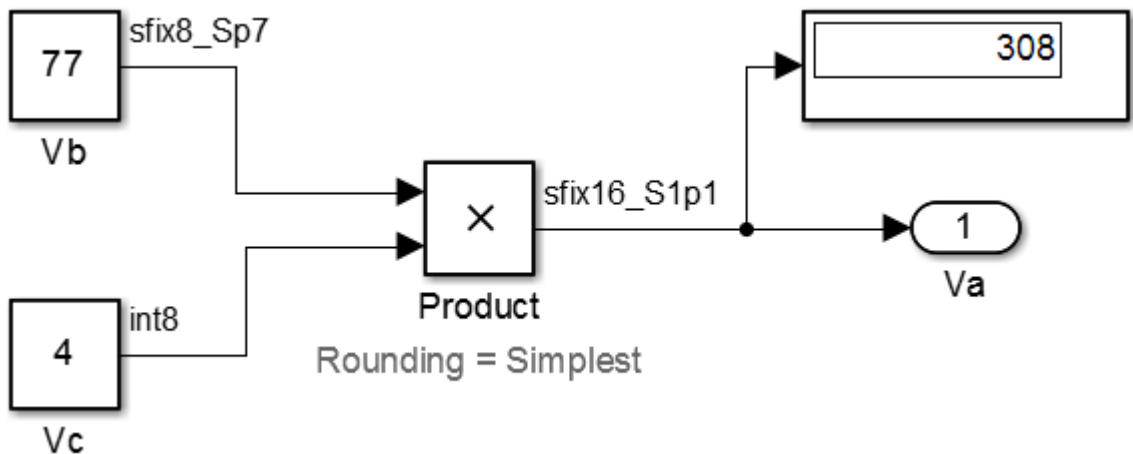


Because the simulation uses multiplication followed by shifts to handle the net slope computation, net slope precision loss occurs. This precision loss results in numerical inaccuracy: the calculated product is 306.9, not 308, as you expect.

Note You can set up the Fixed-Point Designer software to provide alerts when precision loss occurs in fixed-point constants. For more information, see “Net Slope and Net Bias Precision” on page 36-21.

6 Set the **Math and Data Types > Use division for fixed-point net slope computation** to 0n.

Save your model, and simulate again.



The software implements the net slope computation using a rational approximation instead of multiplication followed by shifts. The calculated product is 308, as you expect.

The optimization works for this model because:

- The net slope is representable as a fraction with small integers in the numerator and denominator.
- The **Hardware Implementation > Device details > Signed integer division rounds to** configuration parameter is set to Zero.

Note This setting must match your target hardware rounding mode.

- The **Integer rounding mode** of the Product block in the model is set to Simplest.
- The model does not use multiword operations.

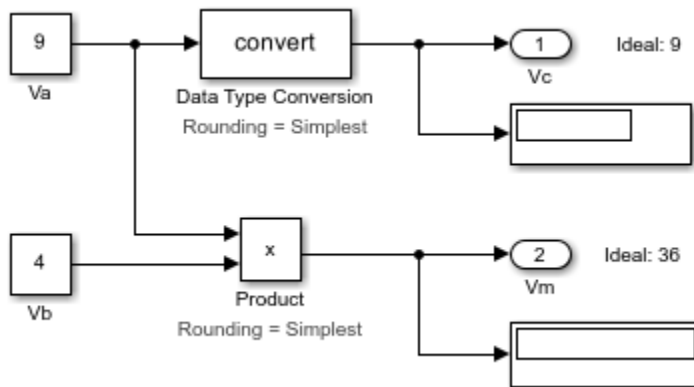
Improve Efficiency of Generated Code with Rational Approximations to Handle Net Slope

This example shows how setting the optimization parameter **Math and Data Types > Use division for fixed-point net slope computation** to On improves the efficiency of generated code.

Note The generated code is more efficient only if division is more efficient than multiplication followed by shifts on your target hardware.

Open ex_net_slope_2 Model

```
open_system("ex_net_slope2.slx")
```



Explore the Model

For the Product block in this model,

$$V_m = V_a \times V_b$$

These values are represented by the general [Slope Bias] encoding scheme described in "Scaling" on page 35-5: $V_i = S_i Q_i + B_i$.

Because there is no bias for the inputs or outputs:

$$S_m Q_m = S_a Q_a \cdot S_b Q_b$$

or

$$Q_m = \frac{S_a S_b}{S_m} \cdot Q_a Q_b$$

where the net slope is:

$$\frac{S_a S_b}{S_m}$$

The net slope for the Product block is 9/10.

Similarly, for the Data Type Conversion block in this model,

$$S_a Q_a + B_a = S_b Q_b + B_b$$

There is no bias. Therefore, the net slope is $\frac{S_b}{S_a}$. The net slope for this block is also 9/10.

Because the net slope can be represented as a fraction, you can set the **Math and Data Types > Use division for fixed-point net slope computation** optimization parameter to On if your model and hardware configuration are suitable. For more information, see “When to Use Division for Fixed-Point Net Slope Computation” on page 48-8.

Set Up the Model and Generate Code

- 1 For the Inport block Va, set the **Output data type** to `fixdt(1, 8, 9/10, 0)`; for the Inport block Vb, set the **Output data type** to `int8`.
- 2 For the Data Type Conversion block, set the **Integer rounding mode** to `Simplest`. Set the **Output data type** to `int16`.
- 3 For the Product block, set the **Integer rounding mode** to `Simplest`. Set the **Output data type** to `int16`.
- 4 Set the **Hardware Implementation > Device details > Signed integer division rounds to** configuration parameter to `Zero`.
- 5 Set the **Math and Data Types > Use division for fixed-point net slope computation** to `Off`.
- 6 From the Simulink **Apps** tab, select **Embedded Coder**. In the **C Code** tab, click **Build**.

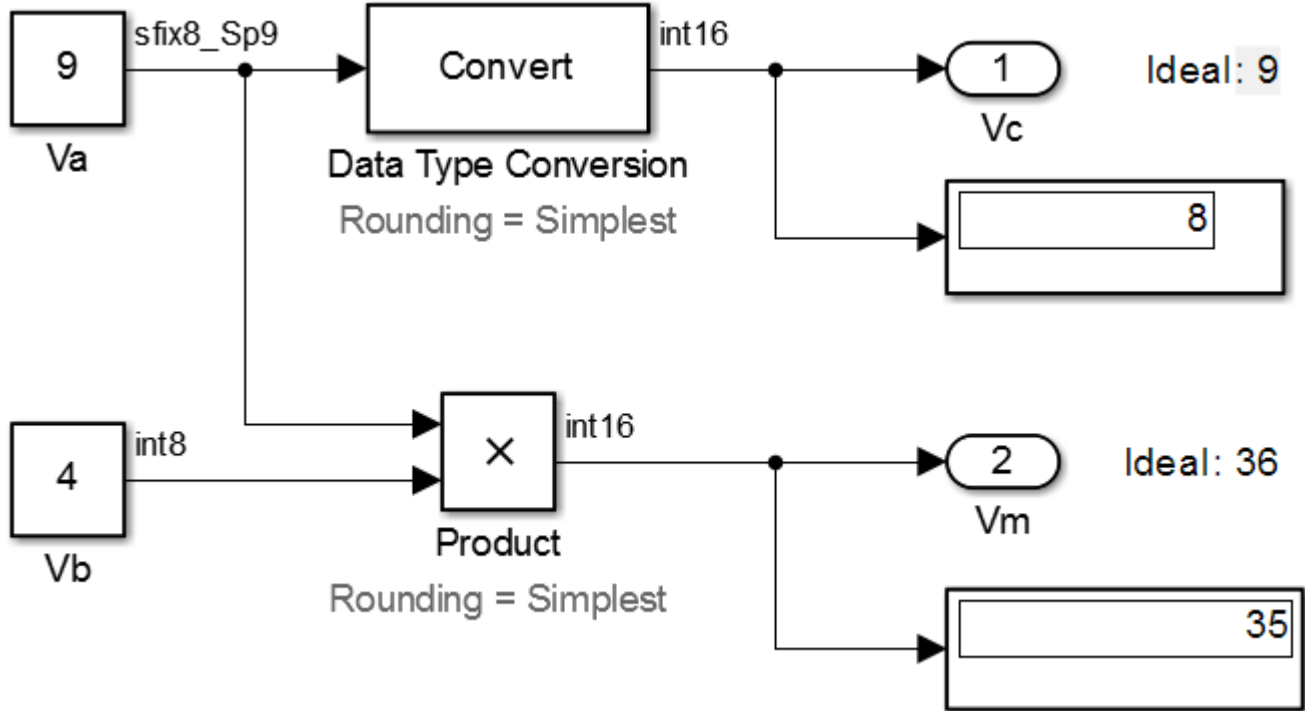
Conceptually, the net slope computation is 9/10 or 0.9:

$$\begin{aligned} V_c &= 0.9 * V_a; \\ V_m &= 0.9 * V_a * V_b; \end{aligned}$$

The generated code uses multiplication with shifts:

```
% For the conversion
Vc = (int16_T)(Va * 115 >> 7);
% For the multiplication
Vm = (int16_T)((Va * Vb >> 1) * 29491 >> 14);
```

The ideal value of the net slope computation is 0.9. In the generated code, the approximate value of the net slope computation is $29491 \gg 15 = 29491/2^{15} = 0.899993896484375$. This approximation introduces numerical inaccuracy. For example, using the same model with constant inputs produces the following results.



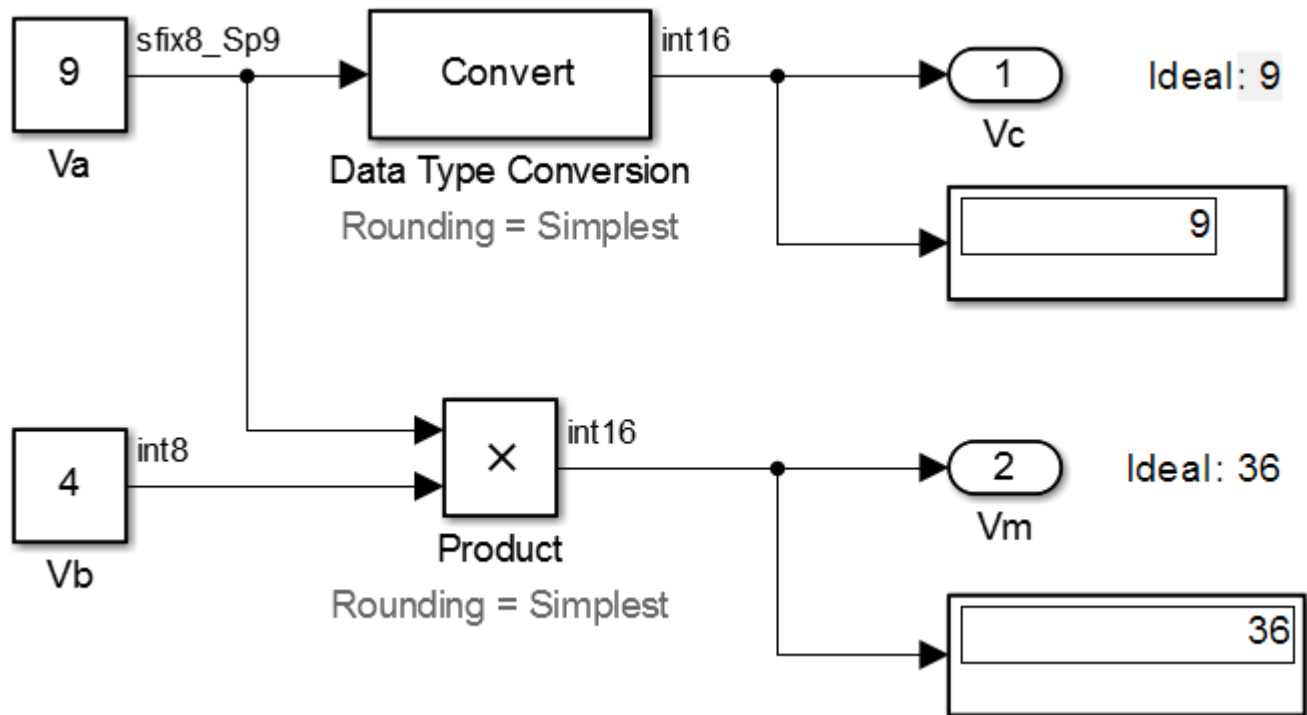
- 7 In the original model with inputs **Va** and **Vb**, set the **Math and Data Types > Use division for fixed-point net slope computation** parameter to **On**, update the diagram, and generate code again.

The generated code now uses integer division instead of multiplication followed by shifts:

```
% For the conversion
Vc = (int16_T)(Va * 9/10);
% For the multiplication
Vm = (int16_T)(Va * Vb * 9/10);
```

- 8 In the generated code, the value of the net slope computation is now the ideal value of 0.9. Using division, the results are numerically accurate.

In the model with constant inputs, set the **Math and Data Types > Use division for fixed-point net slope computation** parameter to **On** and simulate the model.



The optimization works for this model because the:

- Net slope is representable as a fraction with small integers in the numerator and denominator.
- **Hardware Implementation > Device details > Signed integer division rounds to** configuration parameter is set to Zero.

Note This setting must match your target hardware rounding mode.

- For the Product and Data Type Conversion blocks in the model, the **Integer rounding mode** is set to Simplest.
- Model does not use multiword operations.

Use Integer Division to Handle Net Slope Computation

Setting the **Math and Data Types > Use division for fixed-point net slope computation** parameter to Use division for reciprocals of integers only triggers the optimization only in cases where the net slope is the reciprocal of an integer. This setting results in a single integer division to handle net slope computations.

Control the Generation of Fixed-Point Utility Functions

In this section...

“Optimize Generated Code Using Specified Minimum and Maximum Values” on page 48-16

“Eliminate Unnecessary Utility Functions Using Specified Minimum and Maximum Values” on page 48-18

Optimize Generated Code Using Specified Minimum and Maximum Values

The Fixed-Point Designer software uses representable minimum and maximum values and constant values to determine if it is possible to optimize the generated code, for example, by eliminating unnecessary utility functions and saturation code from the generated code.

This optimization results in:

- Reduced ROM and RAM consumption
- Improved execution speed

When you select the **Optimize using specified minimum and maximum values** configuration parameter, the software takes into account input range information, also known as design minimum and maximum, that you specify for signals and parameters in your model. It uses these minimum and maximum values to derive range information for downstream signals in the model and then uses this derived range information to simplify mathematical operations in the generated code whenever possible.

Prerequisites

The **Optimize using specified minimum and maximum values** parameter appears for ERT-based targets only and requires an Embedded Coder license when generating code.

How to Configure Your Model

To make optimization more likely:

- Provide as much design minimum and maximum information as possible. Specify minimum and maximum values for signals and parameters in the model for:
 - Inport and Outport blocks
 - Block outputs
 - Block inputs, for example, for the MATLAB Function and Stateflow Chart blocks
 - `Simulink.Signal` objects
- Before generating code, test the minimum and maximum values for signals and parameters. Otherwise, optimization might result in numerical mismatch with simulation. You can simulate your model with simulation range checking enabled. If errors or warnings occur, fix these issues before generating code.

How to Enable Simulation Range Checking

- 1 In the **Modeling** tab of the Simulink editor, click **Model Settings** to open the Configuration Parameters dialog box.

- 2 In the Configuration Parameters dialog box, select **Diagnostics > Data Validity**.
 - 3 On the **Data Validity** pane, under **Signals**, set **Simulation range checking** to warning or error.
- Use fixed-point data types with binary-point-only (power-of-two) scaling.
 - Provide design minimum and maximum information upstream of blocks as close to the inputs of the blocks as possible. If you specify minimum and maximum values for a block output, these values are most likely to affect the outputs of the blocks immediately downstream. For more information, see “Eliminate Unnecessary Utility Functions Using Specified Minimum and Maximum Values” on page 48-18.

How to Enable Optimization

- 1 In the Configuration Parameters dialog box, set the **Code Generation > System target file** to select an Embedded Real-Time (ERT) target (requires an Embedded Coder license).
- 2 Specify design minimum and maximum values for signals and parameters in your model using the tips in “How to Configure Your Model” on page 48-16.
- 3 Select the **Optimization > Advanced parameters > Optimize using the specified minimum and maximum values** configuration parameter.

Limitations

- This optimization does not occur for:
 - Multiword operations
 - Fixed-point data types with slope and bias scaling
 - Addition unless the fraction length is zero
- This optimization does not take into account minimum and maximum values for:
 - Merge block inputs. To work around this issue, use a `Simulink.Signal` object on the Merge block output and specify the range on this object.
 - Bus elements.
 - Conditionally executed subsystem (such as a triggered subsystem) block outputs that are directly connected to an Output block.

Output blocks in conditionally executed subsystems can have an initial value specified for use only when the system is not triggered. In this case, the optimization cannot use the range of the block output because the range might not cover the initial value of the block.

- There are limitations on precision because you specify the minimum and maximum values as double-precision values. If the true value of a minimum or maximum value cannot be represented as a double, ensure that you round the minimum and maximum values correctly so that they cover the true design range.
- If your model contains multiple instances of a reusable subsystem and each instance uses input signals with different specified minimum and maximum values, this optimization might result in different generated code for each subsystem so code reuse does not occur. Without this optimization, the Simulink Coder software generates code once for the subsystem and shares this code among the multiple instances of the subsystem.

Eliminate Unnecessary Utility Functions Using Specified Minimum and Maximum Values

This example shows how the Fixed-Point Designer software uses the input range for a division operation to determine whether it can eliminate unnecessary utility functions from the generated code. It uses the `fxpdemo_min_max_optimization` model. First, you generate code without using the specified minimum and maximum values to see that the generated code contains utility functions to ensure that division by zero does not occur. You then turn on the optimization, and generate code again. With the optimization, the generated code does not contain the utility function because it is not necessary for the input range.

Generate Code Without Using Minimum and Maximum Values

First, generate code without taking into account the design minimum and maximum values for the first input of the division operation to show the code without the optimization. In this case, the software uses the representable ranges for the two inputs, which are both `uint16`. With these input ranges, it is not possible to implement the division with the specified precision using shifts, so the generated code includes a division utility function.

- 1 Run the example. At the MATLAB command line, enter:

```
fxpdemo_min_max_optimization
```

- 2 In the example window, double-click the **View Optimization Configuration** button.

The Optimization pane of the Configuration Parameters dialog box appears.

Note that the **Optimize using specified minimum and maximum values** parameter is not selected.

- 3 Double-click the **Generate Code** button.

The code generation report appears.

- 4 In the model, right-click the Division with increased fraction length output type block.

The context menu appears.

- 5 From the context menu, select **C/C++ Code > Navigate To C/C++ Code**.

The code generation report highlights the code generated for this block. The generated code includes a call to the `div_repeat_u32` utility function.

```
rtY.Out3 = div_repeat_u32((uint32_T)rtU.In5 << 16,  
    (uint32_T)rtU.In6, 1U);
```

- 6 Click the `div_repeat_u32` link to view the utility function, which contains code for handling division by zero.

Generate Code Using Minimum and Maximum Values

Next, generate code for the same division operation, this time taking into account the design minimum and maximum values for the first input of the Product block. These minimum and maximum values are specified on the Inport block directly upstream of the Product block. With these input ranges, the generated code implements the division by simply using a shift. It does not need to generate a division utility function, reducing both memory usage and execution time.

- 1 Double-click the Inport block labeled 5 to open the block parameters dialog box.
- 2 On the block parameters dialog box, select the **Signal Attributes** pane and note that:
 - The **Minimum** value for this signal is 1.
 - The **Maximum** value for this signal is 100.
- 3 Click **OK** to close the dialog box.
- 4 Double-click the **View Optimization Configuration** button.

The Optimization pane of the Configuration Parameters dialog box appears.

- 5 On this pane, select the **Optimize using specified minimum and maximum values** parameter and click **Apply**.
- 6 Double-click the **Generate Code** button.

The code generation report appears.

- 7 In the model, right-click the Division with increased fraction length output type block.

The context menu appears.

- 8 From the context menu, select **C/C++ Code > Navigate To C/C++ Code**.

The code generation report highlights the code generated for this block. This time, the generated code implements the division with a shift operation and there is no division utility function.

```
tmp = rtU.In6;
rtY.Out3 = (uint32_T)tmp ==
  (uint32_T)0 ? MAX_uint32_T : ((uint32_T)rtU.In5 << 17) /
  (uint32_T)tmp;
```

Modify the Specified Minimum and Maximum Values

Finally, modify the minimum and maximum values for the first input to the division operation so that its input range is too large to guarantee that the value does not overflow when shifted. Here, you cannot shift a 16-bit number 17 bits to the right without overflowing the 32-bit container. Generate code for the division operation, again taking into account the minimum and maximum values. With these input ranges, the generated code includes a division utility function to ensure that no overflow occurs.

- 1 Double-click the Inport block labelled 5 to open the block parameters dialog box.
- 2 On the block parameters dialog box, select the **Signal Attributes** pane and set the **Maximum** value to 40000, then click **OK** to close the dialog box.
- 3 Double-click the **Generate Code** button.

The code generation report appears.

- 4 In the model, right-click the Division with increased fraction length output type block.

The context menu appears.

- 5 From the context menu, select **C/C++ Code > Navigate To C/C++ Code**.

The code generation report highlights the code generated for this block. The generated code includes a call to the `div_repeat_32` utility function.

```
rtY.Out3 = div_repeat_u32((uint32_T)rtU.In5 << 16,  
    (uint32_T)rtU.In6, 1U);
```

Optimize Generated Code with the Model Advisor

In this section...

“Identify Blocks that Generate Expensive Fixed-Point and Saturation Code” on page 48-21


“Identify Questionable Fixed-Point Operations” on page 48-23

“Identify Blocks that Generate Expensive Rounding Code” on page 48-25

You can use the Simulink Model Advisor to help you configure your fixed-point models to achieve a more efficient design and optimize your generated code. To use the Model Advisor to check your fixed-point models:

- 1 In the **Modeling** tab of the model you want to analyze, click **Model Advisor**.
- 2 In the System Selector, select the system to analyze.
- 3 In the Model Advisor left pane, expand the **By Product** node and then the **Embedded Coder** node.
- 4 For fixed-point code generation, the most important check boxes to select are **Identify blocks that generate expensive fixed-point and saturation code**, **Identify questionable fixed-point operations**, **Identify blocks that generate expensive rounding code**, and **Check the hardware implementation**.

To enable all Model Advisor checks associated with the selected node, select the check box of the folder containing the checks. In our case, select the **Embedded Coder** check box.

- 5 Click **Run checks** . Any tips for improving the efficiency of your fixed-point model appear in the Model Advisor window.

The sections that follow discuss fixed-point related checks and sub-checks found in the Model Advisor. These sections explain the checks, discuss their importance in fixed-point code generation, and offer suggestions for tweaking your model to optimize your generated code.

Identify Blocks that Generate Expensive Fixed-Point and Saturation Code

Identify Sum blocks for questionable fixed-point operations

- When the input range of a Sum block exceeds the output range, a range error occurs. You can get any addition or subtraction your application requires by inserting data type conversion blocks before and/or after the Sum block.
- When a Sum block has an input with a slope adjustment factor that does not equal the slope adjustment factor of the output, the mismatch requires the Sum block to perform a multiply operation each time the input is converted to the data type and scaling of the output. The mismatch can be removed by changing the scaling of the output or the input.
- When the net sum of the Sum block input biases does not equal the bias of the output, the generated code includes one extra addition or subtraction instruction to correctly account for the net bias adjustment. Changing the bias of the output scaling can make the net bias adjustment zero and eliminate the need for the extra operation.

Identify Min Max blocks for questionable fixed-point operations

- When the input and output of the MinMax block have different data types, a conversion operation is required every time the block is executed. The model is more efficient with the same data types.
- When the data type and scaling of the input of the MinMax block does not match the data type and scaling of the output, a conversion is required before performing a relational operation. This could result in a range error when casting, or a precision loss each time a conversion is performed. Change the scaling of either the input or output to generate more efficient code.
- When the input of the MinMax block has a different slope adjustment factor than the output, the MinMax block requires a multiply operation each time the block is executed to convert the input to the data type and scaling of the output. You can correct the mismatch by changing the scaling of either the input or output.

Identify Discrete Integrator blocks for questionable fixed-point operations

- When the initial condition for the Discrete-Time Integrator blocks is used to initialize the state and output, the output equation generates excessive code and an extra global variable is required. It is recommended that you set the **Function Block Parameters > Initial condition setting** parameter to State (most efficient).

Identify Compare to Constant blocks for questionable fixed-point operations

- If the input data type of the Compare to Zero block cannot represent zero exactly, the input signal is compared to the closest representable value of zero, resulting in parameter overflow. To avoid this parameter overflow, select an input data type that can represent zero.
- If the **Constant value** parameter of the Compare to Constant is outside the range that the input data type can represent, the input signal is compared to the closest representable value of the constant. This results in parameter overflow. To avoid this parameter overflow, select an input data type that can represent the **Constant value**, or change the **Constant value** to a value that can be accommodated by the input data type.

Identify Lookup Table blocks for questionable fixed-point operations

Efficiency trade-offs related to lookup table data are described in “Effects of Spacing on Speed, Error, and Memory Usage” on page 41-59. Based on these trade-offs, the Model Advisor identifies blocks where there is potential for efficiency improvements, such as:

- Lookup table input data is not evenly spaced.
- Lookup table input data is *not* evenly spaced when quantized, but it is very close to being evenly spaced.
- Lookup table input data is evenly spaced, but the spacing is not a power of two.

For more information on lookup table optimization, see “Lookup Table Optimization” on page 48-27.

Check optimization and hardware implementation settings

- Integer division generated code contains protection against arithmetic exceptions such as division by zero, INT_MIN/-1, and LONG_MIN/-1. If you construct models making it impossible for exception triggering input combinations to reach a division operation, the protection code generated as part of the division operation is redundant.
- The index search method `Evenly-spaced points` requires a division operation, which can be computationally expensive.

Identify blocks that will invoke net slope computation

When a change of fixed-point slope is not a power of two, net slope computation is necessary. Normally, net slope computation is implemented using an integer multiplication followed by shifts. Under some conditions, an alternate implementation requires just an integer division by a constant. One of the conditions is that the net slope can be very accurately represented as the reciprocal of an integer. When this condition is met, the division implementation produces more accurate numerical behavior. Depending on your compiler and embedded hardware, the division implementation might be more desirable than the multiplication and shifts implementation. The generated code might be more efficient in either ROM size or model execution size.

The Model Advisor alerts you when:

- You set the **Use division for fixed-point net slope computation** optimization parameter to 'On', but your model configuration is not compatible with this selection.
- Your model configuration is suitable for using division to handle net slope computation, but you do not set the **Use division for fixed-point net slope computation** optimization parameter to 'On'.

For more information, see “Net Slope Computation” on page 48-8.

Identify product blocks that are less efficient

The number of multiplications and divisions that a block performs can have a significant impact on accuracy and efficiency. The Model Advisor detects some, but not all, situations where rearranging the operations can improve accuracy, efficiency, or both.

One such situation is when a calculation using more than one division operation is computed. A general guideline from the field of numerical analysis is to multiply all the denominator terms together first, then do one and only one division. This improves accuracy and often speed in floating-point and especially fixed-point. This can be accomplished in Simulink by cascading Product blocks. Note that multiple divisions spread over a series of blocks are not detected by the Model Advisor.

Another situation is when a single Product block is configured to do more than one multiplication or division operation. This is supported, but if the output data type is integer or fixed-point, then better results are likely if this operation is split across several blocks each doing one multiplication or one division. Using several blocks allows the user to control the data type and scaling used for intermediate calculations. The choice of data types for intermediate calculations affects precision, range errors, and efficiency.

Check for expensive saturation code

Setting the **Saturate on integer overflow** parameter can produce condition checking code that your application might not require.

Check whether your application requires setting **Block Parameters > Signal Attributes > Saturate on integer overflow**. Otherwise, clear this parameter for the most efficient implementation of the block in the generated code.

Identify Questionable Fixed-Point Operations

This check identifies blocks that generate multiword operations, cumbersome multiplication and division operations, expensive conversion code, inefficiencies in lookup table blocks, and expensive comparison code.

Check for multiword operations

When an operation results in a data type larger than the largest word size of your processor, the generated code contains multiword operations. Multiword operations can be inefficient on hardware. To prevent multiword operations, adjust the word lengths of inputs to operations so that they do not exceed the largest word size of your processor. For more information on controlling multiword operations in generated code, see “Fixed-Point Multiword Operations In Generated Code” on page 54-166.

Check for expensive multiplication code

- “Targeting an Embedded Processor” on page 37-3 discusses the capabilities and limitations of embedded processors. “Design Rules” on page 37-4 recommends that inputs to a multiply operation should not have word lengths larger than the base integer type of your processor. Multiplication with larger word lengths can always be handled in software, but that approach requires much more code and is much slower. The Model Advisor identifies blocks where undesirable software multiplications are required. Visual inspection of the generated code, including the generated multiplication utility function, will make the cost of these operations clear. It is strongly recommended that you adjust the model to avoid these operations.
- “Rules for Arithmetic Operations” on page 36-42 discusses the implementation details of fixed-point multiplication and division. Significant increase in complexity occurs when signals with nonzero biases are involved in multiplication and division. It is strongly recommended that you make changes to eliminate the need for these complicated operations. Extra steps are required to implement the multiplication. Inserting a Data Type Conversion block before and after the block doing the multiplication allows the biases to be removed and allows the user to control data type and scaling for intermediate calculations. In many cases the Data Type Conversion blocks can be moved to the “edges” of a (sub)system. The conversion is only done once and all blocks can benefit from simpler bias-free math.

Check for expensive division code

The rounding behavior of signed integer division is not fully specified by C language standards. Therefore, the generated code for division is too large to provide bit-true agreement between simulation and code generation. To avoid integer division generated code that is too large, in the Configuration Parameters dialog box, on the **Hardware Implementation** pane, set the **Signed integer division rounds to** parameter to the recommended value.

Identify lookup blocks with uneven breakpoint spacing

Efficiency trade-offs related to lookup table data are described in “Effects of Spacing on Speed, Error, and Memory Usage” on page 41-59. Based on these trade-offs, the Model Advisor identifies blocks where there is potential for efficiency improvements, and issues a warning when:

- Lookup table input data is not evenly spaced.
- Lookup table input data is *not* evenly spaced when quantized, but it is very close to being evenly spaced.
- Lookup table input data is evenly spaced, but the spacing is not a power of two.

Check for expensive pre-lookup division

For a Prelookup or n-D Lookup Table block, **Index search method** is Evenly spaced points. Breakpoint data does not have power of 2 spacing.

If breakpoint data is nontunable, it is recommended that you adjust the data to have even, power of 2 spacing. Otherwise, in the block parameter dialog box, specify a different **Index search method** to avoid the computation-intensive division operation.

Check for expensive data type conversions

When a block is configured such that it would generate inefficient code for a data type conversion, the Model Advisor generates a warning and makes suggestions on how to make your model more efficient.

Check for fixed-point comparisons with predetermined results

When you select `isInf`, `isNaN`, or `isFinite` as the operation for the Relational Operator block, the block switches to one-input mode. In this mode, if the input data type is fixed point, boolean, or a built-in integer, the output is `FALSE` for `isInf` and `isNaN`, `TRUE` for `isFinite`. This might result in dead code which will be eliminated by Simulink Coder.

Check for expensive binary comparison operations

- When the input data types of a Relational Operator block are not the same, a conversion operation is required every time the block is executed. If one of the inputs is invariant, then changing the data type and scaling of the invariant input to match the other input improves the efficiency of the model.
- When the inputs of a Relational Operator block have different ranges, there will be a range error when casting, and a precision loss each time a conversion is performed. You can insert Data Type Conversion blocks before the Relational Operator block to convert both inputs to a common data type that has enough range and precision to represent each input.
- When the inputs of a Relational Operator block have different slope adjustment factors, the Relational Operator block is required to perform a multiply operation each time the input with lesser positive range is converted to the data type and scaling of the input with greater positive range. The extra multiplication requires extra code, slows down the speed of execution, and usually introduces additional precision loss. By adjusting the scaling of the inputs, you can eliminate mismatched slopes.

Check for expensive comparison code

When your model is configured such that the generated code contains expensive comparison code, the Model Advisor generates a warning.

Check for expensive fixed-point data types in generated code

When a design contains integer or fixed-point word lengths that do not exist on your target hardware, the generated code can contain extra saturation code, shifts, and multiword operations. By changing the data type to one that is supported by your target hardware, you can improve the efficiency of the generated code. The Model Advisor flags these expensive data types in your model. For example, the Model Advisor would flag a fixed-point data type with a word length of 17 if the target hardware was 32 bits.

Identify Blocks that Generate Expensive Rounding Code

This check alerts you when rounding optimizations are available. To check for blocks that generate expensive rounding code, the Model Advisor performs the following sub-checks:

- Check for expensive rounding operations in multiplication and division
- Check optimization and Hardware Implementation settings (Lookup Blocks)
- Check for expensive rounding in a data type conversion
- Check for expensive rounding modes in the model

Traditional handwritten code, especially for control applications, almost always uses “no effort” rounding. For example, for unsigned integers and two's complement signed integers, shifting right and dropping the bits is equivalent to rounding to floor. To get results comparable to, or better than, what you expect from traditional handwritten code, use the simplest rounding mode. In general the simplest mode provides the minimum cost solution with no overflows. If the simplest mode is not available, round to floor.

The primary exception to this rule is the rounding behavior of signed integer division. The C standard leaves this rounding behavior unspecified, but for most production targets the “no effort” mode is to round to zero. For unsigned division, everything is non-negative, so rounding to floor and rounding to zero are identical. To improve rounding efficiency, set **Model Configuration Parameters > Hardware Implementation > Device details > Signed integer division rounds to** using the mode that your production target uses.

Use the **Integer rounding mode** parameter on your model's blocks to simulate the rounding behavior of the C compiler that you use to compile code generated from the model. This setting appears on the **Signal Attributes** pane of the parameter dialog boxes of blocks that can perform signed integer arithmetic, such as the Product block. To obtain the most efficient generated code, change the **Integer rounding mode** parameter of the block to the recommended setting.

For more information on properties to consider when choosing a rounding mode, see “Choosing a Rounding Method” on page 1-7.

Lookup Table Optimization

A function lookup table is a method by which you can approximate a function using a table with a finite number of points (X, Y). The X values of the lookup table are called the breakpoints. You approximate the value of the ideal function at a point by interpolating between the two breakpoints closest to the point. Because table lookups and simple estimations can be faster than mathematical function evaluations, using lookup table blocks often result in speed gains when simulating a model.

To optimize lookup tables in your model:

- Limit uneven lookup tables.

Unevenly spaced breakpoints require a general-purpose algorithm such as a binary search to determine where the input lies in relation to the breakpoints. This additional computation increases ROM and execution time.

- Prevent evenly spaced lookup tables from being treated as unevenly spaced.

The position search in evenly spaced lookup tables is much faster. In addition, the interpolation requires a simple division.

Sometimes, when a lookup table is converted to fixed-point, a quantization error results. A lookup table that is evenly spaced in floating-point, could be unevenly spaced in the generated fixed-point code. Use the `fixpt_evenspace_cleanup` function to convert the data into an evenly spaced lookup table again.

- Use power of two spaced breakpoints in lookup tables.

In power of two spaced lookup tables, a bit shift replaces the position search, and a bit mask replaces the interpolation making this construct the most efficient regardless of your target language and hardware.

The following table summarizes the effects of lookup table breakpoint spacing.

| Parameter | Even Power of 2 Spaced Data | Evenly Spaced Data | Unevenly Spaced Data |
|-----------------|---|--|---|
| Execution speed | The execution speed is the fastest. The position search and interpolation are the same as for evenly spaced data. However, to increase the speed more, a bit shift replaces the position search, and a bit mask replaces the interpolation. | The execution speed is faster than that for unevenly spaced data, because the position search is faster and the interpolation requires a simple division. | The execution speed is the slowest of the different spacings because the position search is slower, and the interpolation requires more operations. |
| Error | The error can be larger than that for unevenly spaced data because approximating a function with nonuniform curvature requires more points to achieve the same accuracy. | The error can be larger than that for unevenly spaced data because approximating a function with nonuniform curvature requires more points to achieve the same accuracy. | The error can be smaller because approximating a function with nonuniform curvature requires fewer points to achieve the same accuracy. |

| Parameter | Even Power of 2 Spaced Data | Evenly Spaced Data | Unevenly Spaced Data |
|-----------|---|---|---|
| ROM usage | Uses less command ROM, but more data ROM. | Uses less command ROM, but more data ROM. | Uses more command ROM, but less data ROM. |
| RAM usage | Not significant. | Not significant. | Not significant. |

Use the Model Advisor “Identify Questionable Fixed-Point Operations” on page 48-23 check to identify lookup table blocks where there is potential for efficiency improvements.

See Also

More About

- “Effects of Spacing on Speed, Error, and Memory Usage” on page 41-59
- “Create Lookup Tables for a Sine Function” on page 41-5

Selecting Data Types for Basic Operations

In this section...

“Restrict Data Type Word Lengths” on page 48-29

“Avoid Fixed-Point Scalings with Bias” on page 48-29

“Wrap and Round to Floor or Simplest” on page 48-29

“Limit the Use of Custom Storage Classes” on page 48-30

Restrict Data Type Word Lengths

If possible, restrict the fixed-point data type word lengths in your model so that they are equal to or less than the integer size of your target microcontroller. This results in fewer mathematical instructions in the microcontroller, and reduces ROM and execution time.

This recommendation strongly applies to global variables that consume global RAM. For example, Unit Delay blocks have discrete states that have the same word lengths as their input and output signals. These discrete states are global variables that consume global RAM, which is a scarce resource on many embedded systems.

For temporary variables that only occupy a CPU register or stack location briefly, the space consumed by a `long` is less critical. However, depending on the operation, the use of `long` variables in math operations can be expensive. Addition and subtraction of long integers generally requires the same effort as adding and subtracting regular integers, so that operation is not a concern. In contrast, multiplication and division with long integers can require significantly larger and slower code.

Avoid Fixed-Point Scalings with Bias

Whenever possible, avoid using fixed-point numbers with bias. In certain cases, if you choose biases carefully, you can avoid significant increases in ROM and execution time. Refer to “Recommendations for Arithmetic and Scaling” on page 36-31 for more information on how to choose appropriate biases in cases where it is necessary; for example if you are interfacing with a hardware device that has a built-in bias. In general, however, it is safer to avoid using fixed-point numbers with bias altogether.

Inputs to lookup tables are an important exception to this recommendation. If a lookup table input and the associated input data use the same bias, then there is no penalty associated with nonzero bias for that operation.

Wrap and Round to Floor or Simplest

For most fixed-point and integer operations, the Simulink software provides you with options on how overflows are handled and how calculations are rounded. Traditional handwritten code, especially for control applications, almost always uses the “no effort” rounding mode. For example, to reduce the precision of a variable, that variable is shifted right. For unsigned integers and two's complement signed integers, shifting right is equivalent to rounding to floor. To get results comparable to or better than what you expect from traditional handwritten code, you should round to floor in most cases.

The primary exception to this rule is the rounding behavior of signed integer division. The C language leaves this rounding behavior unspecified, but for most targets the “no effort” mode is round to zero. For unsigned division, everything is non-negative, so rounding to floor and rounding to zero are identical.

You can improve code efficiency by setting the value of the **Model Configuration Parameters > Hardware Implementation > Device details > Signed integer division rounds to** parameter to describe how your production target handles rounding for signed division. For Product blocks that are doing only division, setting the **Integer rounding mode** parameter to the rounding mode of your production target gives the best results. You can also use the **Simplest** rounding mode on blocks where it is available. For more information, refer to “Rounding Mode: Simplest” on page 36-14.

The options for overflow handling also have a big impact on the efficiency of your generated code. Using software to detect overflow situations and saturate the results requires the code to be much bigger and slower compared to simply ignoring the overflows. When overflows are ignored for unsigned integers and two's complement signed integers, the results usually wrap around modulo 2^N , where N is the number of bits. Unhandled overflows that wrap around are highly undesirable for many situations.

However, because of code size and speed needs, traditional handwritten code contains very little software saturation. Typically, the fixed-point scaling is very carefully set so that overflow does not occur in most calculations. The code for these calculations safely ignores overflow. To get results comparable to or better than what you would expect from traditional handwritten code, the **Saturate on integer overflow** parameter should not be selected for Simulink blocks doing those calculations.

In a design, there might be a few places where overflow can occur and saturation protection is needed. Traditional handwritten code includes software saturation for these few places where it is needed. To get comparable generated code, the **Saturate on integer overflow** parameter should only be selected for the few Simulink blocks that correspond to these at-risk calculations.

A secondary benefit of using the most efficient options for overflow handling and rounding is that calculations often reduce from multiple statements requiring several lines of C code to small expressions that can be folded into downstream calculations. Expression folding is a code optimization technique that produces benefits such as minimizing the need to store intermediate computations in temporary buffers or variables. This can reduce stack size and make it more likely that calculations can be efficiently handled using only CPU registers. An automatic code generator can carefully apply expression folding across parts of a model and often see optimizations that might not be obvious. Automatic optimizations of this type often allow generated code to exceed the efficiency of typical examples of handwritten code.

Limit the Use of Custom Storage Classes

In addition to the tip mentioned in “Wrap and Round to Floor or Simplest” on page 48-29, to obtain the maximum benefits of expression folding you also need to make sure that the **Storage class** is set to **Auto** for signals in your model. When you choose a setting other than **Auto**, you need to name the signal, and a separate statement is created in the generated code. Therefore, only use a setting other than **Auto** when it is necessary for global variables.

For more information about setting the **Storage class** of signals, see “Configure Signal Data for C Code Generation”.

Use of Shifts by C Code Generation Products

Introduction to Shifts by Code Generation Products

MATLAB Coder, Simulink Coder, and Embedded Coder generate C code that uses the C language's shift left << and shift right >> operators. Modern C compilers provide consistent behavior for shift operators. However, some behaviors of the shift operators are not fully defined by some C standards. When you work with The MathWorks code generation products, you need to know how to manage the use of C shifts.

Two's Complement

Two's complement is a way to interpret a binary number. Most modern processors represent integers using two's complement. MathWorks code generation products require C and C++ compilers to represent signed integers using two's complement. MathWorks toolboxes and documentation use two's complement representation exclusively.

Arithmetic and Logical Shifts

The primary difference between an arithmetic shift and a logical shift is intent. Arithmetic shifts have a mathematical meaning. The intent of logical shifts is to move bits around, making them useful only for unsigned integers being used as collections of bit flags.

The C language does not distinguish between arithmetic and logical shifts and provides only one set of shift operators. When MathWorks code generation products use shifts on signed integers in generated code, the intent is always an arithmetic shift. For unsigned integers, there is no detectable difference in behavior between logical and arithmetic shifts.

Arithmetic Left-Shifts

An arithmetic left-shift represents multiplication by a power of 2.

$$a \ll b = a * 2^b$$

If the value produced by multiplying by 2^b is too big, then an overflow occurs. In case of an overflow, the ideal answer wraps around modulo 2^n to fit in the data type. The C90 standard specifies left-shift behavior. At the bit level, b of the bits are shifted off the left end and discarded. At the right end, b bits of value 0 are shifted in. The standard does not specify a difference between unsigned and signed. For both unsigned and two's complement signed, the bit level behavior provides the intended arithmetic left-shift behavior.

The C99 standard describes the arithmetic interpretation. It also states that for signed types, the behavior is undefined for any negative value or for a positive value that would overflow. A compiler vendor might exploit the C99 standard undefined behavior clause to optimize code in a way that changes the behavior intended by the coder. If your compiler is C99-compliant but not C90-compliant, then turn off the option Replace multiplications by powers of two with signed bitwise shifts (Embedded Coder). Older C++ standards follow the C90 standard with regard to shift left. Newer C++ standards are similar to the C99 standard.

Arithmetic Right-Shifts

An arithmetic right-shift represents division by a power of 2, where the ideal quotient rounds to floor.

$$a \gg b = a / 2^b$$

When a is nonnegative, the C standards state that right-shift must provide this arithmetic behavior. If a is signed and negative, then the standard states that the implementation defines the behavior. The C standard requires that compilers document their implementation behavior. Nearly all compilers implement signed shift right as an arithmetic shift that rounds to floor. This is the simplest and most efficient behavior for the compiler vendor to provide. If you have a compiler that does not provide arithmetic right-shift, or your coding standards do not allow you to use signed right-shift, then you can select options that avoid signed shift right. For example, Allow right shifts on signed integers (Embedded Coder) replaces signed right shifts with a function call.

Out-of-Range Shifts

In C, when shifting an integer of word length n , use a shift amount between 0 and $n - 1$, inclusive. The C standard does not define shifting by other amounts, such as:

- Shifting by a negative constant.
- Shifting by an amount greater than word length.

When the shift amount is constant, the products do not generate out-of-range shifts. The risk of out-of-range shifts comes from explicitly modeled shifts where the shift amount is a non constant variable. When modeling shifts with variable shift amounts, make sure that the shift amount is always in range.

Modeling Sources of Shifts

There are explicit and implicit sources of shifts in models and algorithms.

Explicit

- MATLAB bit-shift functions: `bitsll`, `bitsra`, `bitsrl`, `bitshift`
- Simulink Shift Arithmetic block
- Stateflow bitwise operations

Implicit

- Fixed-point operations that involve a scaling change

When converting fixed-point scaling, if the net slope change is not an exact power of two, then a multiplication followed by a shift approximates the ideal net slope. For more information on net slope computation, see “Handle Net Slope Computation” on page 48-8.

- Underlying higher-level algorithms (for example, FFT algorithms)

Controlling Shifts in Generated Code

Several configuration parameters have an effect on the number and style of shifts that appear in generated code.

- “Signed integer division rounds to”
Set this parameter to `Floor` or `Zero` to avoid extra generated code.
- “Use division for fixed-point net slope computation”

When enabled, this parameter uses division in place of multiplication followed by shifts to perform fixed-point net slope computation.

- Replace multiplications by powers of two with signed bitwise shifts (Embedded Coder)

When this parameter is enabled, multiplications by powers of two are replaced with signed bitwise shifts. Clearing this option supports MISRA C compliance.

- Allow right shifts on signed integers (Embedded Coder)

When this parameter is enabled, generated code can contain right bitwise shifts on signed integers. To prevent right bitwise shifts on signed integers, clear this option.

- “Shift right on a signed integer as arithmetic shift”

Select this parameter if the C compiler implements a signed integer right shift as an arithmetic right shift.

Use Hardware-Efficient Algorithm to Solve Systems of Complex-Valued Linear Equations

This example shows how to solve systems of complex-valued linear equations using hardware-efficient MATLAB® code in Simulink® using a systolic array.

Overview

This example shows a hardware-efficient method of solving the system of simultaneous equations

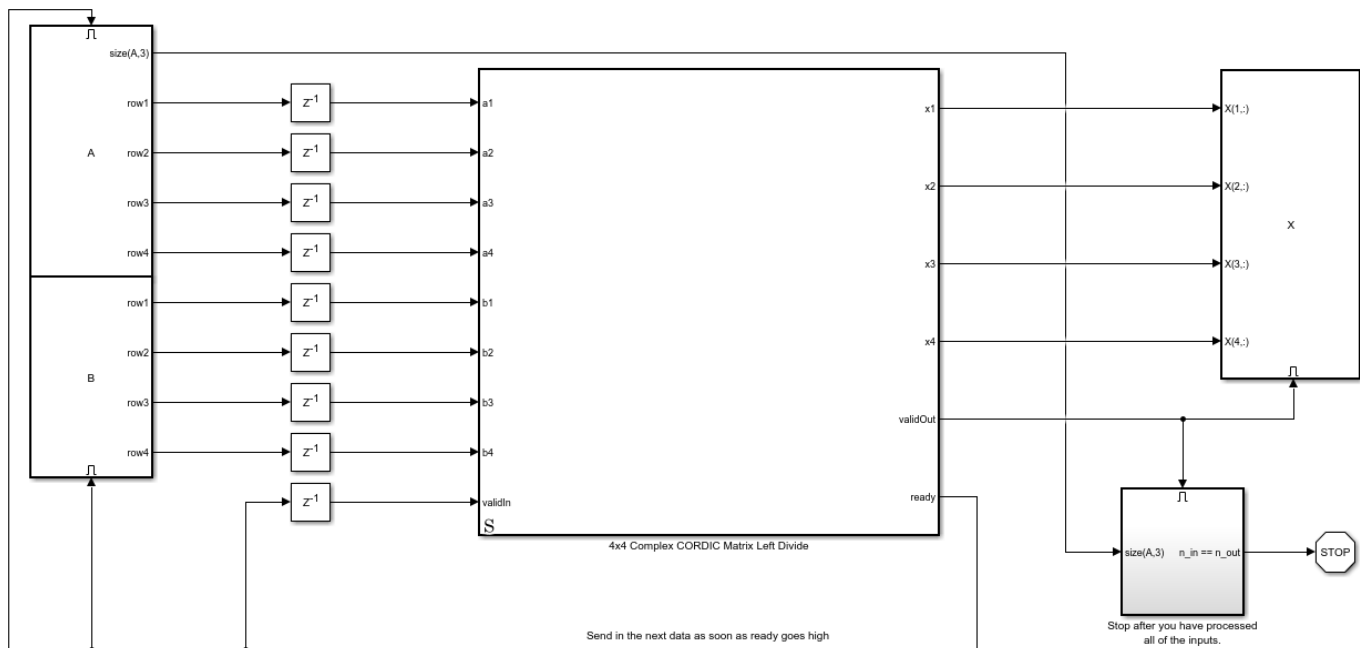
$$AX = B$$

where A is an m-by-n complex matrix, X is an n-by-p complex matrix, and B is an m-by-p complex matrix.

The Simulink model used in this example is:

```
model = 'fi_complex_mldivide_systolic_array_model';
open(model)
```

Solve Systems of Complex-Valued Linear Equations using a Systolic Array



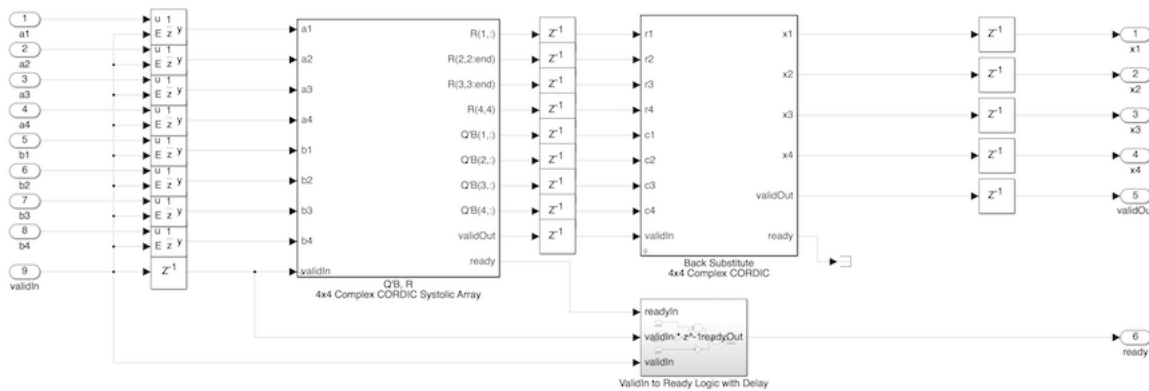
Copyright 2018 The MathWorks, Inc.

The enabled subsystems that send the data contain MATLAB function blocks that keep track of which input to send next, and send the rows of matrices A and B to block "4x4 Complex CORDIC Matrix Left Divide" when the ready signal is high. If you send data when ready is high, you will not invalidate any data already in the pipeline.

The algorithm overwrites matrix A with an upper-triangular matrix R. The algorithm overwrites B with $C = Q'B$ where Q is unitary and $QR = A$. The algorithm uses back-substitution on the upper-triangular matrix equation

$$RX = C.$$

To examine the algorithm, look Under the mask of block "4x4 Complex CORDIC Matrix Left Divide".



X is the solution to the equation $A \cdot X = B$ A is a 4x4 complex matrix B is a 4xp complex matrix X is a 4xp complex matrix
 If $B = \text{eye}(4)$, then $X = \text{inv}(A)$. $A = \begin{bmatrix} a1 \\ a2 \\ a3 \\ a4 \end{bmatrix}$ $B = \begin{bmatrix} b1 \\ b2 \\ b3 \\ b4 \end{bmatrix}$ $X = \begin{bmatrix} x1 \\ x2 \\ x3 \\ x4 \end{bmatrix}$

Define Parameters

Define complex matrices A and B in the base workspace. In this example, matrix A must be 4-by-4, and matrix B must be 4-by-p, where p is the number of right-hand sides.

```
rng('default');
A = complex(randn(4,4), randn(4,4));
B = complex(randn(4,1), randn(4,1));
```

The method uses the CORDIC algorithm, so you must also specify the number of iterations of the CORDIC kernel in the `NumberOfCORDICIterations` parameter, or on the block parameters of the "4x4 Complex CORDIC Matrix Left Divide" block.

When A and B are double-precision floating-point data types, set the number of CORDIC iterations to the number of bits in the mantissa of a double. If the inputs are fixed point, then the number of CORDIC iterations must be less than the word length. The accuracy of the computation improves one bit for each iteration, up to the word length of the inputs. This model will work with fixed-point, double, and single data types.

```
NumberOfCORDICIterations = 52;
```

You must also instantiate variable `BackSubstitutePrototype` to specify the data-type used in back-substitution, or enter a prototype value on the block parameters of the "4x4 Complex CORDIC Matrix Left Divide" block. In this case, matrices A and B are double, so set `BackSubstitutePrototype` value to be a double. This variable is used by the cast function to cast the back-substitute variables using the 'like' syntax.

```
BackSubstitutePrototype = 0;
```

Run the Model

Turn off expected warnings before running the model.

```
warning_state = warning('off', 'Coder:builtins:ConstantFoldingOverflow');
sim(model)
```

After simulation, the model returns matrix X, which is the solution to the matrix equation

$$AX = B$$

Verify the results by checking that AX-B is a small value.

```
err = norm(A*X - B)
```

```
err =
```

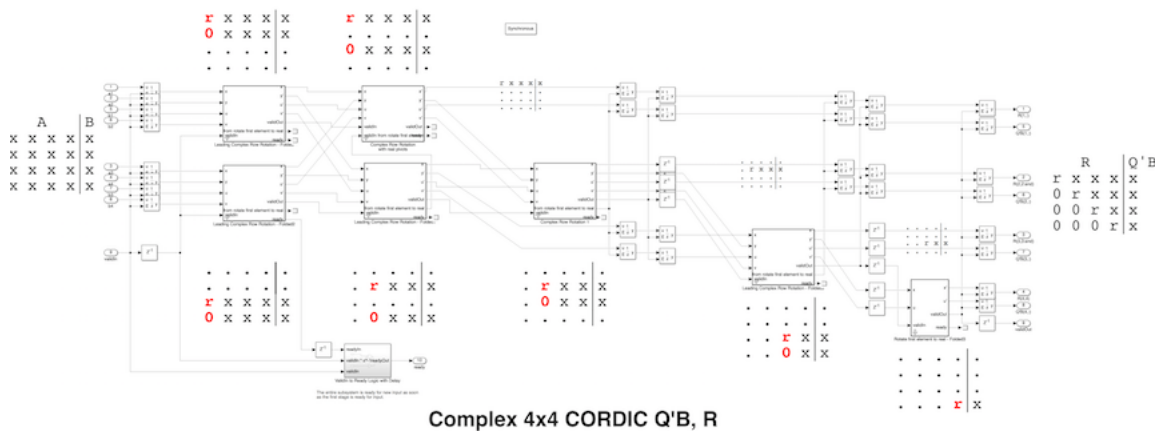
```
6.6743e-15
```

QR Algorithm

The CORDIC algorithm is used to compute the QR decomposition as in the examples Perform QR Factorization Using CORDIC, and Implement Hardware-Efficient QR Decomposition Using CORDIC in a Systolic Array.

This example is of a 4-by-4 matrix A. You can tile the blocks inside this model to build up to any size matrix.

To see the QR algorithm, look under the mask for block "Q'B, R 4x4 Complex CORDIC Systolic Array".



Complex QR

Because the data in this example is complex, an additional step is required to zero out the imaginary parts of the pivots to make them real-valued. This is done by splitting the data into real and imaginary parts and using CORDIC to zero out the imaginary part as if it were two rows of a real matrix. This is equivalent to the complex multiplication

$$e^{-i\theta} x = r$$

where

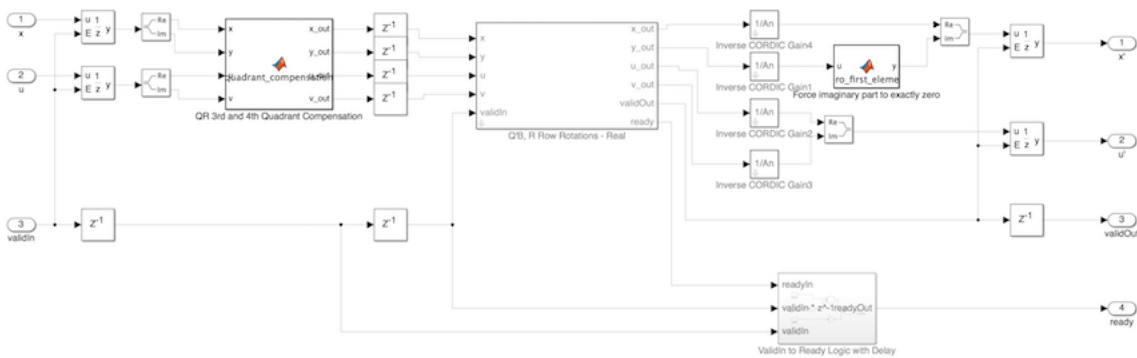
$$\theta = \tan^{-1} \left(\frac{\text{imag}(x)}{\text{real}(x)} \right)$$

and

$$r = \sqrt{\text{real}(x)^2 + \text{imag}(x)^2}$$

except that the computation is done using the CORDIC algorithm without squares or square-roots.

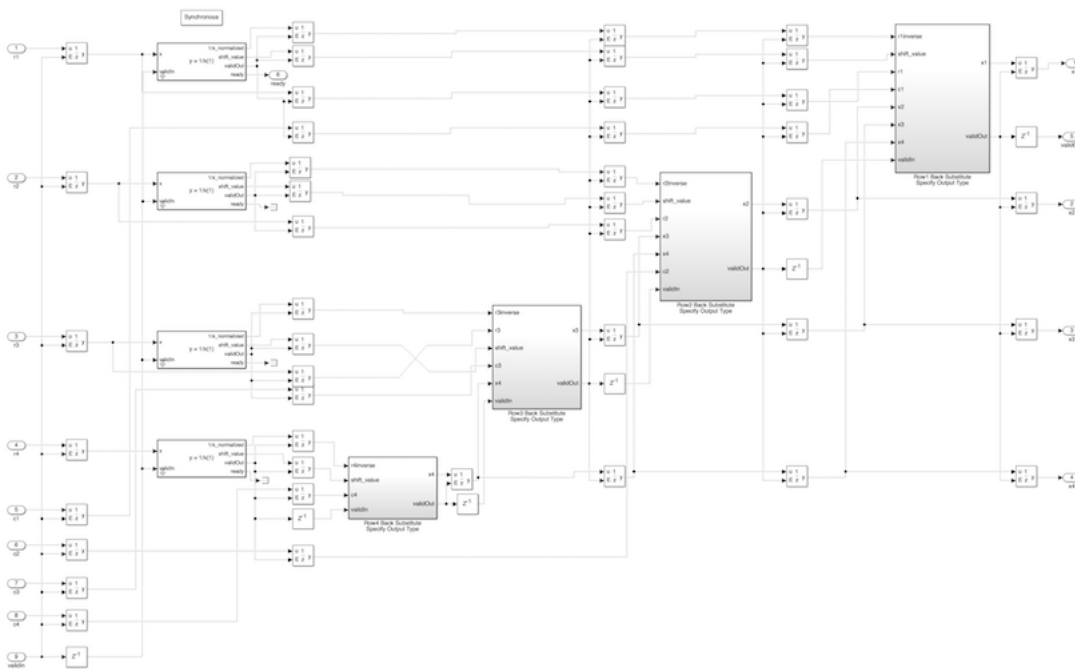
To see the algorithm for $e^{-i\theta}x$, look under the mask for block "Rotate first element to real".



Back-Substitute

Finally, to compute X, compute the reciprocals of the diagonal elements of R and back-substitute into the right-hand side C. The algorithm uses a CORDIC divide implementation to compute the reciprocals.

To see the back-substitute algorithm, look under the mask for block "Back Substitute 4x4 Complex CORDIC".



Matrix Inverse

It is usually unnecessary to explicitly compute the inverse of a matrix [3],[5]. However, if you want to do so, set B equal to the identity matrix I. Then, the solution of the equation

$$AX = I$$

equals

$$X = A^{-1}.$$

```
A = complex(2*rand(4,4)-1,2*rand(4,4)-1);
B = complex(eye(4));
NumberOfCORDICIterations = 52;
BackSubstitutePrototype = 0;
```

Simulate the model

```
sim(model)
```

Verify that AX and XA are close to the identity matrix. There will be small differences due to round-off errors.

```
A*X
```

```
ans =
```

```
1.0000 + 0.0000i    0.0000 - 0.0000i    0.0000 - 0.0000i   -0.0000 - 0.0000i
0.0000 - 0.0000i    1.0000 + 0.0000i   -0.0000 + 0.0000i   -0.0000 + 0.0000i
-0.0000 + 0.0000i    0.0000 - 0.0000i    1.0000 + 0.0000i    0.0000 + 0.0000i
0.0000 + 0.0000i    0.0000 + 0.0000i    0.0000 - 0.0000i    1.0000 - 0.0000i
```


X*A

ans =

```

1.0000 - 0.0000i -0.0000 - 0.0000i -0.0000 + 0.0000i 0.0000 - 0.0000i
0.0000 + 0.0000i 1.0000 + 0.0000i -0.0000 - 0.0000i 0.0000 + 0.0000i
0.0000 + 0.0000i 0.0000 + 0.0000i 1.0000 - 0.0000i -0.0000 + 0.0000i
0.0000 - 0.0000i -0.0000 + 0.0000i -0.0000 + 0.0000i 1.0000 + 0.0000i

```

Fractional Scaling

Normalizing to fractional types makes the computations easier in fixed-point. You can scale the matrices so that their data is in the range [-1, +1] and use fractional fixed-point types because the solution to the matrix equation

$$2^E AX = 2^E B$$

is the same as the solution to

$$AX = B.$$

To convert to fractional scaling with no additional cost in generated code, you can use the `reinterprecast` function in MATLAB or the Data Type Conversion block in Simulink with "Input and output to have equal Stored Integer".

Run Exhaustive Test Points

You can run many test inputs through the model by making A and B three-dimensional arrays.

```

m = 4;
n = 4;
p = 1;
n_test_inputs=100;

```

Create the inputs such that the real and imaginary parts are in the range [-1, +1]

```

A = complex(2*rand(m,n,n_test_inputs)-1, 2*rand(m,n,n_test_inputs)-1);
B = complex(2*rand(m,p,n_test_inputs)-1, 2*rand(m,p,n_test_inputs)-1);

```

In this example, matrices A and B contained real and imaginary parts already in the range [-1, +1], so set the types to be fractional.

```

data_word_length = 24;
data_fraction_length = 23;

```

QR Data Types

The growth in the elements of R in the real QR factorization is \sqrt{m} (see Perform QR Factorization Using CORDIC). The elements are complex in this example, so an additional growth factor of $\sqrt{2}$ is needed because

$$|1 + 1i| = \sqrt{2}.$$

Also, the CORDIC algorithm grows intermediate values by the following gain factor K_N where N is the number of CORDIC iterations before it is normalized out.

$$K_N = \prod_{j=0}^{N-1} \sqrt{1 + 2^{-2j}} \approx 1.6468$$

Therefore, an upperbound for growth in the QR algorithm for m-by-n complex matrix A is the product of all the growth factors:

$$K_N \cdot \sqrt{2m}.$$

In this example, A is 4-by-4, so the maximum growth factor in the QR algorithm is

$$K_n \cdot \sqrt{2m} \approx 1.6468 \cdot \sqrt{2 \cdot 4} = 4.6579.$$

Therefore, the number of additional integer bits to allow in the QR algorithm to avoid overflow when you have m=4 rows is:

```
qr_growth_bits = ceil(log2(1.6468*sqrt(2*m)))
```

```
qr_growth_bits =
```

```
3
```

Cast Matrices A and B to Fixed Point

The model requires that the inputs be signed and the word length of B be the same as the word length of A.

Grow the data wordlength to accommodate the QR growth.

```
qr_input_word_length = data_word_length + qr_growth_bits
```

```
qr_input_word_length =
```

```
27
```

Cast A to fixed point and B to A's type.

```
T.A = fi([], 1, qr_input_word_length, data_fraction_length);
T.B = T.A;
```

```
A = cast(A, 'like', T.A);
B = cast(B, 'like', T.B);
```

Back-Substitute Data Type

If A is an n -by- n invertible complex matrix, and x is the solution of matrix equation $Ax = b$, then using properties of vector and matrix norms (see, for example, reference [6]), it can be shown that

$$\|x\|_{\infty} \leq \frac{\sqrt{n}\|b\|_{\infty}}{\sigma_n}$$

where $\|x\|_\infty \equiv \max(\text{abs}(x))$ and σ_n is the smallest singular value of A . This bound is related to the condition number of A .

In this example, b is a complex vector with the real and imaginary parts bounded by 1. Therefore, $\|b\|_\infty \equiv \max(\text{abs}(b)) \leq \sqrt{2}$.

Therefore, if you know the distribution of singular values for the A matrices in your problem, then you can choose the number of additional integer bits required to avoid overflow in the back-substitute with a given probability (see, for example, reference [1]).

In this example, we compute the singular values for all the matrices in our test bench and choose the number of integer bits to avoid overflow based on that.

```
singular_values = zeros(n,n_test_inputs);
A0 = double(A);
for k = 1:n_test_inputs
    singular_values(:,k) = svd(A0(:,:,k));
end
condition_numbers = singular_values(1,:)./singular_values(end,:);
x_bound = sqrt(2*n)./singular_values(end,:);
```

The number of bits of growth to add is the base-2 logarithm of the maximum bound.

```
integer_bits_for_backsubstitute = ceil(log2(max(x_bound)))
```

```
integer_bits_for_backsubstitute =
```

```
7
```

Subtract the required integer bits for back-substitute from the wordlength, and an additional bit for the sign bit.

```
backsubstitute_fraction_length = T.A.WordLength - integer_bits_for_backsubstitute - 1
```

```
backsubstitute_fraction_length =
```

```
19
```

Therefore, the data type for the back-substitute is:

```
BackSubstitutePrototype = fi(0, 1, T.A.WordLength, backsubstitute_fraction_length)
```

```
BackSubstitutePrototype =
```

```
0
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 27
    FractionLength: 19
```

Set the number of CORDIC iterations to be one less than the fixed-point word length of A .

```
NumberOfCORDICIterations = T.A.WordLength - 1
```

```
NumberOfCORDICIterations =
```

```
    26
```

Run the model with fixed-point inputs.

```
sim(model)
```

Calculate and Plot the Errors

A measure of error is

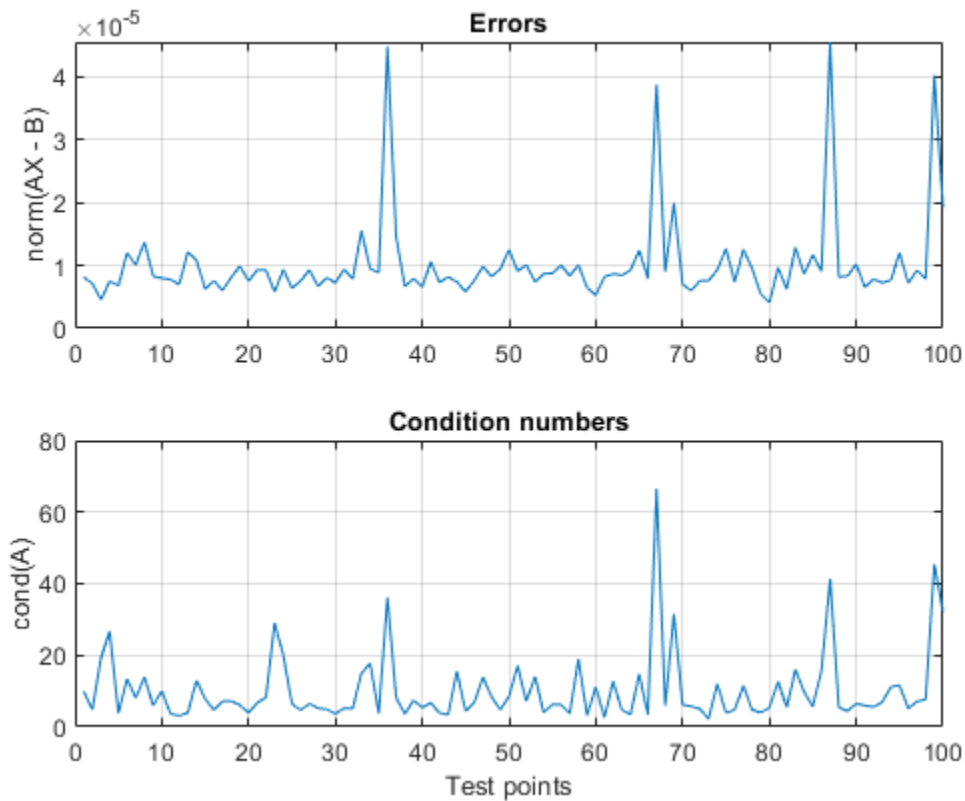
```
norm(A*X - B)
```

for each pair of inputs A and B.

```
norm_error = zeros(1,size(X,3));
B0 = double(B);
X0 = double(X);
for k = 1:size(X,3)
    norm_error(k) = norm(A0(:,:,k)*X0(:,:,k) - B0(:,:,k));
end
```

Plot the errors. The errors are low because of the data types chosen. The errors are typically higher when the condition number of matrix A is high, as the theory predicts.

```
figure(1)
clf
h1 = subplot(2,1,1);
plot(norm_error)
grid on
title('Errors')
ylabel('norm(A*X - B)')
h2 = subplot(2,1,2);
plot(condition_numbers)
grid on
title('Condition numbers')
ylabel('cond(A)')
xlabel('Test points')
linkaxes([h1,h2],'x')
```



Re-enable warnings

```
warning(warning_state);
```

References

- [1] Zizhong Chen and Jack J. Dongarra. "Condition Numbers of Gaussian Random Matrices". SIAM Journal on Matrix Analysis and Applications. 27.3 (July 2005), pp. 603-620.
- [2] George E. Forsythe and Cleve B. Moler. Computer Solution of Linear Algebraic Systems. Englewood Cliffs, N.J.: Prentice-Hall, 1967.
- [3] George E. Forsythe, M.A. Malcom and Cleve B. Moler. Computer Methods for Mathematical Computations. Englewood Cliffs, N.J.: Prentice-Hall, 1977.
- [4] Cleve B. Moler. Cleve's Corner: What is the Condition Number of a Matrix?, The MathWorks, Inc. 2017.
- [5] Cleve B. Moler. Numerical Computing with MATLAB. SIAM, 2004. isbn: 978-0-898716-60-3.
- [6] Gene H. Golub and Charles F. Van Loan. Matrix Computations. The Johns Hopkins University Press.

```
##ok<*NOPTS, *NASGU>
```

Implement Hardware-Efficient QR Decomposition Using CORDIC in a Systolic Array

This example shows how to compute the QR decomposition of matrices using hardware-efficient MATLAB® code in Simulink®.

To solve a system of equations or compute a least-squares solution to the matrix equation $AX = B$ using the QR decomposition, compute R and $Q'B$, where $QR = A$, and $RX = Q'B$. R is an upper triangular matrix and Q is an orthogonal matrix. If you just want Q and R , then set B to be the identity matrix.

In this example, R is computed from matrix A by applying Givens transformations using the CORDIC algorithm. $C = Q'B$ is computed from matrix B by applying the same Givens transformations. The algorithm uses only iterative shifts and additions to perform these computations.

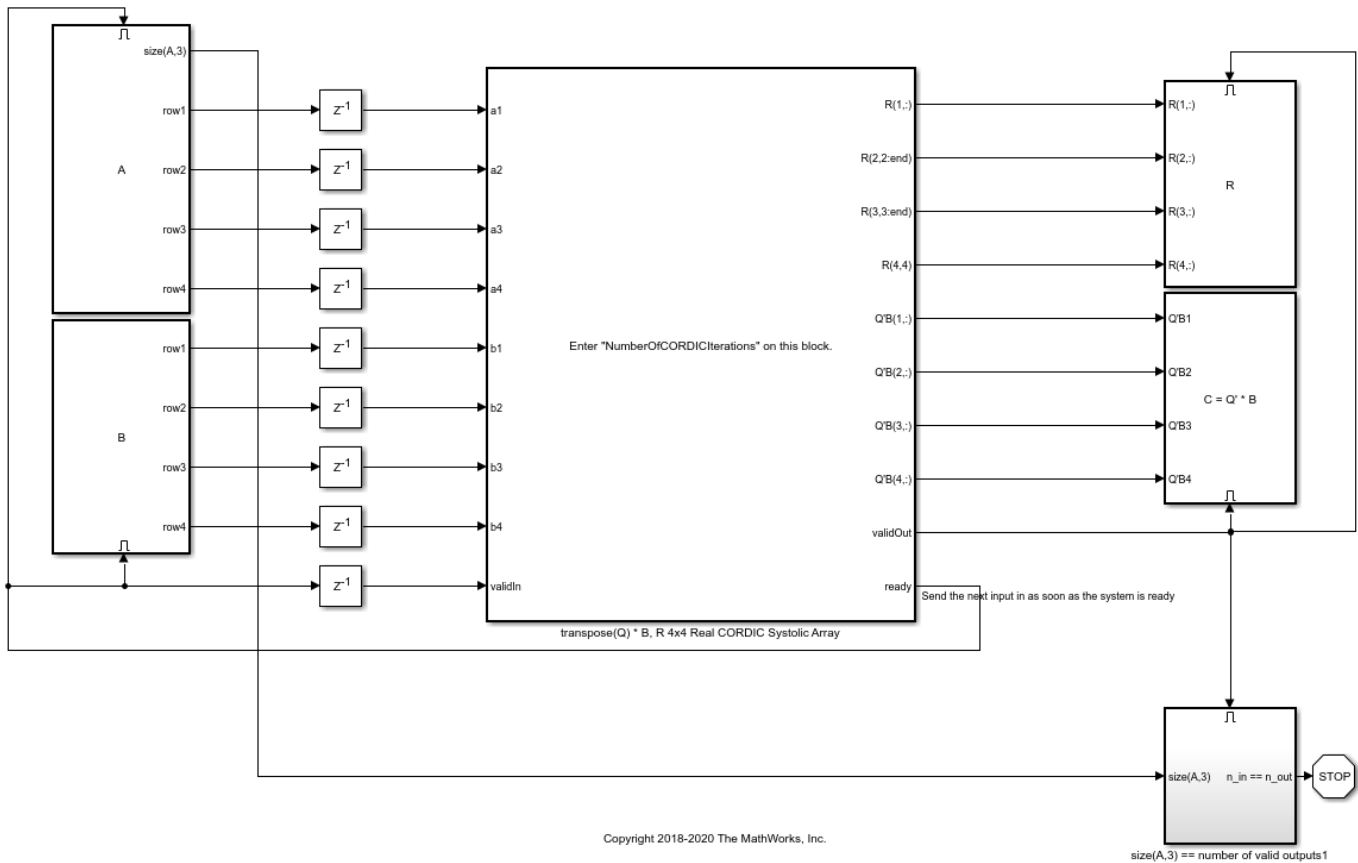
For more information on the algorithm used in this example, see [Perform QR Factorization Using CORDIC](#)

Overview

The Simulink model used in this example is:

```
fxpdemo_real_4x4_systolic_array_QR_model
```

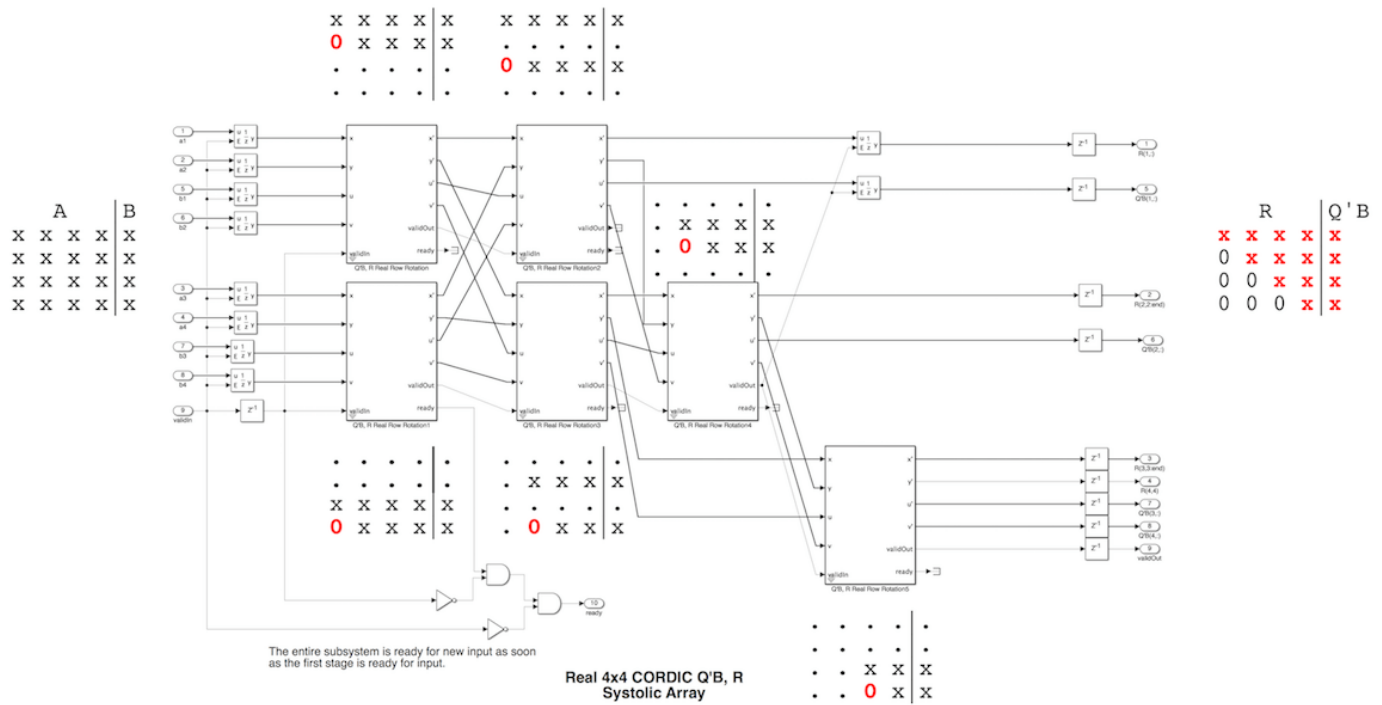
Implement Hardware-Efficient QR Decomposition Using CORDIC in a Systolic Array



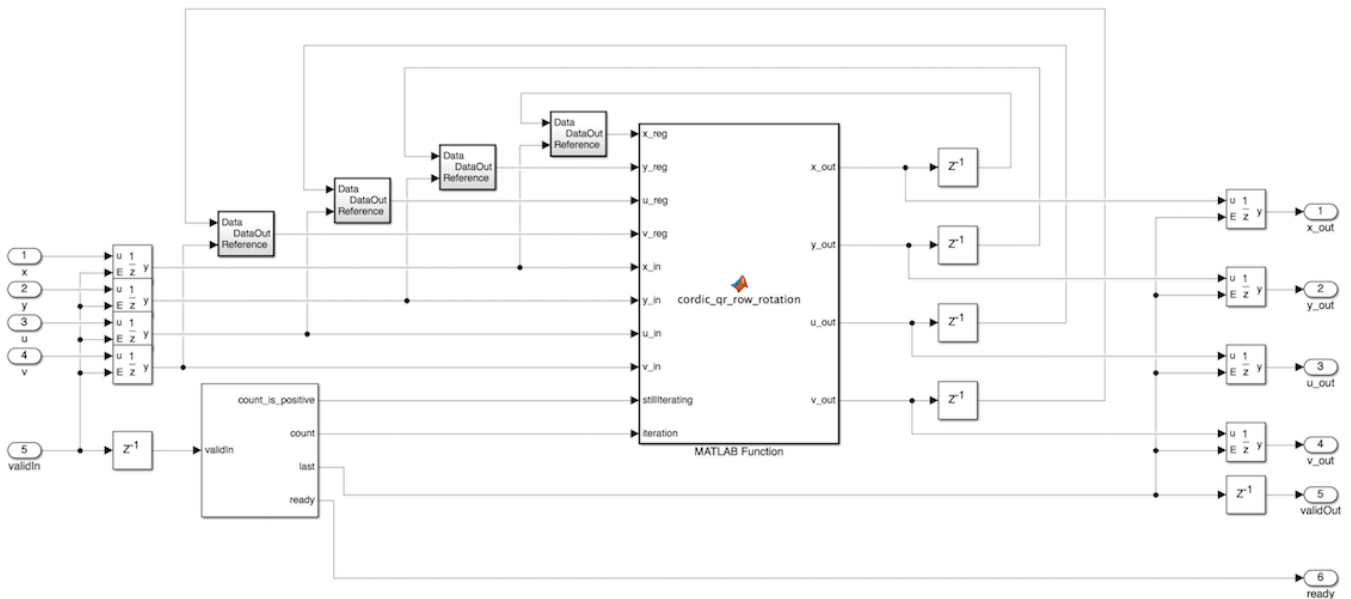
To enter your own input matrices, A and B, open the block parameters of the corresponding enabled subsystem blocks to the left. After simulation, the model returns the computed output matrices, $C = Q'B$ and R to the workspace. You can specify the number of CORDIC iterations in the block parameters of the $\text{transpose}(Q)*B, R\ 4\times 4$ Real CORDIC Systolic Array subsystem. If the inputs are fixed point, then the number of CORDIC iterations must be less than the word length. The accuracy of the computation improves one bit for each iteration, up to the word length of the inputs.

This model will work with fixed-point, double, and single data types.

To see how the algorithm performs the factorization, look under the mask of the $\text{transpose}(Q)*B, R\ 4\times 4$ Real CORDIC Systolic Array subsystem. The annotations indicate which rows of the matrix are being operated on to zero-out the sub-diagonal elements. The systolic array is set up for a 4-by-4 matrix A, but can be extended to any size by following the same pattern. This implementation works only with real input matrices.



To see the MATLAB code in the MATLAB Function block that performs the Givens transformations using CORDIC, continue to look under the block masks.



In this example, the number of rows and columns of A must be 4. Matrix B must have 4 rows and any number of columns.

Use QR to solve matrix equation $Ax = B$

The first step in solving the matrix equation $AX = B$ is to compute $RX = Q'B$, where R is upper-triangular, Q is orthogonal and $Q'R = A$.

The following inputs are double-precision floating-point types, so set the number of iterations to be 52, which is the number of bits in the mantissa of double.

```
format
NumberOfCORDICIterations = 52;
A = 2*rand(4,4)-1;
B = 2*rand(4,4)-1;
```

Simulate the model to compute R and $C = Q'B$.

```
sim fxdemo_real_4x4_systolic_array_QR_model
R
C
```

$R =$

```
1.5149   -0.0519   1.7292   -0.3224
      0    0.9593   -0.0259   -0.0879
      0      0     0.2565    1.0888
      0      0      0      -0.6429
```

$C =$

```
0.5942   -0.2382   0.0676   -0.9370
-0.8887   0.6146   -0.5758   0.3051
0.1725   0.7339   0.5409   0.5374
0.8540   1.1078   -0.2183   -0.5620
```

Verify that back-substituting with R and $C = Q'B$ gives the same results as MATLAB.

```
X = R\C
X_should_be = A\B
```

$X =$

```
-7.1245  -12.1131  -0.6637   1.4236
-0.8779   0.7572  -0.5511   0.3545
 6.3113  10.1759   0.6673  -1.6155
-1.3283  -1.7231   0.3396   0.8741
```

$X_should_be =$

```
-7.1245  -12.1131  -0.6637   1.4236
-0.8779   0.7572  -0.5511   0.3545
 6.3113  10.1759   0.6673  -1.6155
-1.3283  -1.7231   0.3396   0.8741
```

The norm of the difference between built-in MATLAB and the CORDIC QR solution should be small.

```
norm(X - X_should_be)
```

```
ans =
```

```
3.2578e-14
```

Compute Q and R

To compute Q and R, set B equal to the identity matrix. When B equals the identity matrix, then $Q = C'$.

```
NumberOfCORDICIterations = 52;
A = 2*rand(4,4)-1;
B = eye(size(A,1), 'like', A);
sim fxdemo_real_4x4_systolic_array_QR_model
```

```
Q = C';
```

The theoretical QR decomposition is $QR=A$, so the difference between the computed QR and A should be small.

```
norm(Q*R - A)
```

```
ans =
```

```
2.2861e-15
```

QR is not unique

The QR decomposition is only unique up to the signs of the rows of R and the columns of Q. You can make a unique QR decomposition by making the diagonal elements of R all positive.

```
D = diag(sign(diag(R)));
Qunique = Q*D
Runique = D*R
```

```
Qunique =
```

```
-0.3086    0.1224   -0.1033   -0.9376
-0.6277   -0.7636   -0.0952    0.1174
-0.5573    0.3930    0.7146    0.1559
 0.4474   -0.4975    0.6852   -0.2877
```

```
Runique =
```

```
1.4459   -0.8090    0.1547    0.3977
 0         1.1441    0.0809   -0.2494
 0         0         0.8193    0.1894
 0         0         0         0.4836
```

Then you can compare the computed QR from the model to the builtin MATLAB qr function.

```
[Q0,R0] = qr(A);
D0 = diag(sign(diag(R0)));
Q0 = Q0*D0
R0 = D0*R0
```

Q0 =

```
-0.3086    0.1224   -0.1033   -0.9376
-0.6277   -0.7636   -0.0952    0.1174
-0.5573    0.3930    0.7146    0.1559
 0.4474   -0.4975    0.6852   -0.2877
```

R0 =

```
 1.4459   -0.8090    0.1547    0.3977
         0    1.1441    0.0809   -0.2494
         0         0    0.8193    0.1894
         0         0         0    0.4836
```

Use Fixed-Point for Hardware-Efficient Implementation

Use fixed-point input data types to produce efficient HDL code for ASIC and FPGA devices.

For more information on how to choose fixed-point data types that will not overflow, refer to example Perform QR Factorization Using CORDIC.

You can run many test inputs through the model by making A and B 3-dimensional arrays.

```
n_test_inputs=100;
```

The following section defines random inputs for matrices A and B that are scaled between -1 and 1, so set the fixed-point word length to 18 bits and fraction length to 14 bits to allow for growth in the QR factorization and intermediate computations in the CORDIC algorithm.

```
word_length = 18;
fraction_length = 14;
```

The best-precision number of CORDIC iterations is the word length minus one. If the number of CORDIC iterations is set to smaller than `word_length - 1`, then the latency and clock ticks to next ready signal will be shorter, but it will be less accurate. The number of CORDIC iterations should not be set any larger because the generated code does not support shifts greater than the word length of a fixed-point type.

```
NumberOfCORDICIterations = word_length - 1
```

```
NumberOfCORDICIterations =
```

```
17
```

The random test inputs are concatenated so that at time k, the inputs are A(:, :, k) and B(:, :, k). Each element of A and B is a uniform random variable between -1 and +1.

```
A = 2*rand(4,4,n_test_inputs)-1;
```

Choose B to be the identity matrix so $Q=C'$.

```
B = eye(4);
B = repmat(B,1,1,n_test_inputs);
```

Cast A to fixed point, and cast B like A.

```
A = fi(A,1,word_length,fraction_length);
B = cast(B,'like',A);
```

Simulate the model

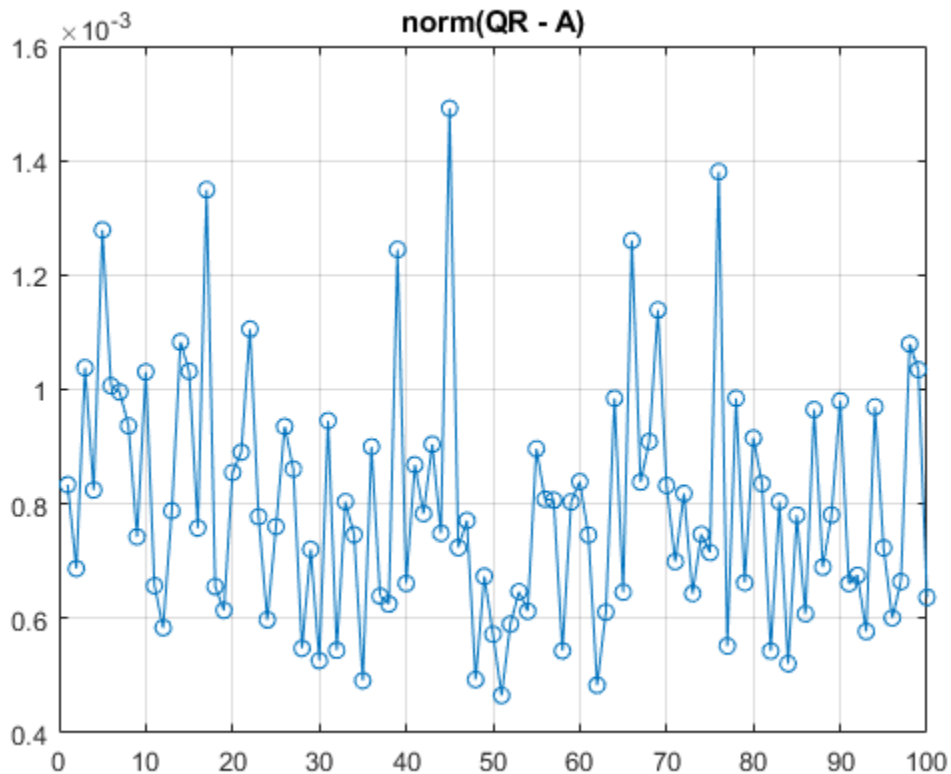
```
sim fxdemo_real_4x4_systolic_array_QR_model
```

Calculate and plot the errors

```
norm_error = zeros(1,size(R,3));
for k = 1:size(R,3)
    Q_times_R_minus_A = double(C(:,:,k))'*double(R(:,:,k)) - double(A(:,:,k));
    norm_error(k) = norm(Q_times_R_minus_A);
end
```

The errors should be on the order of 10^{-3} .

```
clf
plot(norm_error,'o-')
grid on
title('norm(QR - A)')
```



`%#ok<*NASGU,*NOPTS>`

Perform QR Factorization Using CORDIC

This example shows how to write MATLAB® code that works for both floating-point and fixed-point data types. The algorithm used in this example is the QR factorization implemented via CORDIC (Coordinate Rotation Digital Computer).

A good way to write an algorithm intended for a fixed-point target is to write it in MATLAB using builtin floating-point types so you can verify that the algorithm works. When you refine the algorithm to work with fixed-point types, then the best thing to do is to write it so that the same code continues working with floating-point. That way, when you are debugging, then you can switch the inputs back and forth between floating-point and fixed-point types to determine if a difference in behavior is because of fixed-point effects such as overflow and quantization versus an algorithmic difference. Even if the algorithm is not well suited for a floating-point target (as is the case of using CORDIC in the following example), it is still advantageous to have your MATLAB code work with floating-point for debugging purposes.

In contrast, you may have a completely different strategy if your target is floating point. For example, the QR algorithm is often done in floating-point with Householder transformations and row or column pivoting. But in fixed-point it is often more efficient to use CORDIC to apply Givens rotations with no pivoting.

This example addresses the first case, where your target is fixed-point, and you want an algorithm that is independent of data type because it is easier to develop and debug.

In this example you will learn various coding methods that can be applied across systems. The significant design patterns used in this example are the following:

- **Data Type Independence:** the algorithm is written in such a way that the MATLAB code is independent of data type, and will work equally well for fixed-point, double-precision floating-point, and single-precision floating-point.
- **Overflow Prevention:** method to guarantee not to overflow. This demonstrates how to prevent overflows in fixed-point.
- **Solving Systems of Equations:** method to use computational efficiency. Narrow your code scope by isolating what you need to define.

The main part in this example is an implementation of the QR factorization in fixed-point arithmetic using CORDIC for the Givens rotations. The algorithm is written in such a way that the MATLAB code is independent of data type, and will work equally well for fixed-point, double-precision floating-point, and single-precision floating-point.

The QR factorization of M-by-N matrix A produces an M-by-N upper triangular matrix R and an M-by-M orthogonal matrix Q such that $A = Q \cdot R$. A matrix is upper triangular if it has all zeros below the diagonal. An M-by-M matrix Q is orthogonal if $Q' \cdot Q = \text{eye}(M)$, the identity matrix.

The QR factorization is widely used in least-squares problems, such as the recursive least squares (RLS) algorithm used in adaptive filters.

The CORDIC algorithm is attractive for computing the QR algorithm in fixed-point because you can apply orthogonal Givens rotations with CORDIC using only shift and add operations.

Setup

So this example does not change your preferences or settings, we store the original state here, and restore them at the end.

```
originalFormat = get(0, 'format'); format short
originalFipref = get(fipref);      reset(fipref);
originalGlobalFimath = fimath;    resetglobalfimath;
```

Defining the CORDIC QR Algorithm

The CORDIC QR algorithm is given in the following MATLAB function, where A is an M -by- N real matrix, and $niter$ is the number of CORDIC iterations. Output Q is an M -by- M orthogonal matrix, and R is an M -by- N upper-triangular matrix such that $Q \cdot R = A$.

```
function [Q,R] = cordicqr(A,niter)
    Kn = inverse_cordic_growth_constant(niter);
    [m,n] = size(A);
    R = A;
    Q = coder.nullcopy(repmat(A(:,1),1,m)); % Declare type and size of Q
    Q(:) = eye(m); % Initialize Q
    for j=1:n
        for i=j+1:m
            [R(j,j:end),R(i,j:end),Q(:,j),Q(:,i)] = ...
                cordicgivens(R(j,j:end),R(i,j:end),Q(:,j),Q(:,i),niter,Kn);
        end
    end
end
```

This function was written to be independent of data type. It works equally well with builtin floating-point types (double and single) and with the fixed-point `fi` object.

One of the trickiest aspects of writing data-type independent code is to specify data type and size for a new variable. In order to preserve data types without having to explicitly specify them, the output R was set to be the same as input A , like this:

```
R = A;
```

In addition to being data-type independent, this function was written in such a way that MATLAB Coder™ will be able to generate efficient C code from it. In MATLAB, you most often declare and initialize a variable in one step, like this:

```
Q = eye(m)
```

However, `Q=eye(m)` would always produce Q as a double-precision floating point variable. If A is fixed-point, then we want Q to be fixed-point; if A is single, then we want Q to be single; etc.

Hence, you need to declare the type and size of Q in one step, and then initialize it in a second step. This gives MATLAB Coder the information it needs to create an efficient C program with the correct types and sizes. In the finished code you initialize output Q to be an M -by- M identity matrix and the same data type as A , like this:

```
Q = coder.nullcopy(repmat(A(:,1),1,m)); % Declare type and size of Q
Q(:) = eye(m); % Initialize Q
```

The `coder.nullcopy` function declares the size and type of Q without initializing it. The expansion of the first column of A with `repmat` won't appear in code generated by MATLAB; it is only used to

specify the size. The `repmat` function was used instead of `A(:,1:m)` because `A` may have more rows than columns, which will be the case in a least-squares problem. You have to be sure to always assign values to every element of an array when you declare it with `coder.nullcopy`, because if you don't then you will have uninitialized memory.

You will notice this pattern of assignment again and again. This is another key enabler of data-type independent code.

The heart of this function is applying orthogonal Givens rotations in-place to the rows of `R` to zero out sub-diagonal elements, thus forming an upper-triangular matrix. The same rotations are applied in-place to the columns of the identity matrix, thus forming orthogonal `Q`. The Givens rotations are applied using the `cordicgivens` function, as defined in the next section. The rows of `R` and columns of `Q` are used as both input and output to the `cordicgivens` function so that the computation is done in-place, overwriting `R` and `Q`.

```
[R(j,j:end),R(i,j:end),Q(:,j),Q(:,i)] = ...
    cordicgivens(R(j,j:end),R(i,j:end),Q(:,j),Q(:,i),niter,Kn);
```

Defining the CORDIC Givens Rotation

The `cordicgivens` function applies a Givens rotation by performing CORDIC iterations to rows $x=R(j,j:end)$, $y=R(i,j:end)$ around the angle defined by $x(1)=R(j,j)$ and $y(1)=R(i,j)$ where $i>j$, thus zeroing out $R(i,j)$. The same rotation is applied to columns $u = Q(:,j)$ and $v = Q(:,i)$, thus forming the orthogonal matrix `Q`.

```
function [x,y,u,v] = cordicgivens(x,y,u,v,niter,Kn)
    if x(1)<0
        % Compensation for 3rd and 4th quadrants
        x(:) = -x; u(:) = -u;
        y(:) = -y; v(:) = -v;
    end
    for i=0:niter-1
        x0 = x;
        u0 = u;
        if y(1)<0
            % Counter-clockwise rotation
            % x and y form R, u and v form Q
            x(:) = x - bitsra(y, i); u(:) = u - bitsra(v, i);
            y(:) = y + bitsra(x0,i); v(:) = v + bitsra(u0,i);
        else
            % Clockwise rotation
            % x and y form R, u and v form Q
            x(:) = x + bitsra(y, i); u(:) = u + bitsra(v, i);
            y(:) = y - bitsra(x0,i); v(:) = v - bitsra(u0,i);
        end
    end
    % Set y(1) to exactly zero so R will be upper triangular without round off
    % showing up in the lower triangle.
    y(1) = 0;
    % Normalize the CORDIC gain
    x(:) = Kn * x; u(:) = Kn * u;
    y(:) = Kn * y; v(:) = Kn * v;
end
```

The advantage of using CORDIC in fixed-point over the standard Givens rotation is that CORDIC does not use square root or divide operations. Only bit-shifts, addition, and subtraction are needed in the

main loop, and one scalar-vector multiply at the end to normalize the CORDIC gain. Also, CORDIC rotations work well in pipelined architectures.

The bit shifts in each iteration are performed with the bit shift right arithmetic (`bitsra`) function instead of `bitshift`, multiplication by 0.5, or division by 2, because `bitsra`

- generates more efficient embedded code,
- works equally well with positive and negative numbers,
- works equally well with floating-point, fixed-point and integer types, and
- keeps this code independent of data type.

It is worthwhile to note that there is a difference between sub-scripted assignment (`subsasgn`) into a variable `a(:) = b` versus overwriting a variable `a = b`. Sub-scripted assignment into a variable like this

```
x(:) = x + bitsra(y, i);
```

always preserves the type of the left-hand-side argument `x`. This is the recommended programming style in fixed-point. For example fixed-point types often grow their word length in a sum, which is governed by the `SumMode` property of the `fimath` object, so that the right-hand-side `x + bitsra(y, i)` can have a different data type than `x`.

If, instead, you overwrite the left-hand-side like this

```
x = x + bitsra(y, i);
```

then the left-hand-side `x` takes on the type of the right-hand-side sum. This programming style leads to changing the data type of `x` in fixed-point code, and is discouraged.

Defining the Inverse CORDIC Growth Constant

This function returns the inverse of the CORDIC growth factor after `niter` iterations. It is needed because CORDIC rotations grow the values by a factor of approximately 1.6468, depending on the number of iterations, so the gain is normalized in the last step of `cordicgivens` by a multiplication by the inverse $K_n = 1/1.6468 = 0.60725$.

```
function Kn = inverse_cordic_growth_constant(niter)
    Kn = 1/prod(sqrt(1+2.^(-2*(0:double(niter)-1))));
end
```

Exploring CORDIC Growth as a Function of Number of Iterations

The function for CORDIC growth is defined as

```
growth = prod(sqrt(1+2.^(-2*(0:double(niter)-1))))
```

and the inverse is

```
inverse_growth = 1 ./ growth
```

Growth is a function of the number of iterations `niter`, and quickly converges to approximately 1.6468, and the inverse converges to approximately 0.60725. You can see in the following table that the difference from one iteration to the next ceases to change after 27 iterations. This is because the calculation hit the limit of precision in double floating-point at 27 iterations.

| niter | growth | diff(growth) | 1./growth | diff(1./growth) |
|-------|--------------------|--------------|--------------------|-----------------|
| 0 | 1.0000000000000000 | 0 | 1.0000000000000000 | 0 |

```

1  1.414213562373095  0.414213562373095  0.707106781186547  -0.292893218813453
2  1.581138830084190  0.166925267711095  0.632455532033676  -0.074651249152872
3  1.629800601300662  0.048661771216473  0.613571991077896  -0.018883540955780
4  1.642484065752237  0.012683464451575  0.608833912517752  -0.004738078560144
5  1.645688915757255  0.003204850005018  0.607648256256168  -0.001185656261584
6  1.646492278712479  0.000803362955224  0.607351770141296  -0.000296486114872
7  1.646693254273644  0.000200975561165  0.607277644093526  -0.000074126047770
8  1.646743506596901  0.000050252323257  0.607259112298893  -0.000018531794633
9  1.646756070204878  0.000012563607978  0.607254479332562  -0.000004632966330
10 1.646759211139822  0.000003140934944  0.607253321089875  -0.000001158242687
11 1.646759996375617  0.000000785235795  0.607253031529134  -0.000000289560741
12 1.646760192684695  0.000000196309077  0.607252959138945  -0.000000072390190
13 1.646760241761972  0.000000049077277  0.607252941041397  -0.000000018097548
14 1.646760254031292  0.000000012269320  0.607252936517010  -0.000000004524387
15 1.646760257098622  0.000000003067330  0.607252935385914  -0.000000001131097
16 1.646760257865455  0.000000000766833  0.607252935103139  -0.000000000282774
17 1.646760258057163  0.000000000191708  0.607252935032446  -0.000000000070694
18 1.646760258105090  0.000000000047927  0.607252935014772  -0.000000000017673
19 1.646760258117072  0.000000000011982  0.607252935010354  -0.000000000004418
20 1.646760258120067  0.000000000002995  0.607252935009249  -0.000000000001105
21 1.646760258120816  0.000000000000749  0.607252935008973  -0.000000000000276
22 1.646760258121003  0.000000000000187  0.607252935008904  -0.000000000000069
23 1.646760258121050  0.000000000000047  0.607252935008887  -0.000000000000017
24 1.646760258121062  0.000000000000012  0.607252935008883  -0.000000000000004
25 1.646760258121065  0.000000000000003  0.607252935008882  -0.000000000000001
26 1.646760258121065  0.000000000000001  0.607252935008881  -0.000000000000000
27 1.646760258121065  0 0.607252935008881 0
28 1.646760258121065  0 0.607252935008881 0
29 1.646760258121065  0 0.607252935008881 0
30 1.646760258121065  0 0.607252935008881 0
31 1.646760258121065  0 0.607252935008881 0
32 1.646760258121065  0 0.607252935008881 0

```

Comparing CORDIC to the Standard Givens Rotation

The `cordicgivens` function is numerically equivalent to the following standard Givens rotation algorithm from Golub & Van Loan, *Matrix Computations*. In the `cordicqr` function, if you replace the call to `cordicgivens` with a call to `givensrotation`, then you will have the standard Givens QR algorithm.

```

function [x,y,u,v] = givensrotation(x,y,u,v)
    a = x(1); b = y(1);
    if b==0
        % No rotation necessary.  c = 1; s = 0;
        return;
    else
        if abs(b) > abs(a)
            t = -a/b; s = 1/sqrt(1+t^2); c = s*t;
        else
            t = -b/a; c = 1/sqrt(1+t^2); s = c*t;
        end
    end
    x0 = x;          u0 = u;
    % x and y form R,  u and v form Q
    x(:) = c*x0 - s*y;  u(:) = c*u0 - s*v;
    y(:) = s*x0 + c*y;  v(:) = s*u0 + c*v;
end

```

The `givensrotation` function uses division and square root, which are expensive in fixed-point, but good for floating-point algorithms.

Example of CORDIC Rotations

Here is a 3-by-3 example that follows the CORDIC rotations through each step of the algorithm. The algorithm uses orthogonal rotations to zero out the subdiagonal elements of R using the diagonal elements as pivots. The same rotations are applied to the identity matrix, thus producing orthogonal Q such that $Q^*R = A$.

Let A be a random 3-by-3 matrix, and initialize $R = A$, and $Q = \text{eye}(3)$.

$$R = A = \begin{bmatrix} -0.8201 & 0.3573 & -0.0100 \\ -0.7766 & -0.0096 & -0.7048 \\ -0.7274 & -0.6206 & -0.8901 \end{bmatrix}$$

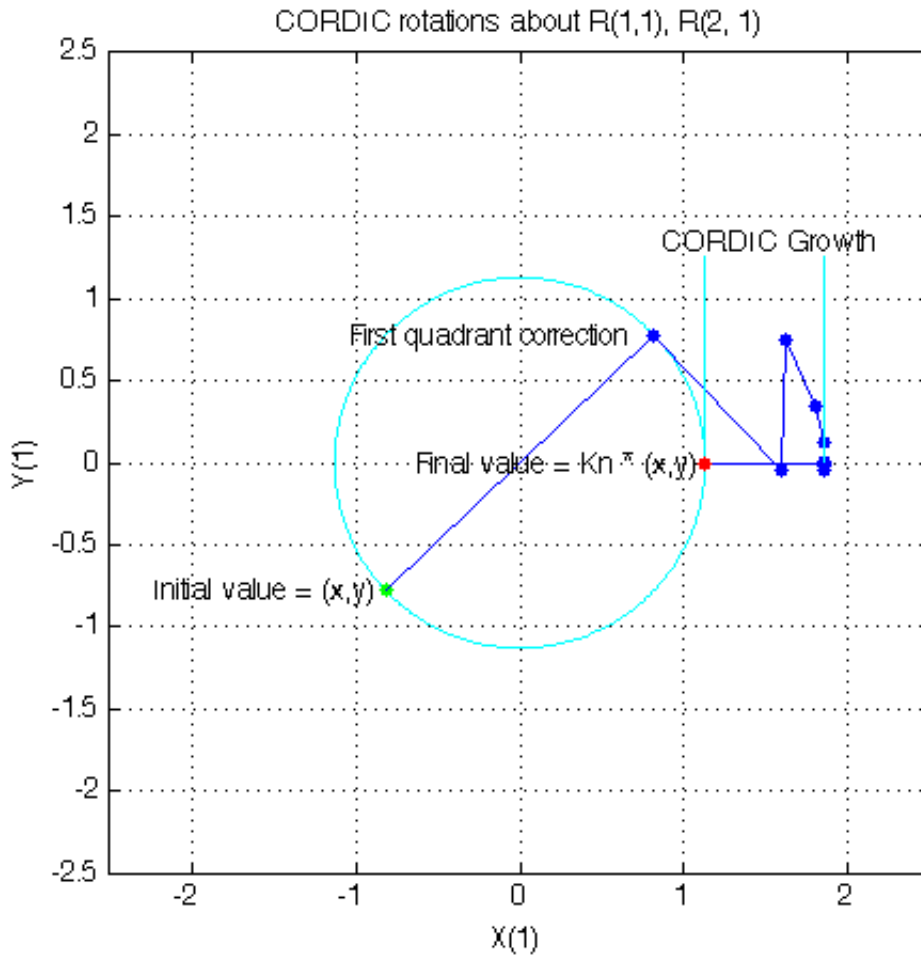
$$Q = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

The first rotation is about the first and second row of R and the first and second column of Q. Element $R(1,1)$ is the pivot and $R(2,1)$ rotates to 0.

$$\begin{array}{l} \text{R before the first rotation} \\ x \begin{bmatrix} -0.8201 & 0.3573 & -0.0100 \end{bmatrix} \\ y \begin{bmatrix} -0.7766 & -0.0096 & -0.7048 \\ -0.7274 & -0.6206 & -0.8901 \end{bmatrix} \end{array} \rightarrow \begin{array}{l} \text{R after the first rotation} \\ x \begin{bmatrix} 1.1294 & -0.2528 & 0.4918 \end{bmatrix} \\ y \begin{bmatrix} 0 & 0.2527 & 0.5049 \\ -0.7274 & -0.6206 & -0.8901 \end{bmatrix} \end{array}$$

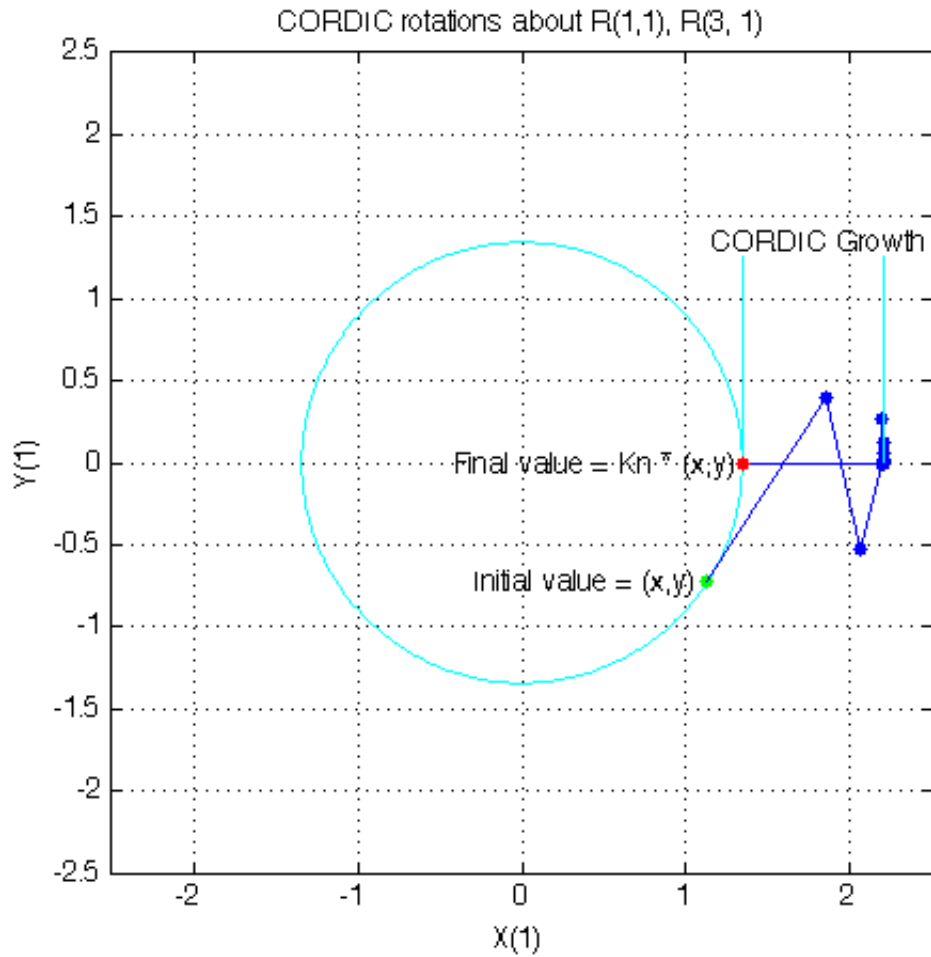
$$\begin{array}{l} \text{Q before the first rotation} \\ u \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \\ v \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \end{array} \rightarrow \begin{array}{l} \text{Q after the first rotation} \\ u \begin{bmatrix} -0.7261 \\ -0.6876 \\ 0 \end{bmatrix} \\ v \begin{bmatrix} 0.6876 \\ -0.7261 \\ 0 \end{bmatrix} \end{array}$$

In the following plot, you can see the growth in x in each of the CORDIC iterations. The growth is factored out at the last step by multiplying it by $K_n = 0.60725$. You can see that $y(1)$ iterates to 0. Initially, the point $[x(1), y(1)]$ is in the third quadrant, and is reflected into the first quadrant before the start of the CORDIC iterations.



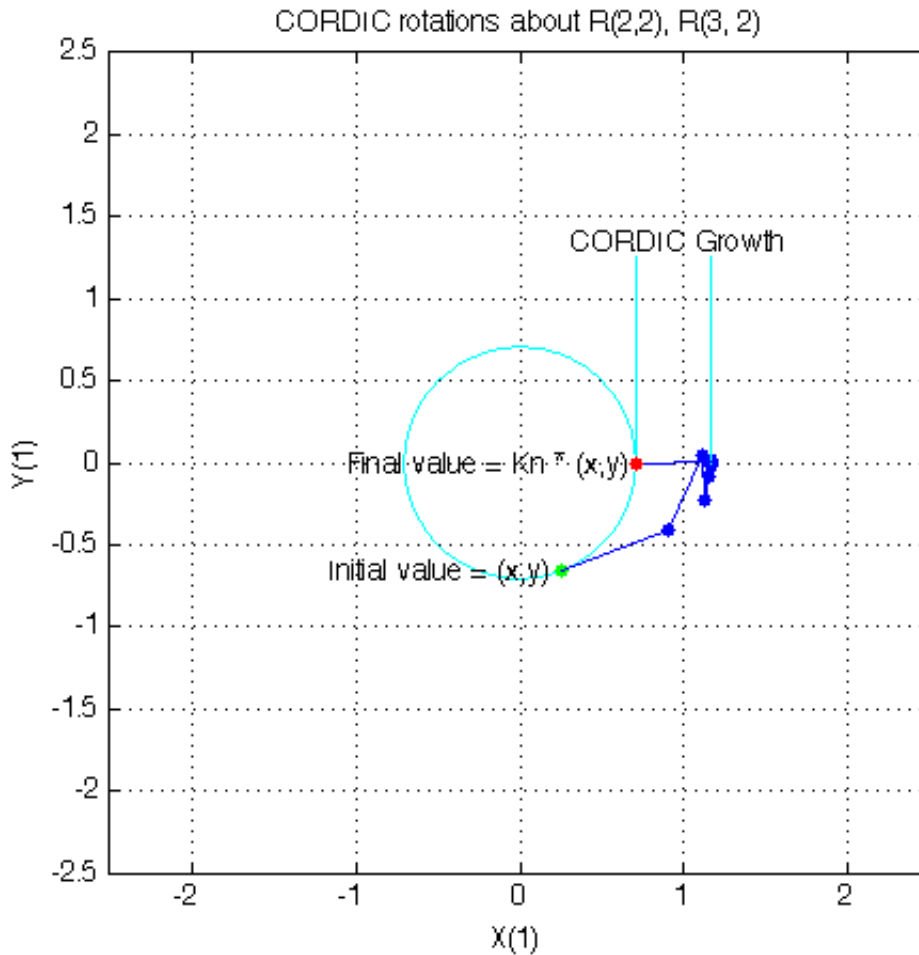
The second rotation is about the first and third row of R and the first and third column of Q. Element $R(1, 1)$ is the pivot and $R(3, 1)$ rotates to 0.

| | | | | | |
|------------------------------|-------------------------------|-----|-----------------------------|----------------------------|-----------|
| R before the second rotation | | -> | R after the second rotation | | |
| x | [1.1294 -0.2528 0.4918] | | x | [1.3434 0.1235 0.8954] | |
| | 0 0.2527 0.5049 | | | 0 0.2527 0.5049 | |
| y | [-0.7274] -0.6206 -0.8901 | -> | y | [0 -0.6586 -0.4820] | |
| Q before the second rotation | | | Q after the second rotation | | |
| u | | | u | | |
| | v | | | v | |
| [-0.7261] | 0.6876 | [0] | [-0.6105] | 0.6876 | [-0.3932] |
| [-0.6876] | -0.7261 | [0] | [-0.5781] | -0.7261 | [-0.3723] |
| [0] | 0 | [1] | [-0.5415] | 0 | [0.8407] |



The third rotation is about the second and third row of R and the second and third column of Q. Element R(2, 2) is the pivot and R(3, 2) rotates to 0.

| | | | | | | | | | |
|---|-----------------------------|-----------|-----------|----|---|----------------------------|-----------|-----------|--|
| | R before the third rotation | | | | | R after the third rotation | | | |
| | 1.3434 | 0.1235 | 0.8954 | | | 1.3434 | 0.1235 | 0.8954 | |
| x | 0 | [0.2527 | 0.5049] | -> | x | 0 | [0.7054 | 0.6308] | |
| y | 0 | [-0.6586 | -0.4820] | -> | y | 0 | [0 | 0.2987] | |
| | Q before the third rotation | | | | | Q after the third rotation | | | |
| | | u | v | | | | u | v | |
| | -0.6105 | [0.6876] | [-0.3932] | | | -0.6105 | [0.6134] | [0.5011] | |
| | -0.5781 | [-0.7261] | [-0.3723] | -> | | -0.5781 | [0.0875] | [-0.8113] | |
| | -0.5415 | [0 | 0.8407] | | | -0.5415 | [-0.7849] | [0.3011] | |



This completes the QR factorization. R is upper triangular, and Q is orthogonal.

$$R = \begin{bmatrix} 1.3434 & 0.1235 & 0.8954 \\ 0 & 0.7054 & 0.6308 \\ 0 & 0 & 0.2987 \end{bmatrix}$$

$$Q = \begin{bmatrix} -0.6105 & 0.6134 & 0.5011 \\ -0.5781 & 0.0875 & -0.8113 \\ -0.5415 & -0.7849 & 0.3011 \end{bmatrix}$$

You can verify that Q is within roundoff error of being orthogonal by multiplying and seeing that it is close to the identity matrix.

$$Q*Q' = \begin{bmatrix} 1.0000 & 0.0000 & 0.0000 \\ 0.0000 & 1.0000 & 0 \\ 0.0000 & 0 & 1.0000 \end{bmatrix}$$

$$Q'*Q = \begin{bmatrix} 1.0000 & 0.0000 & -0.0000 \\ 0.0000 & 1.0000 & -0.0000 \\ -0.0000 & -0.0000 & 1.0000 \end{bmatrix}$$

You can see the error difference by subtracting the identity matrix.

```
Q*Q' - eye(size(Q)) =
           0    2.7756e-16    3.0531e-16
    2.7756e-16    4.4409e-16           0
    3.0531e-16           0    6.6613e-16
```

You can verify that Q^*R is close to A by subtracting to see the error difference.

```
Q*R - A =
-3.7802e-11  -7.2325e-13  -2.7756e-17
-3.0512e-10   1.1708e-12  -4.4409e-16
 3.6836e-10  -4.3487e-13  -7.7716e-16
```

Determining the Optimal Output Type of Q for Fixed Word Length

Since Q is orthogonal, you know that all of its values are between -1 and $+1$. In floating-point, there is no decision about the type of Q : it should be the same floating-point type as A . However, in fixed-point, you can do better than making Q have the identical fixed-point type as A . For example, if A has word length 16 and fraction length 8, and if we make Q also have word length 16 and fraction length 8, then you force Q to be less accurate than it could be and waste the upper half of the fixed-point range.

The best type for Q is to make it have full range of its possible outputs, plus accommodate the 1.6468 CORDIC growth factor in intermediate calculations. Therefore, assuming that the word length of Q is the same as the word length of input A , then the best fraction length for Q is 2 bits less than the word length (one bit for 1.6468 and one bit for the sign).

Hence, our initialization of Q in `cordicqr` can be improved like this.

```
if isfi(A) && (isfixed(A) || isscaleddouble(A))
    Q = fi(one*eye(m), get(A, 'NumericType'), ...
           'FractionLength', get(A, 'WordLength')-2);
else
    Q = coder.nullcopy(repmat(A(:,1), 1, m));
    Q(:) = eye(m);
end
```

A slight disadvantage is that this section of code is dependent on data type. However, you gain a major advantage by picking the optimal type for Q , and the main algorithm is still independent of data type. You can do this kind of input parsing in the beginning of a function and leave the main algorithm data-type independent.

Preventing Overflow in Fixed Point R

This section describes how to determine a fixed-point output type for R in order to prevent overflow. In order to pick an output type, you need to know how much the magnitude of the values of R will grow.

Given real matrix A and its QR factorization computed by Givens rotations without pivoting, an upper-bound on the magnitude of the elements of R is the square-root of the number of rows of A times the magnitude of the largest element in A . Furthermore, this growth will never be greater during an intermediate computation. In other words, let $[m, n] = \text{size}(A)$, and $[Q, R] = \text{givensqr}(A)$. Then

```
max(abs(R(:))) <= sqrt(m) * max(abs(A(:))).
```

This is true because each element of R is formed from orthogonal rotations from its corresponding column in A , so the largest that any element $R(i, j)$ can get is if all of the elements of

its corresponding column $A(:, j)$ were rotated to a single value. In other words, the largest possible value will be bounded by the 2-norm of $A(:, j)$. Since the 2-norm of $A(:, j)$ is equal to the square-root of the sum of the squares of the m elements, and each element is less-than-or-equal-to the largest element of A , then

$$\text{norm}(A(:, j)) \leq \sqrt{m} * \max(\text{abs}(A(:))).$$

That is, for all j

$$\begin{aligned} \text{norm}(A(:, j)) &= \sqrt{A(1, j)^2 + A(2, j)^2 + \dots + A(m, j)^2} \\ &\leq \sqrt{m * \max(\text{abs}(A(:)))^2} \\ &= \sqrt{m} * \max(\text{abs}(A(:))). \end{aligned}$$

and so for all i, j

$$\text{abs}(R(i, j)) \leq \text{norm}(A(:, j)) \leq \sqrt{m} * \max(\text{abs}(A(:))).$$

Hence, it is also true for the largest element of R

$$\max(\text{abs}(R(:))) \leq \sqrt{m} * \max(\text{abs}(A(:))).$$

This becomes useful in fixed-point where the elements of A are often very close to the maximum value attainable by the data type, so we can set a tight upper bound without knowing the values of A . This is important because we want to set an output type for R with a minimum number of bits, only knowing the upper bound of the data type of A . You can use `fi` method `upperbound` to get this value.

Therefore, for all i, j

$$\text{abs}(R(i, j)) \leq \sqrt{m} * \text{upperbound}(A)$$

Note that $\sqrt{m} * \text{upperbound}(A)$ is also an upper bound for the elements of A :

$$\text{abs}(A(i, j)) \leq \text{upperbound}(A) \leq \sqrt{m} * \text{upperbound}(A)$$

Therefore, when picking fixed-point data types, $\sqrt{m} * \text{upperbound}(A)$ is an upper bound that will work for both A and R .

Attaining the maximum is easy and common. The maximum will occur when all elements get rotated into a single element, like the following matrix with orthogonal columns:

$$A = \begin{bmatrix} 7 & -7 & 7 & 7 \\ 7 & 7 & -7 & 7 \\ 7 & -7 & -7 & -7 \\ 7 & 7 & 7 & -7 \end{bmatrix};$$

Its maximum value is 7 and its number of rows is $m=4$, so we expect that the maximum value in R will be bounded by $\max(\text{abs}(A(:))) * \sqrt{m} = 7 * \sqrt{4} = 14$. Since A in this example is orthogonal, each column gets rotated to the max value on the diagonal.

```
niter = 52;
[Q,R] = cordicqr(A,niter)
```

$Q =$

$$\begin{bmatrix} 0.5000 & -0.5000 & 0.5000 & 0.5000 \\ 0.5000 & 0.5000 & -0.5000 & 0.5000 \\ 0.5000 & -0.5000 & -0.5000 & -0.5000 \end{bmatrix}$$


```

0.5000    0.5000    0.5000   -0.5000

```

R =

```

14.0000    0.0000   -0.0000   -0.0000
    0    14.0000   -0.0000    0.0000
    0         0    14.0000    0.0000
    0         0         0    14.0000

```

Another simple example of attaining maximum growth is a matrix that has all identical elements, like a matrix of all ones. A matrix of ones will get rotated into $1*\sqrt{m}$ in the first row and zeros elsewhere. For example, this 9-by-5 matrix will have all $1*\sqrt{9}=3$ in the first row of R.

```

m = 9; n = 5;
A = ones(m,n)
niter = 52;
[Q,R] = cordicqr(A,niter)

```

A =

```

1    1    1    1    1
1    1    1    1    1
1    1    1    1    1
1    1    1    1    1
1    1    1    1    1
1    1    1    1    1
1    1    1    1    1
1    1    1    1    1
1    1    1    1    1

```

Q =

Columns 1 through 7

```

0.3333    0.5567   -0.6784    0.3035   -0.1237    0.0503    0.0158
0.3333    0.0296    0.2498   -0.1702   -0.6336    0.1229   -0.3012
0.3333    0.2401    0.0562   -0.3918    0.4927    0.2048   -0.5395
0.3333    0.0003    0.0952   -0.1857    0.2148    0.4923    0.7080
0.3333    0.1138    0.0664   -0.2263    0.1293   -0.8348    0.2510
0.3333   -0.3973   -0.0143    0.3271    0.4132   -0.0354   -0.2165
0.3333    0.1808    0.3538   -0.1012   -0.2195         0    0.0824
0.3333   -0.6500   -0.4688   -0.2380   -0.2400         0         0
0.3333   -0.0740    0.3400    0.6825   -0.0331         0         0

```

Columns 8 through 9

```

0.0056   -0.0921
-0.5069   -0.1799
0.0359    0.3122
-0.2351   -0.0175
-0.2001    0.0610
-0.0939   -0.6294
0.7646   -0.2849
0.2300    0.2820

```

```

0      0.5485

R =

3.0000  3.0000  3.0000  3.0000  3.0000
0      0.0000  0.0000  0.0000  0.0000
0      0      0.0000  0.0000  0.0000
0      0      0      0.0000  0.0000
0      0      0      0      0.0000
0      0      0      0      0
0      0      0      0      0
0      0      0      0      0
0      0      0      0      0

```

As in the `cordicqr` function, the Givens QR algorithm is often written by overwriting `A` in-place with `R`, so being able to cast `A` into `R`'s data type at the beginning of the algorithm is convenient.

In addition, if you compute the Givens rotations with CORDIC, there is a growth-factor that converges quickly to approximately 1.6468. This growth factor gets normalized out after each Givens rotation, but you need to accommodate it in the intermediate calculations. Therefore, the number of additional bits that are required including the Givens and CORDIC growth are $\log_2(1.6468 * \text{sqrt}(m))$. The additional bits of head-room can be added either by increasing the word length, or decreasing the fraction length.

A benefit of increasing the word length is that it allows for the maximum possible precision for a given word length. A disadvantage is that the optimal word length may not correspond to a native type on your processor (e.g. increasing from 16 to 18 bits), or you may have to increase to the next larger native word size which could be quite large (e.g. increasing from 16 to 32 bits, when you only needed 18).

A benefit of decreasing fraction length is that you can do the computation in-place in the native word size of `A`. A disadvantage is that you lose precision.

Another option is to pre-scale the input by right-shifting. This is equivalent to decreasing the fraction length, with the additional disadvantage of changing the scaling of your problem. However, this may be an attractive option to you if you prefer to only work in fractional arithmetic or integer arithmetic.

Example of Fixed Point Growth in R

If you have a fixed-point input matrix `A`, you can define fixed-point output `R` with the growth defined in the previous section.

Start with a random matrix `X`.

```

X = [0.0513  -0.2097   0.9492   0.2614
      0.8261   0.6252   0.3071  -0.9415
      1.5270   0.1832   0.1352  -0.1623
      0.4669  -1.0298   0.5152  -0.1461];

```

Create a fixed-point `A` from `X`.

```
A = sfi(X)
```

```
A =
```

```

0.0513   -0.2097   0.9492   0.2614
0.8261   0.6252   0.3071  -0.9415
1.5270   0.1832   0.1352  -0.1623
0.4669  -1.0298   0.5152  -0.1461

```

```

        DataTypeMode: Fixed-point: binary point scaling
          Signedness: Signed
          WordLength: 16
        FractionLength: 14

```

```
m = size(A,1)
```

```
m =
```

```
4
```

The growth factor is 1.6468 times the square-root of the number of rows of A. The bit growth is the next integer above the base-2 logarithm of the growth.

```
bit_growth = ceil(log2(cordic_growth_constant * sqrt(m)))
```

```
bit_growth =
```

```
2
```

Initialize R with the same values as A, and a word length increased by the bit growth.

```
R = sfi(A, get(A, 'WordLength')+bit_growth, get(A, 'FractionLength'))
```

```
R =
```

```

0.0513   -0.2097   0.9492   0.2614
0.8261   0.6252   0.3071  -0.9415
1.5270   0.1832   0.1352  -0.1623
0.4669  -1.0298   0.5152  -0.1461

```

```

        DataTypeMode: Fixed-point: binary point scaling
          Signedness: Signed
          WordLength: 18
        FractionLength: 14

```

Use R as input and overwrite it.

```
niter = get(R, 'WordLength') - 1
[Q,R] = cordicqr(R, niter)
```

```
niter =
```

```
17
```

```
Q =
```

```

0.0284  -0.1753   0.9110   0.3723
0.4594   0.4470   0.3507  -0.6828
0.8490   0.0320  -0.2169   0.4808
0.2596  -0.8766  -0.0112  -0.4050

```

```

      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 18
      FractionLength: 16

```

R =

```

1.7989   0.1694   0.4166  -0.6008
   0      1.2251  -0.4764  -0.3438
   0           0    0.9375  -0.0555
   0           0           0    0.7214

```

```

      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 18
      FractionLength: 14

```

Verify that $Q*Q'$ is near the identity matrix.

```
double(Q)*double(Q')
```

ans =

```

1.0000  -0.0001   0.0000   0.0000
-0.0001   1.0001   0.0000  -0.0000
0.0000   0.0000   1.0000  -0.0000
0.0000  -0.0000  -0.0000   1.0000

```

Verify that $Q*R - A$ is small relative to the precision of A.

```
err = double(Q)*double(R) - double(A)
```

err =

```

1.0e-03 *
-0.1048  -0.2355   0.1829  -0.2146
0.3472   0.2949   0.0260  -0.2570
0.2776  -0.1740  -0.1007   0.0966
0.0138  -0.1558   0.0417  -0.0362

```

Increasing Precision in R

The previous section showed you how to prevent overflow in R while maintaining the precision of A. If you leave the fraction length of R the same as A, then R cannot have more precision than A, and your precision requirements may be such that the precision of R must be greater.

An extreme example of this is to define a matrix with an integer fixed-point type (i.e. fraction length is zero). Let matrix X have elements that are the full range for signed 8 bit integers, between -128 and +127.

```
X = [-128 -128 -128 127
      -128 127 127 -128
       127 127 127 127
       127 127 -128 -128];
```

Define fixed-point A to be equivalent to an 8-bit integer.

```
A = sfi(X,8,0)
```

```
A =
```

```
-128 -128 -128 127
-128 127 127 -128
 127 127 127 127
 127 127 -128 -128
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 8
      FractionLength: 0
```

```
m = size(A,1)
```

```
m =
```

```
4
```

The necessary growth is 1.6468 times the square-root of the number of rows of A.

```
bit_growth = ceil(log2(cordic_growth_constant*sqrt(m)))
```

```
bit_growth =
```

```
2
```

Initialize R with the same values as A, and allow for bit growth like you did in the previous section.

```
R = sfi(A, get(A, 'WordLength')+bit_growth, get(A, 'FractionLength'))
```

```
R =
```

```
-128 -128 -128 127
-128 127 127 -128
 127 127 127 127
 127 127 -128 -128
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 10
      FractionLength: 0
```

Compute the QR factorization, overwriting R.

```
niter = get(R, 'WordLength') - 1;
[Q,R] = cordicqr(R, niter)
```

Q =

```
-0.5039  -0.2930  -0.4062  -0.6914
-0.5039   0.8750   0.0039   0.0078
 0.5000   0.2930   0.3984  -0.7148
 0.4922   0.2930  -0.8203   0.0039
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 10
      FractionLength: 8
```

R =

```
257  126  -1  -1
  0  225 151 -148
  0   0 211  104
  0   0   0 -180
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 10
      FractionLength: 0
```

Notice that R is returned with integer values because you left the fraction length of R at 0, the same as the fraction length of A.

The scaling of the least-significant bit (LSB) of A is 1, and you can see that the error is proportional to the LSB.

```
err = double(Q)*double(R)-double(A)
```

err =

```
-1.5039  -1.4102  -1.4531  -0.9336
-1.5039   6.3828   6.4531  -1.9961
 1.5000   1.9180   0.8086  -0.7500
-0.5078   0.9336  -1.3398  -1.8672
```

You can increase the precision in the QR factorization by increasing the fraction length. In this example, you needed 10 bits for the integer part (8 bits to start with, plus 2 bits growth), so when you increase the fraction length you still need to keep the 10 bits in the integer part. For example, you can increase the word length to 32 and set the fraction length to 22, which leaves 10 bits in the integer part.

```
R = sfi(A, 32, 22)
```

R =

```
-128 -128 -128 127
-128 127 127 -128
127 127 127 127
127 127 -128 -128
```

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 32
FractionLength: 22
```

```
niter = get(R, 'WordLength') - 1;
[Q,R] = cordicqr(R, niter)
```

Q =

```
-0.5020 -0.2913 -0.4088 -0.7043
-0.5020 0.8649 0.0000 0.0000
0.4980 0.2890 0.4056 -0.7099
0.4980 0.2890 -0.8176 0.0000
```

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 32
FractionLength: 30
```

R =

```
255.0020 127.0029 0.0039 0.0039
0 220.5476 146.8413 -147.9930
0 0 208.4793 104.2429
0 0 0 -179.6037
```

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 32
FractionLength: 22
```

Now you can see fractional parts in R, and $Q \cdot R - A$ is small.

```
err = double(Q)*double(R)-double(A)
```

err =

```
1.0e-05 *
-0.1234 -0.0014 -0.0845 0.0267
-0.1234 0.2574 0.1260 -0.1094
0.0720 0.0289 -0.0400 -0.0684
0.0957 0.0818 -0.1034 0.0095
```

The number of bits you choose for fraction length will depend on the precision requirements for your particular algorithm.

Picking Default Number of Iterations

The number of iterations is dependent on the desired precision, but limited by the word length of A. With each iteration, the values are right-shifted one bit. After the last bit gets shifted off and the value becomes 0, then there is no additional value in continuing to rotate. Hence, the most precision will be attained by choosing `niter` to be one less than the word length.

For floating-point, the number of iterations is bounded by the size of the mantissa. In double, 52 iterations is the most you can do to continue adding to something with the same exponent. In single, it is 23. See the reference page for `eps` for more information about floating-point accuracy.

Thus, we can make our code more usable by not requiring the number of iterations to be input, and assuming that we want the most precision possible by changing `cordicqr` to use this default for `niter`.

```
function [Q,R] = cordicqr(A,varargin)
    if nargin>=2 && ~isempty(varargin{1})
        niter = varargin{1};
    elseif isa(A,'double') || isfi(A) && isdouble(A)
        niter = 52;
    elseif isa(A,'single') || isfi(A) && issingle(A)
        niter = single(23);
    elseif isfi(A)
        niter = int32(get(A,'WordLength') - 1);
    else
        assert(0,'First input must be double, single, or fi.');
```

A disadvantage of doing this is that this makes a section of our code dependent on data type. However, an advantage is that the function is much more convenient to use because you don't have to specify `niter` if you don't want to, and the main algorithm is still data-type independent. Similar to picking an optimal output type for Q, you can do this kind of input parsing in the beginning of a function and leave the main algorithm data-type independent.

Here is an example from a previous section, without needing to specify an optimal `niter`.

```
A = [7   -7   7   7
      7   7  -7   7
      7  -7  -7  -7
      7   7   7  -7];
```

```
[Q,R] = cordicqr(A)
```

```
Q =
```

```
    0.5000   -0.5000    0.5000    0.5000
    0.5000    0.5000   -0.5000    0.5000
    0.5000   -0.5000   -0.5000   -0.5000
    0.5000    0.5000    0.5000   -0.5000
```

```
R =
```

```
  14.0000    0.0000   -0.0000   -0.0000
         0   14.0000   -0.0000    0.0000
         0         0   14.0000    0.0000
```



```

0      0      0      14.0000

```

Example: QR Factorization Not Unique

When you compare the results from `cordicqr` and the `qr` function in MATLAB, you will notice that the QR factorization is not unique. It is only important that Q is orthogonal, R is upper triangular, and $Q^*R - A$ is small.

Here is a simple example that shows the difference.

```

m = 3;
A = ones(m)

```

A =

```

1      1      1
1      1      1
1      1      1

```

The built-in QR function in MATLAB uses a different algorithm and produces:

```
[Q0,R0] = qr(A)
```

Q0 =

```

-0.5774  -0.5774  -0.5774
-0.5774   0.7887  -0.2113
-0.5774  -0.2113   0.7887

```

R0 =

```

-1.7321  -1.7321  -1.7321
      0      0      0
      0      0      0

```

And the `cordicqr` function produces:

```
[Q,R] = cordicqr(A)
```

Q =

```

0.5774   0.7495   0.3240
0.5774  -0.6553   0.4871
0.5774  -0.0942  -0.8110

```

R =

```

1.7321   1.7321   1.7321
      0   0.0000   0.0000
      0      0  -0.0000

```

Notice that the elements of Q from function `cordicqr` are different from Q0 from built-in QR. However, both results satisfy the requirement that Q is orthogonal:

Q0*Q0'

ans =

```

1.0000    0.0000    0
0.0000    1.0000    0
0         0         1.0000

```

Q*Q'

ans =

```

1.0000    0.0000    0.0000
0.0000    1.0000   -0.0000
0.0000   -0.0000    1.0000

```

And they both satisfy the requirement that $Q^*R - A$ is small:

Q0*R0 - A

ans =

```

1.0e-15 *
-0.1110  -0.1110  -0.1110
-0.1110  -0.1110  -0.1110
-0.1110  -0.1110  -0.1110

```

Q*R - A

ans =

```

1.0e-15 *
-0.2220    0.2220    0.2220
0.4441         0         0
0.2220    0.2220    0.2220

```

Solving Systems of Equations Without Forming Q

Given matrices A and B, you can use the QR factorization to solve for X in the following equation:

$$A^*X = B.$$

If A has more rows than columns, then X will be the least-squares solution. If X and B have more than one column, then several solutions can be computed at the same time. If $A = Q^*R$ is the QR factorization of A, then the solution can be computed by back-solving

$$R^*X = C$$

where $C = Q' * B$. Instead of forming Q and multiplying to get $C = Q' * B$, it is more efficient to compute C directly. You can compute C directly by applying the rotations to the rows of B instead of to the columns of an identity matrix. The new algorithm is formed by the small modification of initializing $C = B$, and operating along the rows of C instead of the columns of Q .

```
function [R,C] = cordicrc(A,B,niter)
    Kn = inverse_cordic_growth_constant(niter);
    [m,n] = size(A);
    R = A;
    C = B;
    for j=1:n
        for i=j+1:m
            [R(j,j:end),R(i,j:end),C(j,:),C(i,:)] = ...
                cordicgivens(R(j,j:end),R(i,j:end),C(j,:),C(i,:),niter,Kn);
        end
    end
end
```

You can verify the algorithm with this example. Let A be a random 3-by-3 matrix, and B be a random 3-by-2 matrix.

```
A = [-0.8201    0.3573   -0.0100
     -0.7766   -0.0096   -0.7048
     -0.7274   -0.6206   -0.8901];
```

```
B = [-0.9286    0.3575
      0.6983    0.5155
      0.8680    0.4863];
```

Compute the QR factorization of A .

```
[Q,R] = cordicqr(A)
```

$Q =$

```
-0.6105    0.6133    0.5012
-0.5781    0.0876   -0.8113
-0.5415   -0.7850    0.3011
```

$R =$

```
1.3434    0.1235    0.8955
  0        0.7054    0.6309
  0         0         0.2988
```

Compute $C = Q' * B$ directly.

```
[R,C] = cordicrc(A,B)
```

$R =$

```
1.3434    0.1235    0.8955
  0        0.7054    0.6309
  0         0         0.2988
```

```
C =  
-0.3068  -0.7795  
-1.1897  -0.1173  
-0.7706  -0.0926
```

Subtract, and you will see that the error difference is on the order of roundoff.

```
Q'*B - C
```

```
ans =  
1.0e-15 *  
-0.0555  0.3331  
  0      0  
0.1110  0.2914
```

Now try the example in fixed-point. Declare A and B to be fixed-point types.

```
A = sfi(A)
```

```
A =  
-0.8201  0.3573  -0.0100  
-0.7766  -0.0096  -0.7048  
-0.7274  -0.6206  -0.8901  
  
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 16  
    FractionLength: 15
```

```
B = sfi(B)
```

```
B =  
-0.9286  0.3575  
0.6983  0.5155  
0.8680  0.4863  
  
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 16  
    FractionLength: 15
```

The necessary growth is 1.6468 times the square-root of the number of rows of A.

```
bit_growth = ceil(log2(cordic_growth_constant*sqrt(m)))
```

```
bit_growth =
```

2

Initialize R with the same values as A, and allow for bit growth.

```
R = sfi(A, get(A, 'WordLength')+bit_growth, get(A, 'FractionLength'))
```

R =

```
-0.8201    0.3573   -0.0100
-0.7766   -0.0096   -0.7048
-0.7274   -0.6206   -0.8901
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 18
    FractionLength: 15
```

The growth in C is the same as R, so initialize C and allow for bit growth the same way.

```
C = sfi(B, get(B, 'WordLength')+bit_growth, get(B, 'FractionLength'))
```

C =

```
-0.9286    0.3575
 0.6983    0.5155
 0.8680    0.4863
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 18
    FractionLength: 15
```

Compute $C = Q^*B$ directly, overwriting R and C.

```
[R,C] = cordicrc(R,C)
```

R =

```
 1.3435    0.1233    0.8954
 0          0.7055    0.6308
 0          0          0.2988
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 18
    FractionLength: 15
```

C =

```
-0.3068   -0.7796
-1.1898   -0.1175
-0.7706   -0.0926
```

```
    DataTypeMode: Fixed-point: binary point scaling
```

```
Signedness: Signed
WordLength: 18
FractionLength: 15
```

An interesting use of this algorithm is that if you initialize B to be the identity matrix, then output argument C is Q'. You may want to use this feature to have more control over the data type of Q. For example,

```
A = [-0.8201    0.3573   -0.0100
      -0.7766   -0.0096   -0.7048
      -0.7274   -0.6206   -0.8901];
B = eye(size(A,1))
```

```
B =
```

```
    1    0    0
    0    1    0
    0    0    1
```

```
[R,C] = cordicrc(A,B)
```

```
R =
```

```
    1.3434    0.1235    0.8955
         0    0.7054    0.6309
         0         0    0.2988
```

```
C =
```

```
   -0.6105   -0.5781   -0.5415
    0.6133    0.0876   -0.7850
    0.5012   -0.8113    0.3011
```

Then C is orthogonal

```
C'*C
```

```
ans =
```

```
    1.0000    0.0000    0.0000
    0.0000    1.0000   -0.0000
    0.0000   -0.0000    1.0000
```

```
and R = C*A
```

```
R - C*A
```

```
ans =
```

```
    1.0e-15 *
```

```

0.6661   -0.0139   -0.1110
0.5551   -0.2220    0.6661
-0.2220   -0.1110    0.2776

```

Links to the Documentation

Fixed-Point Designer™

- `bitsra` Bit shift right arithmetic
- `fi` Construct fixed-point numeric object
- `fimath` Construct `fimath` object
- `fipref` Construct `fipref` object
- `get` Property values of object
- `globalfimath` Configure global `fimath` and return handle object
- `isfi` Determine whether variable is `fi` object
- `sfi` Construct signed fixed-point numeric object
- `upperbound` Upper bound of range of `fi` object
- `fiaccel` Accelerate fixed-point code

MATLAB

- `bitshift` Shift bits specified number of places
- `ceil` Round toward positive infinity
- `double` Convert to double precision floating point
- `eps` Floating-point relative accuracy
- `eye` Identity matrix
- `log2` Base 2 logarithm and dissect floating-point numbers into exponent and mantissa
- `prod` Product of array elements
- `qr` Orthogonal-triangular factorization
- `repmat` Replicate and tile array
- `single` Convert to single precision floating point
- `size` Array dimensions
- `sqrt` Square root
- `subsasgn` Subscripted assignment

Functions Used in this Example

These are the MATLAB functions used in this example.

CORDICQR computes the QR factorization using CORDIC.

- `[Q,R] = cordicqr(A)` chooses the number of CORDIC iterations based on the type of `A`.
- `[Q,R] = cordicqr(A,niter)` uses `niter` number of CORDIC iterations.

CORDICRC computes `R` from the QR factorization of `A`, and also returns `C = Q'*B` without computing `Q`.

- `[R,C] = cordicrc(A,B)` chooses the number of CORDIC iterations based on the type of A.
- `[R,C] = cordicrc(A,B,niter)` uses `niter` number of CORDIC iterations.

CORDIC_GROWTH_CONSTANT returns the CORDIC growth constant.

- `cordic_growth = cordic_growth_constant(niter)` returns the CORDIC growth constant as a function of the number of CORDIC iterations, `niter`.

GIVENSQR computes the QR factorization using standard Givens rotations.

- `[Q,R] = givensqr(A)`, where A is M-by-N, produces an M-by-N upper triangular matrix R and an M-by-M orthogonal matrix Q so that $A = Q \cdot R$.

CORDICQR_MAKEPLOTS makes the plots in this example by executing the following from the MATLAB command line.

```
load A_3_by_3_for_cordicqr_demo.mat
niter=32;
[Q,R] = cordicqr_makeplots(A,niter)
```

References

- 1 Ray Andraka, "A survey of CORDIC algorithms for FPGA based computers," 1998, ACM 0-89791-978-5/98/01.
- 2 Anthony J Cox and Nicholas J Higham, "Stability of Householder QR factorization for weighted least squares problems," in Numerical Analysis, 1997, Proceedings of the 17th Dundee Conference, Griffiths DF, Higham DJ, Watson GA (eds). Addison-Wesley, Longman: Harlow, Essex, U.K., 1998; 57-73.
- 3 Gene H. Golub and Charles F. Van Loan, *Matrix Computations*, 3rd ed, Johns Hopkins University Press, 1996, section 5.2.3 Givens QR Methods.
- 4 Daniel V. Rabinkin, William Song, M. Michael Vai, and Huy T. Nguyen, "Adaptive array beamforming with fixed-point arithmetic matrix inversion using Givens rotations," Proceedings of Society of Photo-Optical Instrumentation Engineers (SPIE) -- Volume 4474 Advanced Signal Processing Algorithms, Architectures, and Implementations XI, Franklin T. Luk, Editor, November 2001, pp. 294--305.
- 5 Jack E. Volder, "The CORDIC Trigonometric Computing Technique," Institute of Radio Engineers (IRE) Transactions on Electronic Computers, September, 1959, pp. 330-334.
- 6 Musheng Wei and Qiaohua Liu, "On growth factors of the modified Gram-Schmidt algorithm," Numerical Linear Algebra with Applications, Vol. 15, issue 7, September 2008, pp. 621-636.

Cleanup

```
fipref(originalFipref);
globalfimath(originalGlobalFimath);
close all
set(0, 'format', originalFormat);
%#ok<*MNEFF,*NASGU,*NOPTS,*ASGLU>
```


Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using QR Decomposition

This example shows how to implement a hardware-efficient least-squares solution to the real-valued matrix equation $AX=B$ using the Real Partial-Systolic Matrix Solve Using QR Decomposition block.

Define Matrix Dimensions

Specify the number of rows in matrices A and B, the number of columns in matrix A, and the number of columns in matrix B.

```
m = 300; % Number of rows in matrices A and B
n = 10;  % Number of columns in matrix A
p = 1;  % Number of columns in matrix B
```

Generate Random Least-Squares Matrices

For this example, use the helper function `realRandomLeastSquaresMatrices` to generate random matrices A and B for the least-squares problem $AX=B$. The matrices are generated such that the elements of A and B are between -1 and +1, and A is full rank.

```
rng('default')
[A,B] = fixed.example.realRandomLeastSquaresMatrices(m,n,p);
```

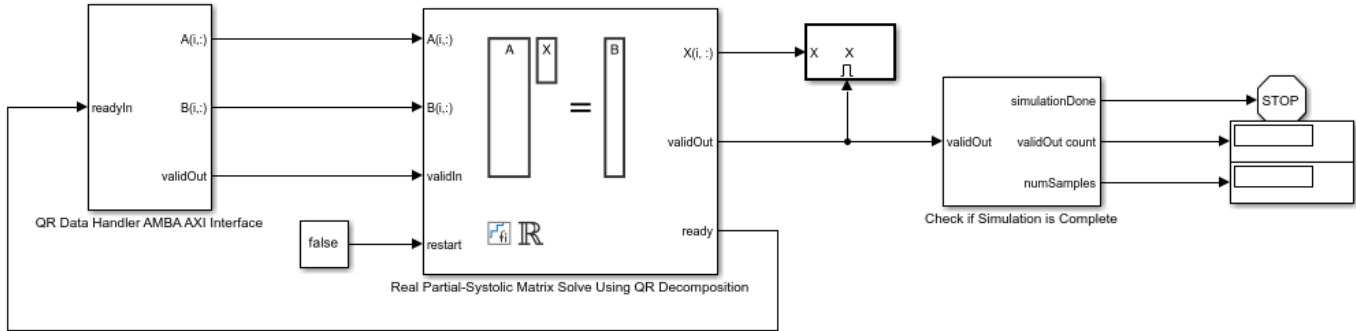
Select Fixed-Point Data Types

Use the helper function `realQRMatrixSolveFixedpointTypes` to select fixed-point data types for input matrices A and B, and output X such that there is a low probability of overflow during the computation.

```
max_abs_A = 1;      % Upper bound on max(abs(A(:)))
max_abs_B = 1;      % Upper bound on max(abs(B(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.realQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,precisionBits);
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);
```

Open the Model

```
model = 'RealPartialSystolicQRMatrixSolveModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

The Data Handler subsystem in this model takes real matrices A and B as inputs. It sends rows of A and B to QR block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set a delay between the feeding in rows of A and B in the Data Handler to emulate the processing time of the upstream block. `validIn` remains high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Real Partial-Systolic Matrix Solve Using QR Decomposition block.

```
numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A and B
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
    'regularizationParameter', 0, ...
    'numSamples', numSamples, 'rowDelay', rowDelay, 'OutputType', OutputType);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Real Partial-Systolic Matrix Solve Using QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of X are output in the order they are computed, last row first, so you must reconstruct the data to interpret the results. To reconstruct the matrix X from the output data, use the helper function `matrixSolveModelOutputToArray`.

```
X = fixed.example.matrixSolveModelOutputToArray(out.X, n, p, numSamples);
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Real Partial-Systolic Matrix Solve Using QR Decomposition block, compute the relative error.

```
relative_error = norm(double(A*X - B))/norm(double(B)) %#ok<NOPTS>
```

```
relative_error =
```

2.3952e-05

Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition

This example shows how to implement a hardware-efficient solution to the real-valued matrix equation $A'AX=B$ using the Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition block.

Forward and Backward Substitution

When an upper triangular factor is ready, then forward and backward substitution are computed with the current input B to produce output X .

$$X = R_k \setminus (R'_k \setminus B)$$

Define Matrix Dimensions

Specify the number of rows in matrix A , the number of columns in matrix A and rows in B , and the number of columns in matrix B .

```
m = 30; % Number of rows in A
n = 10; % Number of columns in A and rows in B
p = 1; % Number of columns in B
numInputs = 3; % Number of A and B matrices
```

Generate Matrices

For this example, use the helper function `realRandomQlessQRMatrices` to generate random matrices A and B for the problem $A'AX=B$. The matrices are generated such that the elements of A and B are between -1 and $+1$, and A is full rank.

```
rng('default')
[A,B] = fixed.example.realRandomQlessQRMatrices(m,n,p);
if numInputs > 1
    for i = 2:numInputs
        [Atemp,Btemp] = fixed.example.realRandomQlessQRMatrices(m,n,p);
        A = cat(3,A,Atemp);
        B = cat(3,B,Btemp);
    end
end
```

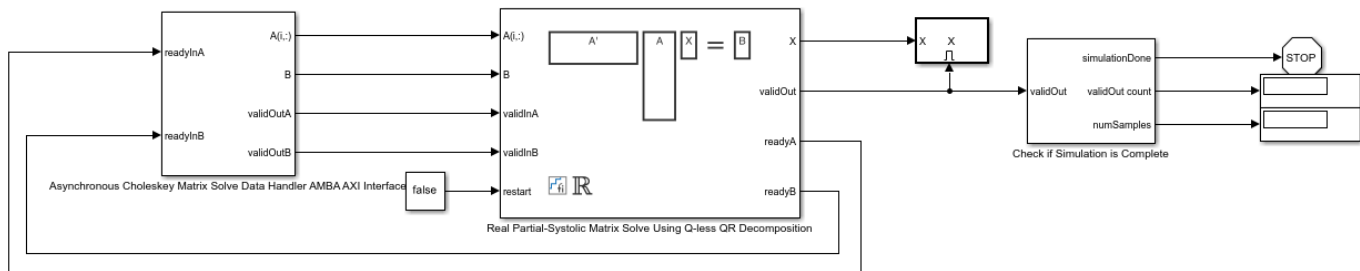
Select Fixed-Point Data Types

Use the helper function `realQlessQRMatrixSolveFixedpointTypes` to select fixed-point data types for input matrices A and B , and output X such that there is a low probability of overflow during the computation.

```
max_abs_A = 1; % Upper bound on max(abs(A(:)))
max_abs_B = 1; % Upper bound on max(abs(B(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.realQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,precisionBits);
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);
```

Open the Model

```
model = 'RealPartialSystolicQlessQRMatrixSolveModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

AMBA AXI Handshaking Process

The Data Handler subsystem in this model takes real matrices A and B as inputs. It sends rows of A and full matrix of B to the QR Decomposition block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set delays for the feeding in rows of A and the feeding in of B matrices in the Data Handler to emulate the processing time of the upstream block. `validInA` and `validInB` remain high when `aDelay` and `bDelay` are set to 0 because this indicates the Data Handler always has data available. When all matrices A and B are sent, the Data Handler loops back to the first A and B matrices.

Asynchronous Matrix Solver

This block operates asynchronously. First, Q-less QR decomposition is performed on the input A matrix and the resulting R matrix is put into a buffer. Then, the Forward Backward Substitute block uses the input B matrix and the buffered R matrix to compute $R'RX = B$. Because the R and B matrices are stored separately in buffers, the upstream Q-less QR decomposition block and the downstream Forward Backward Substitute block can run independently. The Forward Backward Substitute block starts processing when the first R and B matrices are available. Then it runs continuously using the latest buffered R and B matrices, regardless of the status of the Q-less QR Decomposition block. For example, if the upstream block stops providing A and B matrices, the Forward Backward Substitute block continues to generate the same output using the last pair of R and B matrices.

The Data Handler sends A and B matrices to the QR decomposition block iteratively. After sending out the last A matrix, the Data Handler resets its internal counter and sends out first A matrix. The B matrix is handled in a similar fashion.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition block.

```
numOutputs = 10; % Number of recorded outputs
aDelay = 1; % Delay of clock cycles between feeding in rows of A
bDelay = 1; % Delay of clock cycles between feeding in B matrices
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
```

```
'regularizationParameter',0,...
'aDelay',aDelay,'bDelay',bDelay,...
'numOutputs',numOutputs,'OutputType',OutputType);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition block outputs matrix X at each time step. When a valid result matrix is output, the block sets `validOut` to true.

```
X = out.X;
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition block, compute the relative error. Choose the last output of the simulation.

```
X = double(X(:,:,end));
```

Synchronize the last output X with the input by finding the inputs A and B that produced it.

```
A = double(A);
B = double(B);
relative_errors = zeros(size(A,3),size(B,3));
for k = 1:size(A,3)
    for g = 1:size(B,3)
        relative_errors(k,g) = norm(A(:,:,k)'*A(:,:,k)*X - B(:,:,g))/norm(B(:,:,g));
    end
end
```

```
[AUsed,Bused] = find(relative_errors==min(relative_errors,[],'all')) %#ok<NOPTS>
```

```
relative_error = norm(double(A(:,:,AUsed)'*A(:,:,AUsed)*X - B(:,:,Bused)))/norm(double(B(:,:,Bused))
```

```
AUsed =
```

```
3
```

```
Bused =
```

```
2
```

```
relative_error =
```

```
8.1385e-05
```

Implement Hardware-Efficient Real Burst Matrix Solve Using QR Decomposition

This example shows how to implement a hardware-efficient least-squares solution to the real-valued matrix equation $AX=B$ using the Real Burst Matrix Solve Using QR Decomposition block.

Define Matrix Dimensions

Specify the number of rows in matrices A and B, the number of columns in matrix A, and the number of columns in matrix B.

```
m = 50; % Number of rows in matrices A and B
n = 10; % Number of columns in matrix A
p = 1; % Number of columns in matrix B
```

Generate Random Least-Squares Matrices

For this example, use the helper function `realRandomLeastSquaresMatrices` to generate random matrices A and B for the least-squares problem $AX=B$. The matrices are generated such that the elements of A and B are between -1 and +1, and A has rank r.

```
rng('default')
r = 3; % Rank of A
[A,B] = fixed.example.realRandomLeastSquaresMatrices(m,n,p,r);
```

Select Fixed-Point Data Types

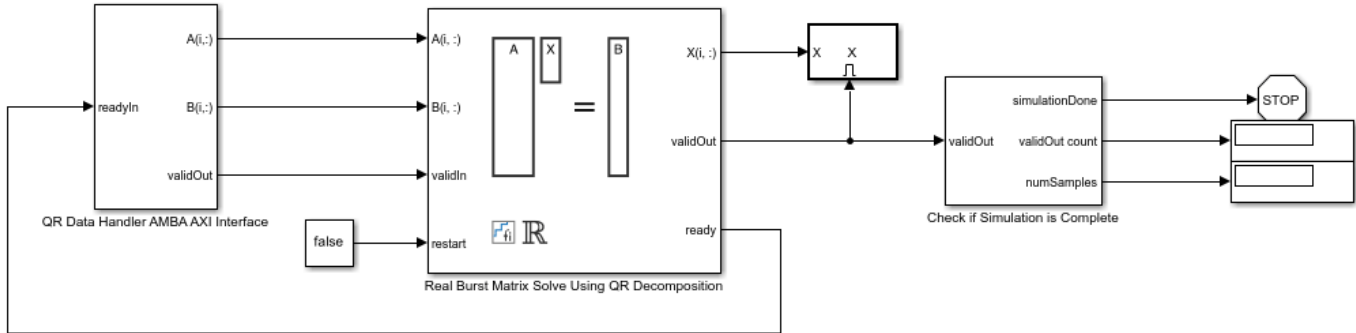
Use the helper function `realQRMatrixSolveFixedpointTypes` to select fixed-point data types for input matrices A and B, and output X such that there is a low probability of overflow during the computation.

```
max_abs_A = 1; % Upper bound on max(abs(A(:)))
max_abs_B = 1; % Upper bound on max(abs(B(:)))
precisionBits = 24; % Number of bits of precision

T = fixed.realQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,precisionBits);
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);
```

Open the Model

```
model = 'RealBurstQRMatrixSolveModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

The Data Handler subsystem in this model takes real matrices A and B as inputs. It sends rows of A and B to QR block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set a delay between the feeding in rows of A and B in the Data Handler to emulate the processing time of the upstream block. `validIn` remains high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Real Burst Matrix Solve Using QR Decomposition block.

```
numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A and B
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
    'regularizationParameter', 0, ...
    'numSamples', numSamples, 'rowDelay', rowDelay, 'OutputType', OutputType);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Real Burst Matrix Solve Using QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of X are output in the order they are computed, last row first, so you must reconstruct the data to interpret the results. To reconstruct the matrix X from the output data, use the helper function `matrixSolveModelOutputToArray`.

```
X = fixed.example.matrixSolveModelOutputToArray(out.X, n, p, numSamples);
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Real Burst Matrix Solve Using QR Decomposition block, compute the relative error.

```
relative_error = norm(double(A*X - B))/norm(double(B)) %#ok<NOPTS>
```

```
relative_error =
```


3.4664e-06

Implement Hardware-Efficient Real Burst Matrix Solve Using Q-less QR Decomposition

This example shows how to implement a hardware-efficient solution to the real-valued matrix equation $A'AX=B$ using the Real Burst Matrix Solve Using Q-less QR Decomposition block.

Define Matrix Dimensions

Specify the number of rows in matrix A, the number of columns in matrix A and rows in B, and the number of columns in matrix B.

```
m = 100; % Number of rows in A
n = 10;  % Number of columns in A and rows in B
p = 1;  % Number of columns in B
```

Generate Matrices

For this example, use the helper function `realRandomQlessQRMatrices` to generate random matrices A and B for the problem $A'AX=B$. The matrices are generated such that the elements of A and B are between -1 and +1, and A is full rank.

```
rng('default')
[A,B] = fixed.example.realRandomQlessQRMatrices(m,n,p);
```

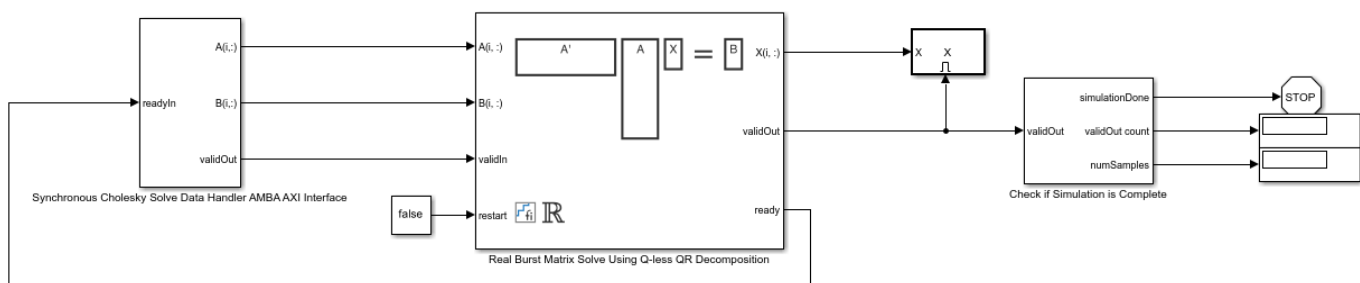
Select Fixed-Point Data Types

Use the helper function `realQlessQRMatrixSolveFixedpointTypes` to select fixed-point data types for input matrices A and B, and output X such that there is a low probability of overflow during the computation.

```
max_abs_A = 1;      % Upper bound on max(abs(A(:)))
max_abs_B = 1;      % Upper bound on max(abs(B(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.realQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,precisionBits);
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);
```

Open the Model

```
model = 'RealBurstQlessQRMatrixSolveModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

The Data Handler subsystem in this model takes real matrices A and B as inputs. It sends rows of A and B to QR block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set a delay between the feeding in rows of A and B in the Data Handler to emulate the processing time of the upstream block. `validIn` remains high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Real Burst Matrix Solve Using Q-less QR Decomposition block.

```
numSamples = 1; % Number of samples
rowDelay = 1; % Delay of clock cycles between feeding in rows of A and B
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
    'regularizationParameter', 0, ...
    'numSamples', numSamples, 'rowDelay', rowDelay, 'OutputType', OutputType);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Real Burst Matrix Solve Using Q-less QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of X are output in the order they are computed, last row first, so you must reconstruct the data to interpret the results. To reconstruct the matrix X from the output data, use the helper function `matrixSolveModelOutputToArray`.

```
X = fixed.example.matrixSolveModelOutputToArray(out.X, n, p);
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Real Burst Matrix Solve Using Q-less QR Decomposition block, compute the relative error.

```
relative_error = norm(double(A'*A*X - B))/norm(double(B)) %#ok<NOPTS>
```

```
relative_error =
    5.5603e-04
```

Implement Hardware-Efficient Real Partial-Systolic QR Decomposition

This example shows how to implement a hardware-efficient QR decomposition using the Real Partial-Systolic QR Decomposition block.

Economy Size QR Decomposition

The Real Partial-Systolic QR Decomposition block performs the first step of solving the least-squares matrix equation $AX = B$ which transforms A in-place to R and B in-place to $C = Q'B$, then solves the transformed system $RX = C$, where QR is the orthogonal-triangular decomposition of A.

To compute the stand-alone QR decomposition, this example sets B to be the identity matrix so that the output of the Real Partial-Systolic QR Decomposition block is the upper-triangular R and $C = Q'$.

Define Matrix Dimensions

Specify the number of rows in matrices A and B, the number of columns in matrix A, and the number of columns in matrix B. This example sets B to be the identity matrix the same size as the number of rows of A.

```
m = 10; % Number of rows in matrices A and B
n = 3;  % Number of columns in matrix A
p = m;  % Number of columns in matrix B
```

Generate Matrices A and B

Use the helper function `realUniformRandomArray` to generate a random matrix A such that the elements of A are between -1 and +1, and A is full rank. Matrix B is the identity matrix.

```
rng('default')
A = fixed.example.realUniformRandomArray(-1,1,m,n);
B = eye(m);
```

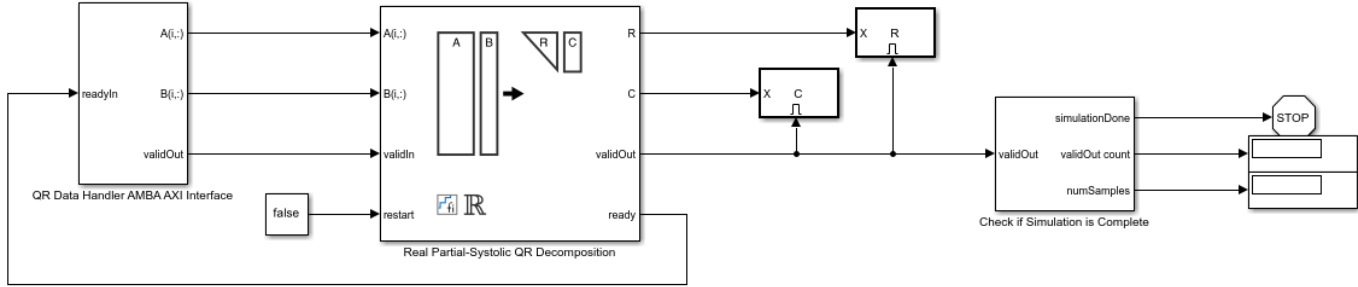
Select Fixed-Point Data Types

Use the helper function `qrFixedpointTypes` to select fixed-point data types for matrices A and B that guarantee no overflow will occur in the transformation of A in-place to R and B in-place to $C = Q'B$.

```
max_abs_A = 1; % Upper bound on max(abs(A(:)))
max_abs_B = 1; % Upper bound on max(abs(B(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.qrFixedpointTypes(m,max_abs_A,max_abs_B,precisionBits);
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
```

Open the Model

```
model = 'RealPartialSystolicQRModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

The Data Handler subsystem in this model takes real matrices A and B as inputs. It sends rows of A and B to QR block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set a delay between the feeding in rows of A and B in the Data Handler to emulate the processing time of the upstream block. `validIn` remains high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Real Partial-Systolic QR Decomposition block.

```
numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A and B
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
    'numSamples', numSamples, 'rowDelay', rowDelay);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Real Partial-Systolic QR Decomposition block outputs matrices R and C at each time step. When valid result matrices are output, the block sets `validOut` to true.

```
R = out.R;
C = out.C;
```

Extract the Economy-Size Q

The block computes $C = Q'B$. In this example, B is the identity matrix, so $Q = C'$ is the economy-size orthogonal factor of the QR decomposition.

```
Q = C';
```

Verify that Q is Orthogonal and R is Upper-Triangular

Q is orthogonal, so $Q'Q$ is the identity matrix within roundoff.

```
I = Q'*Q
```

```
I =  
  
    1.0000    -0.0000    -0.0000  
   -0.0000     1.0000    -0.0000  
   -0.0000    -0.0000     1.0000  
  
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 62  
    FractionLength: 48
```

R is an upper-triangular matrix.

R

```
R =  
  
    2.2180    0.8559   -0.5607  
     0      2.0578   -0.4017  
     0         0      1.7117  
  
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 29  
    FractionLength: 24
```

```
isequal(R, triu(R))
```

```
ans =  
  
    logical  
  
     1
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Real Partial-Systolic QR Decomposition block, compute the relative error.

```
relative_error = norm(double(Q*R - A))/norm(double(A))
```

```
relative_error =  
  
    1.5886e-06
```

Suppress mlint warnings.

```
 %#ok<*NOPTS>
```

Implement Hardware-Efficient Real Partial-Systolic Q-less QR Decomposition

This example shows how to implement a hardware-efficient Q-less QR decomposition using the Real Partial-Systolic Q-less QR Decomposition block.

Economy Size Q-less QR Decomposition

The Real Partial-Systolic Q-less QR Decomposition block performs the first step of solving the matrix equation $A'AX = B$ which transforms A in-place to upper-triangular R, then solves the transformed system $R'RX = B$, where $R'R = A'A$.

Define Matrix Dimensions

Specify the number of rows and columns in matrix A.

```
m = 5; % Number of rows in matrix A
n = 3; % Number of columns in matrix A
```

Generate Matrix A

Use the helper function `realUniformRandomArray` to generate a random matrix A such that the elements of A are between -1 and +1, and A is full rank.

```
rng('default')
A = fixed.example.realUniformRandomArray(-1,1,m,n);
```

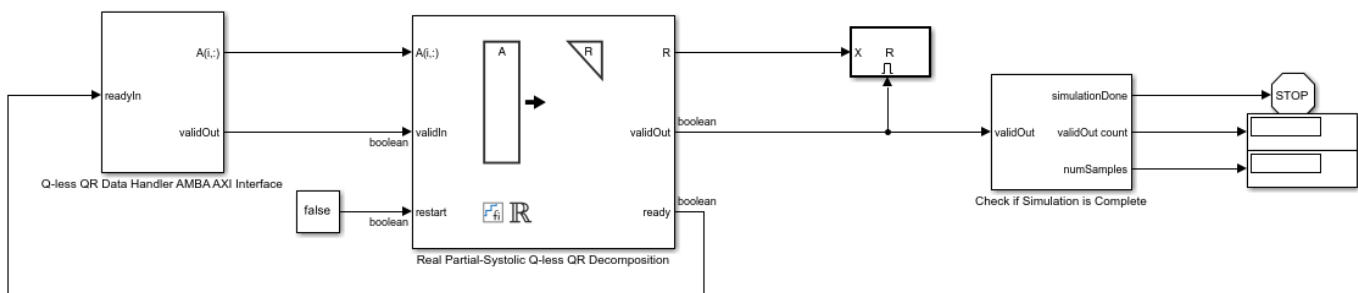
Select Fixed-Point Data Types

Use the helper function `qlessqrFixedpointTypes` to select fixed-point data types for matrix A that guarantee no overflow will occur in the transformation of A in-place to R.

```
max_abs_A = 1; % Upper bound on max(abs(A(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.qlessqrFixedpointTypes(m,max_abs_A,precisionBits);
A = cast(A,'like',T.A);
```

Open the Model

```
model = 'RealPartialSystolicQlessQRModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

AMBA AXI Handshaking Process

The Data Handler subsystem in this model takes real matrix A as input. It sends rows of A to the QR Decomposition block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set delay for the feeding in rows of A in the Data Handler to emulate the processing time of the upstream block. `validOut` signal of the Data Handler remain high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Real Partial-Systolic Q-less QR Decomposition block.

```
numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A
fixed.example.setModelWorkspace(model, 'A', A, 'm', m, 'n', n, ...
    'numSamples', numSamples, 'rowDelay', rowDelay);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Real Partial-Systolic QR Decomposition block outputs matrix R at each time step. When a valid result matrix is output, the block sets `validOut` to true.

```
R = out.R;
```

R is an upper-triangular matrix.

```
R
```

```
R =
```

```
    1.5379    0.0432   -0.1395
         0    1.5978    0.4742
         0         0    1.5192
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 28
    FractionLength: 24
```

```
isequal(R, triu(R))
```

```
ans =
```

```
    logical
```

```
    1
```


Verify the Accuracy of the Output

To evaluate the accuracy of the Real Partial-Systolic Q-less QR Decomposition block, compute the relative error.

```
relative_error = norm(double(R'*R - A'*A))/norm(double(A'*A))
```

```
relative_error =
```

```
8.2641e-07
```

Suppress mlint warnings.

```
 %#ok<*NOPTS>
```

Implement Hardware-Efficient Real Burst QR Decomposition

This example shows how to implement a hardware-efficient QR decomposition using the Real Burst QR Decomposition block.

Economy Size QR Decomposition

The Real Burst QR Decomposition block performs the first step of solving the least-squares matrix equation $AX = B$ which transforms A in-place to R and B in-place to $C = Q'B$, then solves the transformed system $RX = C$, where QR is the orthogonal-triangular decomposition of A.

To compute the stand-alone QR decomposition, this example sets B to be the identity matrix so that the output of the Real Burst QR Decomposition block is the upper-triangular R and $C = Q'$.

Define Matrix Dimensions

Specify the number of rows in matrices A and B, the number of columns in matrix A, and the number of columns in matrix B. This example sets B to be the identity matrix the same size as the number of rows of A.

```
m = 10; % Number of rows in matrices A and B
n = 3;  % Number of columns in matrix A
p = m; % Number of columns in matrix B
```

Generate Matrices A and B

Use the helper function `realUniformRandomArray` to generate a random matrix A such that the elements of A are between -1 and +1, and A is full rank. Matrix B is the identity matrix.

```
rng('default')
A = fixed.example.realUniformRandomArray(-1,1,m,n);
B = eye(m);
```

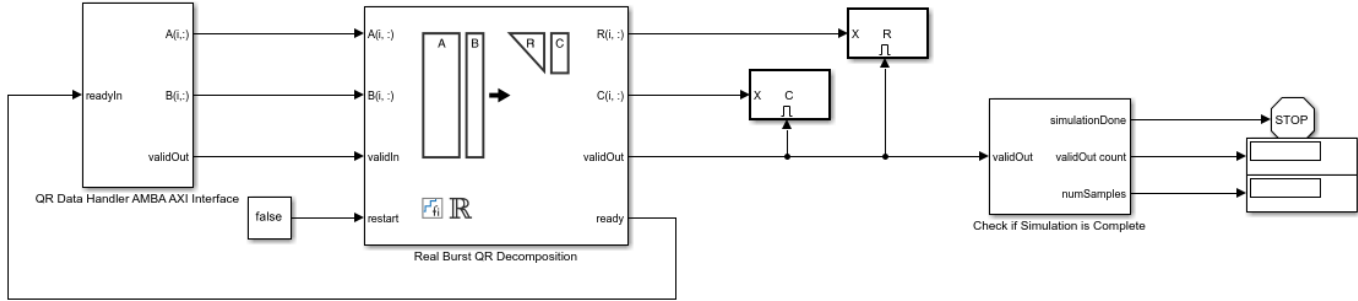
Select Fixed-Point Data Types

Use the helper function `qrFixedpointTypes` to select fixed-point data types for matrices A and B that guarantee no overflow will occur in the transformation of A in-place to R and B in-place to $C = Q'B$.

```
max_abs_A = 1; % Upper bound on max(abs(A(:)))
max_abs_B = 1; % Upper bound on max(abs(B(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.qrFixedpointTypes(m,max_abs_A,max_abs_B,precisionBits);
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
```

Open the Model

```
model = 'RealBurstQRModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

The Data Handler subsystem in this model takes real matrices A and B as inputs. It sends rows of A and B to QR block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set a delay between the feeding in rows of A and B in the Data Handler to emulate the processing time of the upstream block. `validIn` remains high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Real Burst QR Decomposition block.

```
numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A and B
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
    'numSamples', numSamples, 'rowDelay', rowDelay);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Real Burst QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of matrices R and C are output in reverse order to accommodate back-substitution, so you must reconstruct the data to interpret the results. To reconstruct the matrices R and C from the output data, use the helper function `qrModelOutputToArray`.

```
[C,R] = fixed.example.qrModelOutputToArray(out.C,out.R,m,n,p,numSamples);
```

Extract the Economy-Size Q

The block computes $C = Q'B$. In this example, B is the identity matrix, so $Q = C'$ is the economy-size orthogonal factor of the QR decomposition.

```
Q = C';
```

Verify that Q is Orthogonal and R is Upper-Triangular

Q is orothogonal, so $Q'Q$ is the identity matrix within roundoff.

```
I = Q'*Q
```

```
I =
```

```
    1.0000    -0.0000    -0.0000  
   -0.0000     1.0000    -0.0000  
   -0.0000    -0.0000     1.0000
```

```
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 62  
    FractionLength: 48
```

R is an upper-triangular matrix.

```
R
```

```
R =
```

```
    2.2180    0.8559   -0.5607  
         0    2.0578   -0.4017  
         0         0    1.7117
```

```
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 29  
    FractionLength: 24
```

```
isequal(R, triu(R))
```

```
ans =
```

```
logical
```

```
1
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Real Burst QR Decomposition block, compute the relative error.

```
relative_error = norm(double(Q*R - A))/norm(double(A))
```

```
relative_error =
```

```
1.5886e-06
```

Suppress mlint warnings.

```
 %#ok<*NOPTS>
```

Implement Hardware-Efficient Real Burst Q-less QR Decomposition

This example shows how to implement a hardware-efficient Q-less QR decomposition using the Real Burst Q-less QR Decomposition block.

Economy Size Q-less QR Decomposition

The Real Burst Q-less QR Decomposition block performs the first step of solving the matrix equation $A'AX = B$ which transforms A in-place to upper-triangular R, then solves the transformed system $R'R X = B$, where $R'R = A'A$.

Define Matrix Dimensions

Specify the number of rows and columns in matrix A.

```
m = 5; % Number of rows in matrix A
n = 3; % Number of columns in matrix A
```

Generate Matrix A

Use the helper function `realUniformRandomArray` to generate a random matrix A such that the elements of A are between -1 and +1, and A is full rank.

```
rng('default')
A = fixed.example.realUniformRandomArray(-1,1,m,n);
```

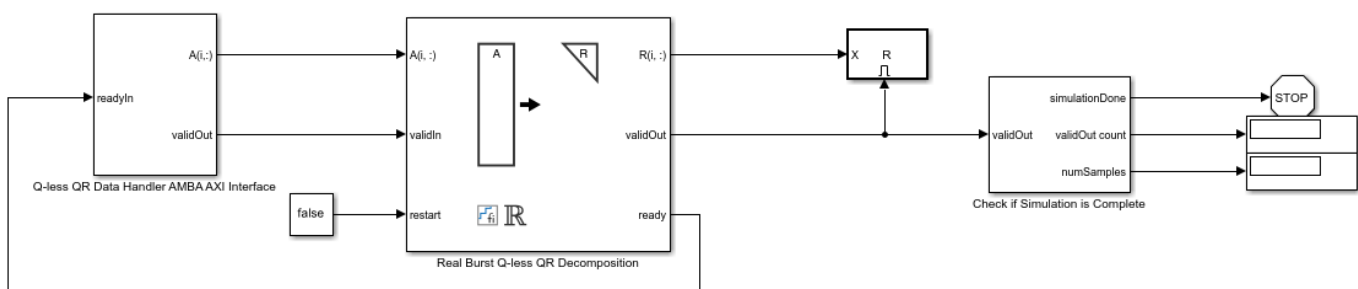
Select Fixed-Point Data Types

Use the helper function `qllessqrFixedpointTypes` to select fixed-point data types for matrix A that guarantee no overflow will occur in the transformation of A in-place to R.

```
max_abs_A = 1; % Upper bound on max(abs(A(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.qllessqrFixedpointTypes(m,max_abs_A,precisionBits);
A = cast(A,'like',T.A);
```

Open the Model

```
model = 'RealBurstQlessQRModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

AMBA AXI Handshaking Process

The Data Handler subsystem in this model takes real matrix A as input. It sends rows of A to the QR Decomposition block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set delay for the feeding in rows of A in the Data Handler to emulate the processing time of the upstream block. `validOut` signal of the Data Handler remain high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Real Burst Q-less QR Decomposition block.

```
numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A
fixed.example.setModelWorkspace(model, 'A', A, 'm', m, 'n', n, ...
    'regularizationParameter', 0, ...
    'numSamples', numSamples, 'rowDelay', rowDelay);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Real Burst Q-less QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of matrix R are output in reverse order to accommodate back-substitution, so you must reconstruct the data to interpret the results. To reconstruct the matrix R from the output data, use the helper function `qlessqrModelOutputToArray`.

```
R = fixed.example.qlessqrModelOutputToArray(out.R, m, n, numSamples);
```

R is an upper-triangular matrix.

R

R =

```
    1.5379    0.0432   -0.1395
         0    1.5978    0.4742
         0         0    1.5192
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 28
    FractionLength: 24
```

```
isequal(R, triu(R))
```

ans =

```
    logical
```

1

Verify the Accuracy of the Output

To evaluate the accuracy of the Real Burst Q-less QR Decomposition block, compute the relative error.

```
relative_error = norm(double(R'*R - A'*A))/norm(double(A'*A))
```

```
relative_error =
```

```
8.2641e-07
```

Suppress mlint warnings.

```
 %#ok<*NOPTS>
```

Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using QR Decomposition

This example shows how to implement a hardware-efficient least-squares solution to the complex-valued matrix equation $AX=B$ using the Complex Partial-Systolic Matrix Solve Using QR Decomposition block.

Define Matrix Dimensions

Specify the number of rows in matrices A and B, the number of columns in matrix A, and the number of columns in matrix B.

```
m = 300; % Number of rows in matrices A and B
n = 10;  % Number of columns in matrix A
p = 1;  % Number of columns in matrix B
```

Generate Random Least-Squares Matrices

For this example, use the helper function `complexRandomLeastSquaresMatrices` to generate random matrices A and B for the least-squares problem $AX=B$. The matrices are generated such that the real and imaginary parts of the elements of A and B are between -1 and +1, and A is full rank.

```
rng('default')
[A,B] = fixed.example.complexRandomLeastSquaresMatrices(m,n,p);
```

Select Fixed-Point Data Types

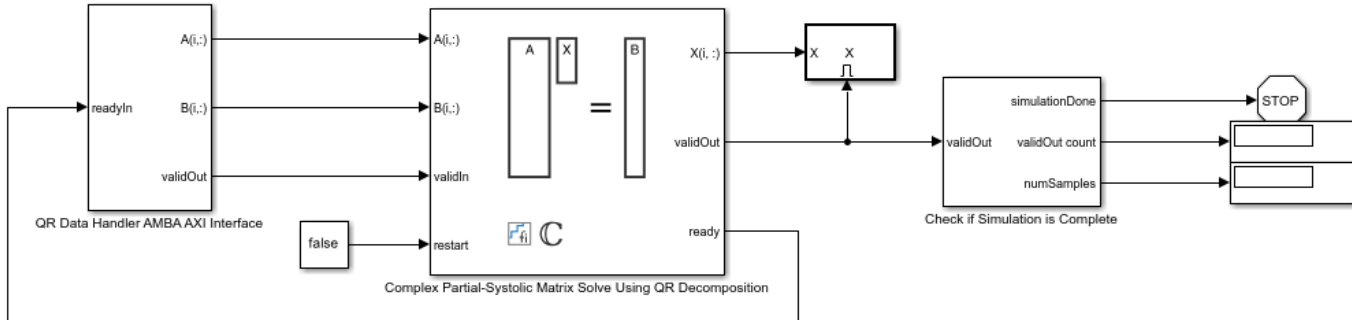
Use the helper function `complexQRMatrixSolveFixedpointTypes` to select fixed-point data types for input matrices A and B, and output X such that there is a low probability of overflow during the computation.

The real and imaginary parts of the elements of A and B are between -1 and 1, so the maximum possible absolute value of any element is $\sqrt{2}$.

```
max_abs_A = sqrt(2); % Upper bound on max(abs(A(:)))
max_abs_B = sqrt(2); % Upper bound on max(abs(B(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.complexQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,precisionBits);
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);
```

Open the Model

```
model = 'ComplexPartialSystolicQRMatrixSolveModel';
open_system(model);
```

Copyright 2021 The MathWorks, Inc.

The Data Handler subsystem in this model takes complex matrices A and B as inputs. It sends rows of A and B to QR block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set a delay between the feeding in rows of A and B in the Data Handler to emulate the processing time of the upstream block. `validIn` remains high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Complex Partial-Systolic Matrix Solve Using QR Decomposition block.

```
numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A and B
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
    'regularizationParameter', 0, ...
    'numSamples', numSamples, 'rowDelay', rowDelay, 'OutputType', OutputType);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Complex Partial-Systolic Matrix Solve Using QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of X are output in the order they are computed, last row first, so you must reconstruct the data to interpret the results. To reconstruct the matrix X from the output data, use the helper function `matrixSolveModelOutputToArray`.

```
X = fixed.example.matrixSolveModelOutputToArray(out.X, n, p, numSamples);
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Complex Partial-Systolic Matrix Solve Using QR Decomposition block, compute the relative error.

```
relative_error = norm(double(A*X - B))/norm(double(B)) %#ok<NOPTS>
```

```
relative_error =
```

3.9937e-05

Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition

This example shows how to implement a hardware-efficient solution to the complex-valued matrix equation $A'AX=B$ using the Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition block.

Forward and Backward Substitution

When an upper triangular factor is ready, then forward and backward substitution are computed with the current input B to produce output X.

$$X = R_k \setminus (R'_k \setminus B)$$

Define Matrix Dimensions

Specify the number of rows in matrix A, the number of columns in matrix A and rows in B, and the number of columns in matrix B.

```
m = 30; % Number of rows in A
n = 10; % Number of columns in A and rows in B
p = 1; % Number of columns in B
numInputs = 3; % Number of A and B matrices
```

Generate Matrices

For this example, use the helper function `complexRandomQlessQRMatrices` to generate random matrices A and B for the problem $A'AX=B$. The matrices are generated such that the real and imaginary parts of the elements of A and B are between -1 and +1, and A is full rank.

```
rng('default')
[A,B] = fixed.example.complexRandomQlessQRMatrices(m,n,p);
if numInputs > 1
    for i = 2:numInputs
        [Atemp,Btemp] = fixed.example.complexRandomQlessQRMatrices(m,n,p);
        A = cat(3,A,Atemp);
        B = cat(3,B,Btemp);
    end
end
```

Select Fixed-Point Data Types

Use the helper function `complexQlessQRMatrixSolveFixedpointTypes` to select fixed-point data types for input matrices A and B, and output X such that there is a low probability of overflow during the computation.

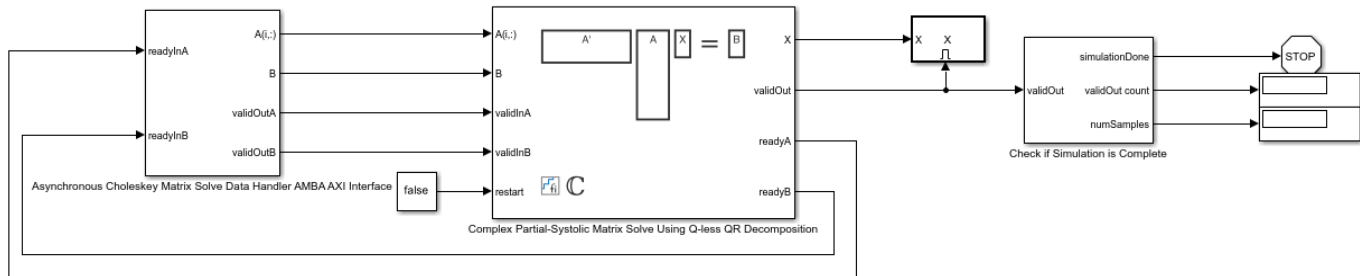
The real and imaginary parts of the elements of A and B are between -1 and 1, so the maximum possible absolute value of any element is $\sqrt{2}$.

```
max_abs_A = sqrt(2); % Upper bound on max(abs(A(:)))
max_abs_B = sqrt(2); % Upper bound on max(abs(B(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.complexQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,precisionBits);
A = cast(A,'like',T.A);
```

```
B = cast(B, 'like', T.B);
OutputType = fixed.extractNumericType(T.X);
```

Open the Model

```
model = 'ComplexPartialSystolicQlessQRMatrixSolveModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

AMBA AXI Handshaking Process

The Data Handler subsystem in this model takes complex matrices A and B as inputs. It sends rows of A and full matrix of B to the QR Decomposition block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set delays for the feeding in rows of A and the feeding in of B matrices in the Data Handler to emulate the processing time of the upstream block. `validInA` and `validInB` remain high when `aDelay` and `bDelay` are set to 0 because this indicates the Data Handler always has data available. When all matrices A and B are sent, the Data Handler loops back to the first A and B matrices.

Asynchronous Matrix Solver

This block operates asynchronously. First, Q-less QR decomposition is performed on the input A matrix and the resulting R matrix is put into a buffer. Then, the Forward Backward Substitute block uses the input B matrix and the buffered R matrix to compute $R'RX = B$. Because the R and B matrices are stored separately in buffers, the upstream Q-less QR decomposition block and the downstream Forward Backward Substitute block can run independently. The Forward Backward Substitute block starts processing when the first R and B matrices are available. Then it runs continuously using the latest buffered R and B matrices, regardless of the status of the Q-less QR Decomposition block. For example, if the upstream block stops providing A and B matrices, the Forward Backward Substitute block continues to generate the same output using the last pair of R and B matrices.

The Data Handler sends A and B matrices to the QR decomposition block iteratively. After sending out the last A matrix, the Data Handler resets its internal counter and sends out first A matrix. The B matrix is handled in a similar fashion.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition block.

```
numOutputs = 1; % Number of recorded outputs
aDelay = 1; % Delay of clock cycles between feeding in rows of A
```

```

bDelay = 1; % Delay of clock cycles between feeding in B matrices
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
    'regularizationParameter', 0, ...
    'aDelay', aDelay, 'bDelay', bDelay, ...
    'numOutputs', numOutputs, 'OutputType', OutputType);

```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition block outputs matrix X at each time step. When a valid result matrix is output, the block sets `validOut` to true.

```
X = out.X;
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition block, compute the relative error. Choose the last output of the simulation.

```
X = double(X(:, :, end));
```

Synchronize the last output X with the input by finding the inputs A and B that produced it.

```

A = double(A);
B = double(B);
relative_errors = zeros(size(A,3),size(B,3));
for k = 1:size(A,3)
    for g = 1:size(B,3)
        relative_errors(k,g) = norm(A(:,:,k)'*A(:,:,k)*X - B(:,:,g))/norm(B(:,:,g));
    end
end
[AUsed,Bused] = find(relative_errors==min(relative_errors,[],'all')) %#ok<NOPTS>

```

```
relative_error = norm(double(A(:,:,AUsed)'*A(:,:,AUsed)*X - B(:,:,Bused)))/norm(double(B(:,:,Bused)))
```

```
AUsed =
```

```
1
```

```
Bused =
```

```
2
```

```
relative_error =
```

```
6.1162e-05
```

Implement Hardware-Efficient Complex Burst Matrix Solve Using QR Decomposition

This example shows how to implement a hardware-efficient least-squares solution to the complex-valued matrix equation $AX=B$ using the Complex Burst Matrix Solve Using QR Decomposition block.

Define Matrix Dimensions

Specify the number of rows in matrices A and B, the number of columns in matrix A, and the number of columns in matrix B.

```
m = 50; % Number of rows in matrices A and B
n = 10; % Number of columns in matrix A
p = 1; % Number of columns in matrix B
```

Generate Random Least-Squares Matrices

For this example, use the helper function `complexRandomLeastSquaresMatrices` to generate random matrices A and B for the least-squares problem $AX=B$. The matrices are generated such that the real and imaginary parts of the elements of A and B are between -1 and +1, and A has rank r.

```
rng('default')
r = 3; % Rank of A
[A,B] = fixed.example.complexRandomLeastSquaresMatrices(m,n,p,r);
```

Select Fixed-Point Data Types

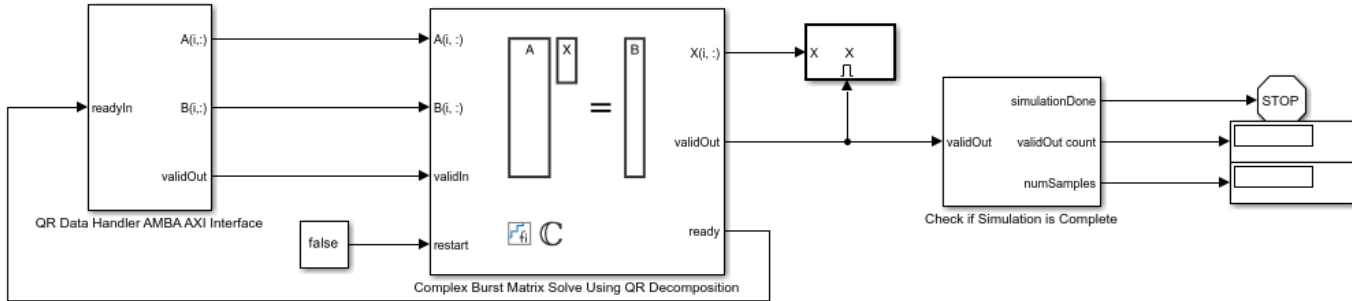
Use the helper function `complexQRMatrixSolveFixedpointTypes` to select fixed-point data types for input matrices A and B, and output X such that there is a low probability of overflow during the computation.

The real and imaginary parts of the elements of A and B are between -1 and 1, so the maximum possible absolute value of any element is $\sqrt{2}$.

```
max_abs_A = sqrt(2); % Upper bound on max(abs(A(:)))
max_abs_B = sqrt(2); % Upper bound on max(abs(B(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.complexQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,precisionBits);
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);
```

Open the Model

```
model = 'ComplexBurstQRMatrixSolveModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

The Data Handler subsystem in this model takes complex matrices A and B as inputs. It sends rows of A and B to QR block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set a delay between the feeding in rows of A and B in the Data Handler to emulate the processing time of the upstream block. `validIn` remains high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Complex Burst Matrix Solve Using QR Decomposition block.

```
numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A and B
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
    'regularizationParameter', 0, ...
    'numSamples', numSamples, 'rowDelay', rowDelay, 'OutputType', OutputType);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Complex Burst Matrix Solve Using QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of X are output in the order they are computed, last row first, so you must reconstruct the data to interpret the results. To reconstruct the matrix X from the output data, use the helper function `matrixSolveModelOutputToArray`.

```
X = fixed.example.matrixSolveModelOutputToArray(out.X, n, p, numSamples);
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Complex Burst Matrix Solve Using QR Decomposition block, compute the relative error.

```
relative_error = norm(double(A*X - B))/norm(double(B)) %#ok<NOPTS>
```

```
relative_error =
```

3.0528e-06

Implement Hardware-Efficient Complex Burst Matrix Solve Using Q-less QR Decomposition

This example shows how to implement a hardware-efficient solution to the complex-valued matrix equation $A'AX=B$ using the Complex Burst Matrix Solve Using Q-less QR Decomposition block.

Define Matrix Dimensions

Specify the number of rows in matrix A, the number of columns in matrix A and rows in B, and the number of columns in matrix B.

```
m = 100; % Number of rows in A
n = 10;  % Number of columns in A and rows in B
p = 1;  % Number of columns in B
```

Generate Matrices

For this example, use the helper function `complexRandomQlessQRMatrices` to generate random matrices A and B for the problem $A'AX=B$. The matrices are generated such that the real and imaginary parts of the elements of A and B are between -1 and +1, and A is full rank.

```
rng('default')
[A,B] = fixed.example.complexRandomQlessQRMatrices(m,n,p);
```

Select Fixed-Point Data Types

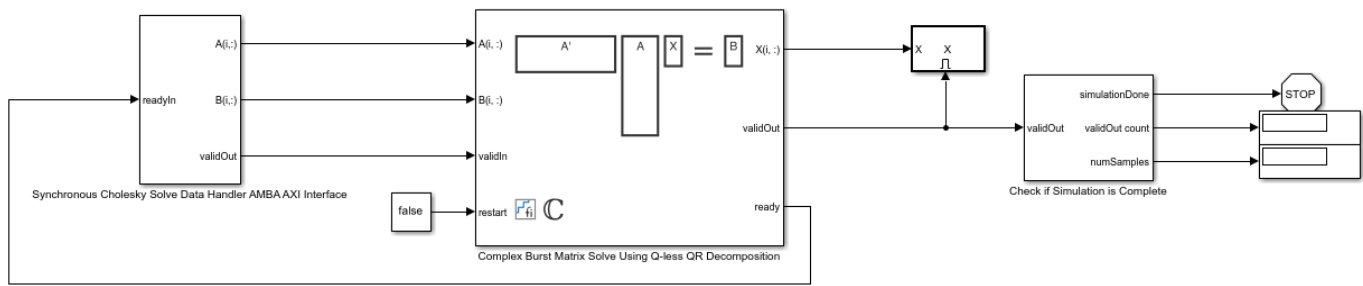
Use the helper function `complexQlessQRMatrixSolveFixedpointTypes` to select fixed-point data types for input matrices A and B, and output X such that there is a low probability of overflow during the computation.

The real and imaginary parts of the elements of A and B are between -1 and 1, so the maximum possible absolute value of any element is $\sqrt{2}$.

```
max_abs_A = sqrt(2); % Upper bound on max(abs(A(:)))
max_abs_B = sqrt(2); % Upper bound on max(abs(B(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.complexQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,precisionBits);
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);
```

Open the Model

```
model = 'ComplexBurstQlessQRMatrixSolveModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

The Data Handler subsystem in this model takes complex matrices A and B as inputs. It sends rows of A and B to QR block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set a delay between the feeding in rows of A and B in the Data Handler to emulate the processing time of the upstream block. `validIn` remains high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Complex Burst Matrix Solve Using Q-less QR Decomposition block.

```
numSamples = 1; % Number of samples
rowDelay = 1; % Delay of clock cycles between feeding in rows of A and B
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
    'regularizationParameter', 0, ...
    'numSamples', numSamples, 'rowDelay', rowDelay, 'OutputType', OutputType);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Complex Burst Matrix Solve Using Q-less QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of X are output in the order they are computed, last row first, so you must reconstruct the data to interpret the results. To reconstruct the matrix X from the output data, use the helper function `matrixSolveModelOutputToArray`.

```
X = fixed.example.matrixSolveModelOutputToArray(out.X, n, p);
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Complex Burst Matrix Solve Using Q-less QR Decomposition block, compute the relative error.

```
relative_error = norm(double(A'*A*X - B))/norm(double(B)) %#ok<NOPTS>
```

```
relative_error =
```

1.1496e-04

Implement Hardware-Efficient Complex Partial-Systolic QR Decomposition

This example shows how to implement a hardware-efficient QR decomposition using the Complex Partial-Systolic QR Decomposition block.

Economy Size QR Decomposition

The Complex Partial-Systolic QR Decomposition block performs the first step of solving the least-squares matrix equation $AX = B$ which transforms A in-place to R and B in-place to $C = Q'B$, then solves the transformed system $RX = C$, where QR is the orthogonal-triangular decomposition of A.

To compute the stand-alone QR decomposition, this example sets B to be the identity matrix so that the output of the Complex Partial-Systolic QR Decomposition block is the upper-triangular R and $C = Q'$.

Define Matrix Dimensions

Specify the number of rows in matrices A and B, the number of columns in matrix A, and the number of columns in matrix B. This example sets B to be the identity matrix the same size as the number of rows of A.

```
m = 10; % Number of rows in matrices A and B
n = 3;  % Number of columns in matrix A
p = m;  % Number of columns in matrix B
```

Generate Matrices A and B

Use the helper function `complexUniformRandomArray` to generate a random matrix A such that the real and imaginary parts of the elements of A are between -1 and +1, and A is full rank. Matrix B is the identity matrix.

```
rng('default')
A = fixed.example.complexUniformRandomArray(-1,1,m,n);
B = eye(m);
```

Select Fixed-Point Data Types

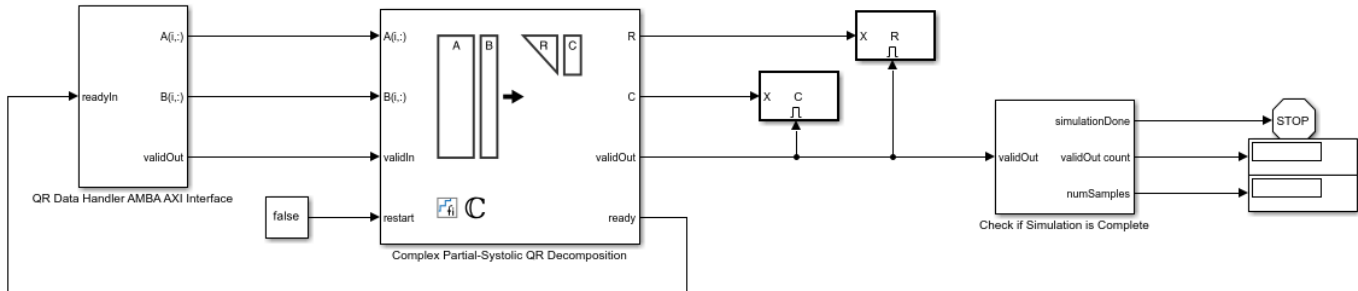
Use the helper function `qrFixedpointTypes` to select fixed-point data types for input matrices A and B that guarantee no overflow will occur in the transformation of A in-place to R and B in-place to $C = Q'B$.

The real and imaginary parts of the elements of A are between -1 and 1, so the maximum possible absolute value of any element is $\sqrt{2}$.

```
max_abs_A = sqrt(2); % Upper bound on max(abs(A(:)))
max_abs_B = 1;      % Upper bound on max(abs(B(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.qrFixedpointTypes(m,max_abs_A,max_abs_B,precisionBits);
A = cast(A,'like',T.A);
B = complex(cast(B,'like',T.B));
```

Open the Model

```
model = 'ComplexPartialSystolicQRModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

The Data Handler subsystem in this model takes complex matrices A and B as inputs. It sends rows of A and B to QR block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set a delay between the feeding in rows of A and B in the Data Handler to emulate the processing time of the upstream block. `validIn` remains high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Complex Partial-Systolic QR Decomposition block.

```
numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A and B
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
    'numSamples', numSamples, 'rowDelay', rowDelay);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Complex Partial-Systolic QR Decomposition block outputs matrices R and C at each time step. When valid result matrices are output, the block sets `validOut` to true.

```
R = out.R;
C = out.C;
```

Extract the Economy-Size Q

The block computes $C = Q'B$. In this example, B is the identity matrix, so $Q = C'$ is the economy-size orthogonal factor of the QR decomposition.

```
Q = C';
```

Verify that Q is Orthogonal and R is Upper-Triangular

Q is orthogonal, so $Q'Q$ is the identity matrix within roundoff.

```
I = Q'*Q
```

```
I =
```

```
 1.0000 + 0.0000i  -0.0000 + 0.0000i  -0.0000 + 0.0000i
-0.0000 - 0.0000i  1.0000 + 0.0000i  -0.0000 + 0.0000i
-0.0000 - 0.0000i  -0.0000 - 0.0000i  1.0000 + 0.0000i
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 62
      FractionLength: 48
```

R is an upper-triangular matrix.

```
R
```

```
R =
```

```
 3.1655 + 0.0000i  0.4870 + 1.1980i  0.1466 - 0.9092i
 0.0000 + 0.0000i  2.2184 + 0.0000i  -0.2159 - 0.0972i
 0.0000 + 0.0000i  0.0000 + 0.0000i  2.2903 + 0.0000i
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 29
      FractionLength: 24
```

```
isequal(R, triu(R))
```

```
ans =
```

```
logical
```

```
1
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Complex Partial-Systolic QR Decomposition block, compute the relative error.

```
relative_error = norm(double(Q*R - A))/norm(double(A))
```

```
relative_error =
```

```
1.2460e-06
```

Suppress mlint warnings.

%#ok<*NOPTS>

Implement Hardware-Efficient Complex Partial-Systolic Q-less QR Decomposition

This example shows how to implement a hardware-efficient Q-less QR decomposition using the Complex Partial-Systolic Q-less QR Decomposition block.

Economy Size Q-less QR Decomposition

The Complex Partial-Systolic Q-less QR Decomposition block performs the first step of solving the matrix equation $A'AX = B$ which transforms A in-place to upper-triangular R, then solves the transformed system $R'RX = B$, where $R'R = A'A$.

Define Matrix Dimensions

Specify the number of rows and columns in matrix A.

```
m = 5; % Number of rows in matrix A
n = 3; % Number of columns in matrix A
```

Generate Matrix A

Use the helper function `complexUniformRandomArray` to generate a random matrix A such that the real and imaginary parts of the elements of A are between -1 and +1, and A is full rank.

```
rng('default')
A = fixed.example.complexUniformRandomArray(-1,1,m,n);
```

Select Fixed-Point Data Types

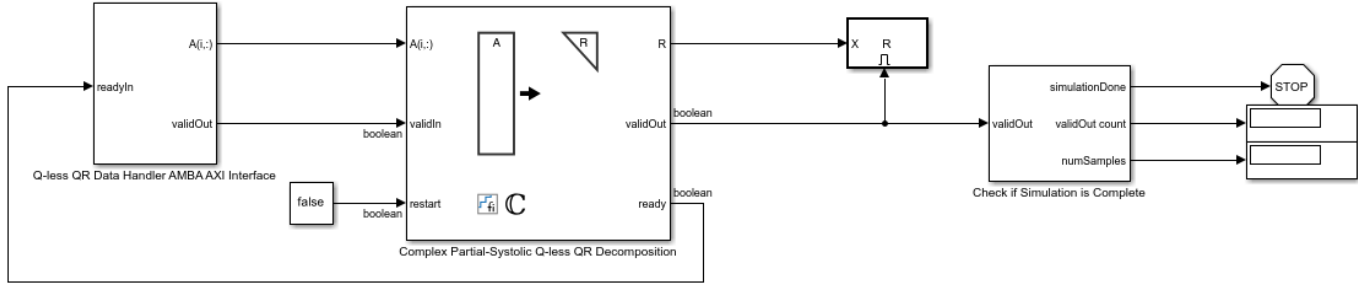
Use the helper function `qlessqrFixedpointTypes` to select fixed-point data types for matrix A that guarantee no overflow will occur in the transformation of A in-place to R.

The real and imaginary parts of the elements of A are between -1 and 1, so the maximum possible absolute value of any element is $\sqrt{2}$.

```
max_abs_A = sqrt(2); % Upper bound on max(abs(A(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.qlessqrFixedpointTypes(m,max_abs_A,precisionBits);
A = cast(A,'like',T.A);
```

Open the Model

```
model = 'ComplexPartialSystolicQlessQRModel';
open_system(model);
```

Copyright 2021 The MathWorks, Inc.

AMBA AXI Handshaking Process

The Data Handler subsystem in this model takes real matrix A as input. It sends rows of A to the QR Decomposition block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set delay for the feeding in rows of A in the Data Handler to emulate the processing time of the upstream block. `validOut` signal of the Data Handler remain high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Complex Partial-Systolic Q-less QR Decomposition block.

```
numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A
fixed.example.setModelWorkspace(model, 'A', A, 'm', m, 'n', n, ...
    'numSamples', numSamples, 'rowDelay', rowDelay);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Complex Partial-Systolic QR Decomposition block outputs matrix R at each time step. When a valid result matrix is output, the block sets `validOut` to true.

```
R = out.R;
```

R is an upper-triangular matrix.

R

$R =$

```
2.1863 + 0.0000i   0.6427 - 1.0882i   -0.5771 - 0.3089i
0.0000 + 0.0000i   1.8126 + 0.0000i    0.2095 + 0.0599i
0.0000 + 0.0000i   0.0000 + 0.0000i    1.7760 + 0.0000i
```

DataTypeMode: Fixed-point: binary point scaling

```
Signedness: Signed  
WordLength: 29  
FractionLength: 24
```

```
isequal(R, triu(R))
```

```
ans =
```

```
logical
```

```
1
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Complex Partial-Systolic Q-less QR Decomposition block, compute the relative error.

```
relative_error = norm(double(R'*R - A'*A))/norm(double(A'*A))
```

```
relative_error =
```

```
9.1285e-07
```

Suppress mlint warnings.

```
 %#ok<*NOPTS>
```

Implement Hardware-Efficient Complex Burst QR Decomposition

This example shows how to implement a hardware-efficient QR decomposition using the Complex Burst QR Decomposition block.

Economy Size QR Decomposition

The Complex Burst QR Decomposition block performs the first step of solving the least-squares matrix equation $AX = B$ which transforms A in-place to R and B in-place to $C = Q'B$, then solves the transformed system $RX = C$, where QR is the orthogonal-triangular decomposition of A.

To compute the stand-alone QR decomposition, this example sets B to be the identity matrix so that the output of the Complex Burst QR Decomposition block is the upper-triangular R and $C = Q'$.

Define Matrix Dimensions

Specify the number of rows in matrices A and B, the number of columns in matrix A, and the number of columns in matrix B. This example sets B to be the identity matrix the same size as the number of rows of A.

```
m = 10; % Number of rows in matrices A and B
n = 3;  % Number of columns in matrix A
p = m; % Number of columns in matrix B
```

Generate Matrices A and B

Use the helper function `complexUniformRandomArray` to generate a random matrix A such that the real and imaginary parts of the elements of A are between -1 and +1, and A is full rank. Matrix B is the identity matrix.

```
rng('default')
A = fixed.example.complexUniformRandomArray(-1,1,m,n);
B = eye(m);
```

Select Fixed-Point Data Types

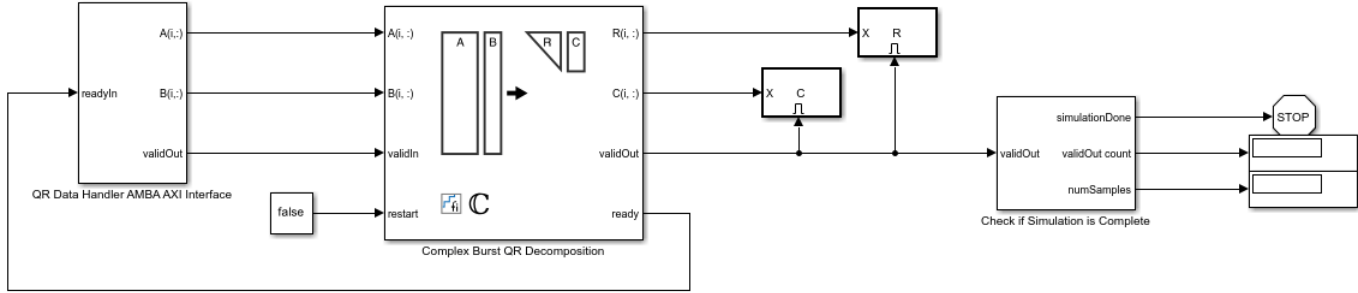
Use the helper function `qrFixedpointTypes` to select fixed-point data types for matrices A and B that guarantee no overflow will occur in the transformation of A in-place to R and B in-place to $C = Q'B$.

The real and imaginary parts of the elements of A are between -1 and 1, so the maximum possible absolute value of any element is $\sqrt{2}$.

```
max_abs_A = sqrt(2); % Upper bound on max(abs(A(:)))
max_abs_B = 1;      % Upper bound on max(abs(B(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.qrFixedpointTypes(m,max_abs_A,max_abs_B,precisionBits);
A = cast(A,'like',T.A);
B = complex(cast(B,'like',T.B));
```

Open the Model

```
model = 'ComplexBurstQRModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

The Data Handler subsystem in this model takes complex matrices A and B as inputs. It sends rows of A and B to QR block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set a delay between the feeding in rows of A and B in the Data Handler to emulate the processing time of the upstream block. `validIn` remains high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Complex Burst QR Decomposition block.

```
numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A and B
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
    'numSamples', numSamples, 'rowDelay', rowDelay);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Complex Burst QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of matrices R and C are output in reverse order to accommodate back-substitution, so you must reconstruct the data to interpret the results. To reconstruct the matrices R and C from the output data, use the helper function `qrModelOutputToArray`.

```
[C,R] = fixed.example.qrModelOutputToArray(out.C,out.R,m,n,p,numSamples);
```

Extract the Economy-Size Q

The block computes $C = Q'B$. In this example, B is the identity matrix, so $Q = C'$ is the economy-size orthogonal factor of the QR decomposition.

```
Q = C';
```

Verify that Q is Orthogonal and R is Upper-Triangular

Q is orothogonal, so $Q'Q$ is the identity matrix within roundoff.

```
I = Q'*Q
```

```
I =
```

```
 1.0000 + 0.0000i  -0.0000 + 0.0000i  -0.0000 + 0.0000i
-0.0000 - 0.0000i  1.0000 + 0.0000i  -0.0000 + 0.0000i
-0.0000 - 0.0000i  -0.0000 - 0.0000i  1.0000 + 0.0000i
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 62
      FractionLength: 48
```

R is an upper-triangular matrix.

```
R
```

```
R =
```

```
 3.1655 + 0.0000i  0.4870 + 1.1980i  0.1466 - 0.9092i
 0.0000 + 0.0000i  2.2184 + 0.0000i  -0.2159 - 0.0972i
 0.0000 + 0.0000i  0.0000 + 0.0000i  2.2903 + 0.0000i
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 29
      FractionLength: 24
```

```
isequal(R, triu(R))
```

```
ans =
```

```
logical
```

```
1
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Complex Burst QR Decomposition block, compute the relative error.

```
relative_error = norm(double(Q*R - A))/norm(double(A))
```

```
relative_error =
```

```
1.2460e-06
```

Suppress mlint warnings.

```
 %#ok< *NOPTS>
```

Implement Hardware-Efficient Complex Burst Q-less QR Decomposition

This example shows how to implement a hardware-efficient Q-less QR decomposition using the Complex Burst Q-less QR Decomposition block.

Economy Size Q-less QR Decomposition

The Complex Burst Q-less QR Decomposition block performs the first step of solving the matrix equation $A'AX = B$ which transforms A in-place to upper-triangular R, then solves the transformed system $R'RX = B$, where $R'R = A'A$.

Define Matrix Dimensions

Specify the number of rows and columns in matrix A.

```
m = 5; % Number of rows in matrix A
n = 3; % Number of columns in matrix A
```

Generate Matrix A

Use the helper function `complexUniformRandomArray` to generate a random matrix A such that the real and imaginary parts of the elements of A are between -1 and +1, and A is full rank.

```
rng('default')
A = fixed.example.complexUniformRandomArray(-1,1,m,n);
```

Select Fixed-Point Data Types

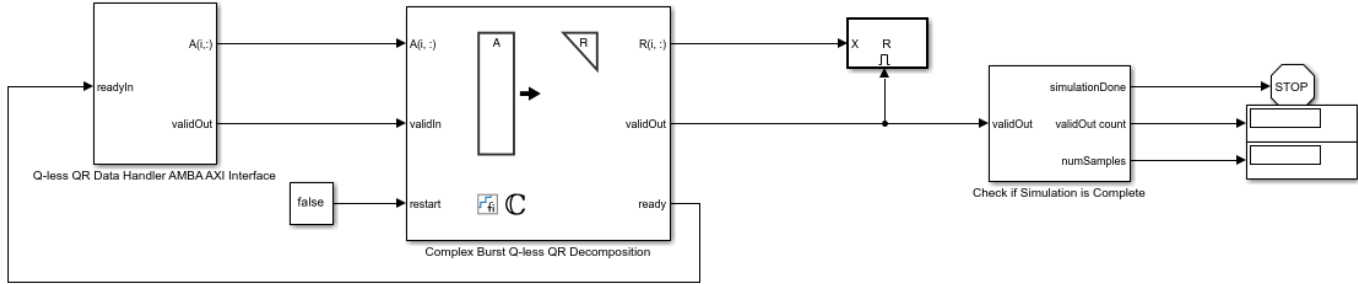
Use the helper function `qllessqrFixedpointTypes` to select fixed-point data types for matrix A that guarantee no overflow will occur in the transformation of A in-place to R.

The real and imaginary parts of the elements of A are between -1 and 1, so the maximum possible absolute value of any element is $\sqrt{2}$.

```
max_abs_A = sqrt(2); % Upper bound on max(abs(A(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.qllessqrFixedpointTypes(m,max_abs_A,precisionBits);
A = cast(A,'like',T.A);
```

Open the Model

```
model = 'ComplexBurstQlessQRModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

AMBA AXI Handshaking Process

The Data Handler subsystem in this model takes complex matrix A as input. It sends rows of A to the QR Decomposition block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set delay for the feeding in rows of A in the Data Handler to emulate the processing time of the upstream block. `validOut` signal of the Data Handler remain high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Complex Burst Q-less QR Decomposition block.

```
numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A
fixed.example.setModelWorkspace(model, 'A', A, 'm', m, 'n', n, ...
    'numSamples', numSamples, 'rowDelay', rowDelay);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Complex Burst Q-less QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of matrix R are output in reverse order to accommodate back-substitution, so you must reconstruct the data to interpret the results. To reconstruct the matrix R from the output data, use the helper function `qlessqrModelOutputToArray`.

```
R = fixed.example.qlessqrModelOutputToArray(out.R, m, n, numSamples);
```

R is an upper-triangular matrix.

R

$R =$

```
2.1863 + 0.0000i    0.6427 - 1.0882i   -0.5771 - 0.3089i
0.0000 + 0.0000i    1.8126 + 0.0000i    0.2095 + 0.0599i
```

```
0.0000 + 0.0000i  0.0000 + 0.0000i  1.7760 + 0.0000i
```

```
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 29  
    FractionLength: 24
```

```
isequal(R, triu(R))
```

```
ans =
```

```
    logical
```

```
    1
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Complex Burst Q-less QR Decomposition block, compute the relative error.

```
relative_error = norm(double(R'*R - A'*A))/norm(double(A'*A))
```

```
relative_error =
```

```
    8.3255e-07
```

Suppress mlint warnings.

```
 %#ok<*NOPTS>
```


Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using QR Decomposition with Diagonal Loading

This example shows how to implement a hardware-efficient least-squares solution to the real-valued matrix equation

$$\begin{bmatrix} \lambda I \\ A \end{bmatrix} X = \begin{bmatrix} O \\ B \end{bmatrix}$$

using the Real Partial-Systolic Matrix Solve Using QR Decomposition block. This method is known as diagonal loading, and λ is known as a regularization parameter.

Diagonal Loading Method

When you set the regularization parameter to a non-zero value in block Real Partial-Systolic Matrix Solve Using QR Decomposition, then the block computes the the least-squares solution to

$$\begin{bmatrix} \lambda I \\ A \end{bmatrix} X = \begin{bmatrix} O \\ B \end{bmatrix}$$

where I is an n-by-n identity matrix and O is an array of zeros of size n-by-p.

This method is known as diagonal loading, and λ is known as a regularization parameter. The familiar textbook least-squares solution for this equation is the following, which is derived by multiplying both sides of the equation by $[\lambda I, A']$ and taking the inverse of the matrix on the left-hand side.

$$X_{LS} = (\lambda^2 I + A' A)^{-1} A' B.$$

Real Partial-Systolic Matrix Solve Using QR Decomposition block computes the solution efficiently without computing an inverse by computing the QR decomposition, transforming

$$\begin{bmatrix} \lambda I \\ A \end{bmatrix}$$

in-place to upper-triangular R, and transforming

$$\begin{bmatrix} O \\ B \end{bmatrix}$$

in-place to

$$C = Q' \begin{bmatrix} O \\ B \end{bmatrix}$$

and solving the transformed equation $RX = C$.

Define Matrix Dimensions

Specify the number of rows in matrices A and B, the number of columns in matrix A, and the number of columns in matrix B.

```
m = 300; % Number of rows in matrices A and B
n = 10; % Number of columns in matrix A
p = 1; % Number of columns in matrix B
```

Define the Regularization Parameter

When the regularization parameter is non-zero in block Real Partial-Systolic Matrix Solve Using QR Decomposition, then the diagonal-loading method is used. When the regularization parameter is zero, then the equations reduce to the regular least-squares formula $AX=B$.

```
regularizationParameter = 1e-3;
```

Block Parameters

Block Parameters: Real Partial-Systolic Matrix Solve U...

Real Partial-Systolic Matrix Solve Using QR Decomposition (mask) (link)

Compute the value of x in the equation $Ax = B$, where A and B are real-valued matrices.

Use the partial-systolic implementation to minimize system latency and increase the throughput. Partial systolic-implementations require more hardware resources than burst implementations.

Parameters

Number of rows in matrices A and B

Number of columns in matrix A

Number of columns in matrix B

Regularization parameter

Output datatype >>

OK Cancel Help Apply

Generate Random Least-Squares Matrices

For this example, use the helper function `realRandomLeastSquaresMatrices` to generate random matrices A and B for the least-squares problem $AX=B$. The matrices are generated such that the elements of A and B are between -1 and +1, and A has rank r.

```
rng('default')
r = 3; % Rank of matrix A
[A,B] = fixed.example.realRandomLeastSquaresMatrices(m,n,p,r);
```

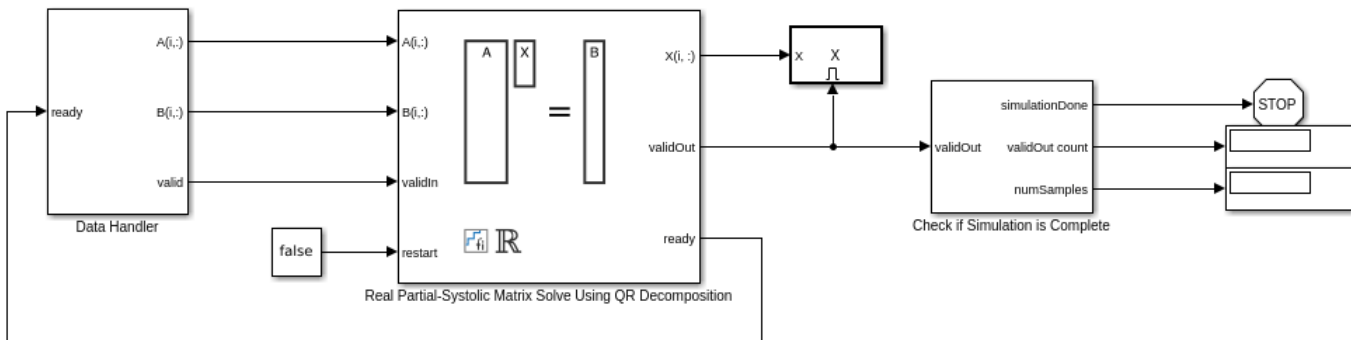
Select Fixed-Point Data Types

Use the helper function `realQRMatrixSolveFixedpointTypes` to select fixed-point data types for input matrices A and B, and output X such that there is a low probability of overflow during the computation.

```
max_abs_A = 1; % Upper bound on max(abs(A(:)))
max_abs_B = 1; % Upper bound on max(abs(B(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.realQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,precisionBits);
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);
```

Open the Model

```
model = 'RealPartialSystolicQRDiagonalLoadingMatrixSolveModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

The Data Handler subsystem in this model takes real matrices A and B as inputs. The `ready` port triggers the Data Handler. After sending a true `validIn` signal, there may be some delay before `ready` is set to false. When the Data Handler detects the leading edge of the `ready` signal, the block sets `validIn` to true and sends the next row of A and B. This protocol allows data to be sent whenever a leading edge of the `ready` signal is detected, ensuring that all data is processed.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Real Partial-Systolic Matrix Solve Using QR Decomposition block.

```
numSamples = 1; % Number of sample matrices
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
```

```
'regularizationParameter', regularizationParameter, ...  
'numSamples', numSamples, 'OutputType', OutputType);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Real Partial-Systolic Matrix Solve Using QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of `X` are output in the order they are computed, last row first, so you must reconstruct the data to interpret the results. To reconstruct the matrix `X` from the output data, use the helper function `matrixSolveModelOutputToArray`.

```
X = fixed.example.matrixSolveModelOutputToArray(out.X, n, p, numSamples);
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Real Partial-Systolic Matrix Solve Using QR Decomposition block, compute the relative error.

```
A_lambda = [regularizationParameter * eye(n); A];  
B_0 = [zeros(n, p); B];  
relative_error = norm(double(A_lambda*X - B_0))/norm(double(B_0)) %#ok<NOPTS>
```

```
relative_error =
```

```
1.2189e-04
```

Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using QR Decomposition with Diagonal Loading

This example shows how to implement a hardware-efficient least-squares solution to the complex-valued matrix equation

$$\begin{bmatrix} \lambda I \\ A \end{bmatrix} X = \begin{bmatrix} O \\ B \end{bmatrix}$$

using the Complex Partial-Systolic Matrix Solve Using QR Decomposition block. This method is known as diagonal loading, and λ is known as a regularization parameter.

Diagonal Loading Method

When you set the regularization parameter to a non-zero value in block Complex Partial-Systolic Matrix Solve Using QR Decomposition, then the block computes the the least-squares solution to

$$\begin{bmatrix} \lambda I \\ A \end{bmatrix} X = \begin{bmatrix} O \\ B \end{bmatrix}$$

where I is an n-by-n identity matrix and O is an array of zeros of size n-by-p.

This method is known as diagonal loading, and λ is known as a regularization parameter. The familiar textbook least-squares solution for this equation is the following, which is derived by multiplying both sides of the equation by $[\lambda I, A']$ and taking the inverse of the matrix on the left-hand side.

$$X_{LS} = (\lambda^2 I + A' A)^{-1} A' B.$$

Complex Partial-Systolic Matrix Solve Using QR Decomposition block computes the solution efficiently without computing an inverse by computing the QR decomposition, transforming

$$\begin{bmatrix} \lambda I \\ A \end{bmatrix}$$

in-place to upper-triangular R, and transforming

$$\begin{bmatrix} O \\ B \end{bmatrix}$$

in-place to

$$C = Q' \begin{bmatrix} O \\ B \end{bmatrix}$$

and solving the transformed equation $RX = C$.

Define Matrix Dimensions

Specify the number of rows in matrices A and B, the number of columns in matrix A, and the number of columns in matrix B.

```
m = 300; % Number of rows in matrices A and B
n = 10; % Number of columns in matrix A
p = 1; % Number of columns in matrix B
```

Define the Regularization Parameter

When the regularization parameter is non-zero in block Complex Partial-Systolic Matrix Solve Using QR Decomposition, then the diagonal-loading method is used. When the regularization parameter is zero, then the equations reduce to the regular least-squares formula $AX=B$.

```
regularizationParameter = 1e-3;
```

Block Parameters

Block Parameters: Complex Partial-Systolic Matrix Solve...

Complex Partial-Systolic Matrix Solve Using QR Decomposition (mask) (link)

Compute the value of x in the equation $Ax = B$, where A and B are complex-valued matrices.

Use the partial-systolic implementation to minimize system latency and increase the throughput. Partial systolic-implementations require more hardware resources than burst implementations.

Parameters

Number of rows in matrices A and B

Number of columns in matrix A

Number of columns in matrix B

Regularization parameter

Output datatype >>

OK Cancel Help Apply

Generate Random Least-Squares Matrices

For this example, use the helper function `complexRandomLeastSquaresMatrices` to generate random matrices A and B for the least-squares problem $AX=B$. The matrices are generated such that the elements of A and B are between -1 and +1, and A has rank r.

```
rng('default')
r = 3; % Rank of matrix A
[A,B] = fixed.example.complexRandomLeastSquaresMatrices(m,n,p,r);
```

Select Fixed-Point Data Types

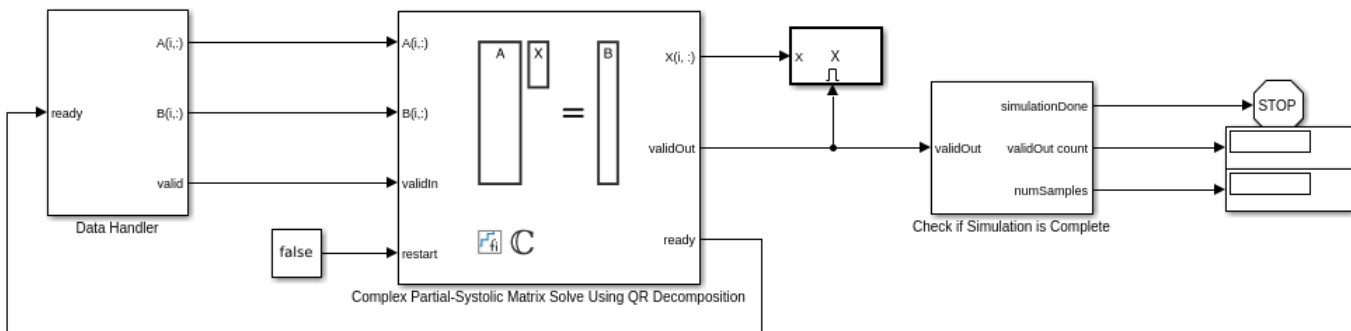
Use the helper function `complexQRMatrixSolveFixedpointTypes` to select fixed-point data types for input matrices A and B, and output X such that there is a low probability of overflow during the computation.

The real and imaginary parts of the elements of A and B are between -1 and 1, so the maximum possible absolute value of any element is $\sqrt{2}$.

```
max_abs_A = sqrt(2); % Upper bound on max(abs(A(:)))
max_abs_B = sqrt(2); % Upper bound on max(abs(B(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.complexQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,precisionBits);
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);
```

Open the Model

```
model = 'ComplexPartialSystolicQRDiagonalLoadingMatrixSolveModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

The Data Handler subsystem in this model takes complex matrices A and B as inputs. The `ready` port triggers the Data Handler. After sending a true `validIn` signal, there may be some delay before `ready` is set to false. When the Data Handler detects the leading edge of the `ready` signal, the block sets `validIn` to true and sends the next row of A and B. This protocol allows data to be sent whenever a leading edge of the `ready` signal is detected, ensuring that all data is processed.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Complex Partial-Systolic Matrix Solve Using QR Decomposition block.

```
numSamples = 1; % Number of sample matrices
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
    'regularizationParameter', regularizationParameter, ...
    'numSamples', numSamples, 'OutputType', OutputType);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Complex Partial-Systolic Matrix Solve Using QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of `X` are output in the order they are computed, last row first, so you must reconstruct the data to interpret the results. To reconstruct the matrix `X` from the output data, use the helper function `matrixSolveModelOutputToArray`.

```
X = fixed.example.matrixSolveModelOutputToArray(out.X, n, p, numSamples);
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Complex Partial-Systolic Matrix Solve Using QR Decomposition block, compute the relative error.

```
A_lambda = [regularizationParameter * eye(n); A];
B_0 = [zeros(n, p); B];
relative_error = norm(double(A_lambda*X - B_0))/norm(double(B_0)) %#ok<NOPTS>
```

```
relative_error =
```

```
1.0386e-04
```


Determine Fixed-Point Types for QR Decomposition

This example shows how to use `fixed.qrFixedpointTypes` to analytically determine fixed-point types for the computation of the QR decomposition.

Define Matrix Dimensions

Specify the number of rows in matrices A and B , the number of columns in matrix A , and the number of columns in matrix B . This example sets B to be the identity matrix the same size as the number of rows of A .

```
m = 10; % Number of rows in matrices A and B
n = 3;  % Number of columns in matrix A
```

Generate Matrices A and B

Use the helper function `realUniformRandomArray` to generate a random matrix A such that the elements of A are between -1 and $+1$. Matrix B is the identity matrix.

```
rng('default')
A = fixed.example.realUniformRandomArray(-1,1,m,n);
B = eye(m);
```

Select Fixed-Point Types

Use `fixed.qrFixedpointTypes` to select fixed-point data types for matrices A and B that guarantee no overflow will occur in the transformation of A in-place to $R = Q'A$ and B in-place to $C = Q'B$.

```
max_abs_A = 1; % Upper bound on max(abs(A(:)))
max_abs_B = 1; % Upper bound on max(abs(B(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.qrFixedpointTypes(m,max_abs_A,max_abs_B,precisionBits)

T = struct with fields:
    A: [0x0 embedded.fi]
    B: [0x0 embedded.fi]
```

$T.A$ is the type computed for transforming A to $R = Q'A$ in-place so that it does not overflow.

$T.A$

ans =

[]

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 29
    FractionLength: 24
```

$T.B$ is the type computed for transforming B to $C = Q'B$ in-place so that it does not overflow.

$T.B$

ans =

```

[]
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 29
    FractionLength: 24

```

Use the Specified Types to Compute the QR Decomposition

Cast the inputs to the types determined by `fixed.qrFixedpointTypes`.

```

A = cast(A,'like',T.A);
B = cast(B,'like',T.B);

```

Accelerate `fixed.qrAB` by using `fiaccl` to generate a MATLAB executable (MEX) function.

```

fiaccl fixed.qrAB -args {A,B} -o qrAB_mex

```

Compute the QR decomposition.

```

[C,R] = qrAB_mex(A,B);

```

Extract the Economy-Size Q

The function `fixed.qrAB` transforms A to $R = Q'A$ and B to $C = Q'B$. In this example, B is the identity matrix, so $Q = C'$ is the economy-size orthogonal factor of the QR decomposition.

```

Q = C';

```

Verify that Q is Orthogonal and R is Upper-Triangular

Q is orthogonal, so $Q'Q$ is the identity matrix within rounding error.

```

I = Q'*Q

```

```

I =
    1.0000    -0.0000    -0.0000
   -0.0000    1.0000    -0.0000
   -0.0000    -0.0000    1.0000

```

```

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 62
    FractionLength: 48

```

R is an upper-triangular matrix.

```

R

```

```

R =
    2.2180    0.8559   -0.5607
         0    2.0578   -0.4017
         0         0    1.7117

```

```

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 29
    FractionLength: 24

```

```

isequal(R, triu(R))

```

```
ans = logical  
     1
```

Verify the Accuracy of the Output

To evaluate the accuracy of the `fixed.qrAB` function, compute the relative error.

```
relative_error = norm(double(Q*R - A))/norm(double(A))
```

```
relative_error = 1.5886e-06
```

Suppress `mlint` warnings.

```
 %#ok<*NOPTS>
```

See Also

Functions

`fixed.qrFixedpointTypes` | `fixed.qrAB` | `qr`

Blocks

Real Burst QR Decomposition | Complex Burst QR Decomposition | Real Partial-Systolic QR Decomposition | Complex Partial-Systolic QR Decomposition

Determine Fixed-Point Types for Q-less QR Decomposition

This example shows how to use `fixed.qlessqrFixedpointTypes` to analytically determine a fixed-point type for the computation of the Q-less QR decomposition.

Define Matrix Dimensions

Specify the number of rows and columns in matrix A .

```
m = 10; % Number of rows in matrix A
n = 3;  % Number of columns in matrix A
```

Generate Matrix A

Use the helper function `realUniformRandomArray` to generate a random matrix A such that the elements of A are between -1 and $+1$.

```
rng('default')
A = fixed.example.realUniformRandomArray(-1,1,m,n);
```

Select Fixed-Point Type

Use the `fixed.qlessqrFixedpointTypes` function to select the fixed-point data type for matrix A that guarantees no overflow will occur in the transformation of A in-place to $R = Q'A$.

```
max_abs_A = 1; % Upper bound on max(abs(A(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.qlessqrFixedpointTypes(m,max_abs_A,precisionBits)

T = struct with fields:
    A: [0x0 embedded.fi]
```

$T.A$ is the type computed for transforming A to $R = Q'A$ in-place so that it does not overflow.

$T.A$

ans =

[]

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 29
    FractionLength: 24
```

Use the Specified Type to Compute the Q-less QR Decomposition

Cast the input to the type determined by `fixed.qlessqrFixedpointTypes`.

```
A = cast(A,'like',T.A);
```

Accelerate `fixed.qlessQR` by using `fiaccel` to generate a MATLAB executable (MEX) function.

```
fiaccel fixed.qlessQR -args {A} -o qlessQR_mex
```

Compute the QR decomposition.

```
R = qlessQR_mex(A);
```

Verify that R is Upper-Triangular

R is an upper-triangular matrix.

R

```
R =
    2.2180    0.8559   -0.5607
         0    2.0578   -0.4017
         0         0    1.7117
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 29
    FractionLength: 24
```

```
isequal(R, triu(R))
```

```
ans = logical
      1
```

Verify the Accuracy of the Output

To evaluate the accuracy of the `fixed.qlessQR` function, compute the relative error.

$R = Q'A$, and Q is orthogonal, so $R'R = A'QQ'A = A'A$, within rounding error.

```
relative_error = norm(double(R'*R - A'*A))/norm(double(A'*A))
```

```
relative_error = 9.3865e-07
```

Suppress `mLint` warnings.

```
 %#ok<*NOPTS>
```

See Also

Functions

`fixed.qlessqrFixedpointTypes` | `fixed.qlessQR`

Blocks

Real Burst Q-less QR Decomposition | Complex Burst QR Decomposition | Real Partial-Systolic Q-less QR Decomposition | Complex Partial-Systolic Q-less QR Decomposition | Real Partial-Systolic Q-less QR Decomposition with Forgetting Factor | Complex Partial-Systolic Q-less QR Decomposition with Forgetting Factor

Algorithms to Determine Fixed-Point Types for Complex Q-less QR Matrix Solve $A'AX=B$

This example shows the algorithms that the `fixed.complexQlessQRMatrixSolveFixedpointTypes` function uses to analytically determine fixed-point types for the solution of the complex matrix equation $A'AX = B$, where A is an m -by- n matrix with $m \geq n$, B is n -by- p , and X is n -by- p .

Overview

You can solve the fixed-point matrix equation $A'AX = B$ using QR decomposition. Using a sequence of orthogonal transformations, QR decomposition transforms matrix A in-place to upper triangular R , where $QR = A$ is the economy-size QR decomposition. This reduces the equation to an upper-triangular system of equations $R'RX = B$. To solve for X , compute $X = R \setminus (R' \setminus B)$ through forward- and backward-substitution of R into B .

You can determine appropriate fixed-point types for the matrix equation $A'AX = B$ by selecting the fraction length based on the number of bits of precision defined by your requirements. The `fixed.complexQlessQRMatrixSolveFixedpointTypes` function analytically computes the following upper bounds on R , and X to determine the number of integer bits required to avoid overflow [1,2,3].

The upper bound for the magnitude of the elements of $R = Q'A$ is

$$\max(|R(:)|) \leq \sqrt{m} \max(|A(:)|).$$

The upper bound for the magnitude of the elements of $X = (A'A) \setminus B$ is

$$\max(|X(:)|) \leq \frac{\sqrt{n} \max(|B(:)|)}{\min(\text{svd}(A))^2}.$$

Since computing `svd(A)` is more computationally expensive than solving the system of equations, the `fixed.complexQlessQRMatrixSolveFixedpointTypes` function estimates a lower bound of `min(svd(A))`.

Fixed-point types for the solution of the matrix equation $(A'A)X = B$ are generally well-bounded if the number of rows, m , of A are much greater than the number of columns, n (i.e. $m \gg n$), and A is full rank. If A is not inherently full rank, then it can be made so by adding random noise. Random noise naturally occurs in physical systems, such as thermal noise in radar or communications systems. If $m = n$, then the dynamic range of the system can be unbounded, for example in the scalar equation $x = a^2/b$ and $a, b \in [-1, 1]$, then x can be arbitrarily large if b is close to 0.

Proofs of the Bounds

Properties and Definitions of Vector and Matrix Norms

The proofs of the bounds use the following properties and definitions of matrix and vector norms, where Q is an orthogonal matrix, and v is a vector of length m [6].

$$\begin{aligned} \|Av\|_2 &\leq \|A\|_2 \|v\|_2 \\ \|Q\|_2 &= 1 \\ \|v\|_\infty &= \max(|v(:)|) \\ \|v\|_\infty &\leq \|v\|_2 \leq \sqrt{m} \|v\|_\infty \end{aligned}$$

If A is an m -by- n matrix and $QR = A$ is the economy-size QR decomposition of A , where Q is orthogonal and m -by- n and R is upper-triangular and n -by- n , then the singular values of R are equal to the singular values of A . If A is nonsingular, then

$$\|R^{-1}\|_2 = \|(R')^{-1}\|_2 = \frac{1}{\min(\text{svd}(R))} = \frac{1}{\min(\text{svd}(A))}$$

Upper Bound for $R = Q'A$

The upper bound for the magnitude of the elements of R is

$$\max(|R(:)|) \leq \sqrt{m} \max(|A(:)|).$$

Proof of Upper Bound for $R = Q'A$

The j th column of R is equal to $R(:, j) = Q'A(:, j)$, so

$$\begin{aligned} \max(|R(:, j)|) &= \|R(:, j)\|_\infty \\ &\leq \|R(:, j)\|_2 \\ &= \|Q'A(:, j)\|_2 \\ &\leq \|Q'\|_2 \|A(:, j)\|_2 \\ &= \|A(:, j)\|_2 \\ &\leq \sqrt{m} \|A(:, j)\|_\infty \\ &= \sqrt{m} \max(|A(:, j)|) \\ &\leq \sqrt{m} \max(|A(:)|). \end{aligned}$$

Since $\max(|R(:, j)|) \leq \sqrt{m} \max(|A(:)|)$ for all $1 \leq j$, then

$$\max(|R(:)|) \leq \sqrt{m} \max(|A(:)|).$$

Upper Bound for $X = (A'A) \setminus B$

The upper bound for the magnitude of the elements of $X = (A'A) \setminus B$ is

$$\max(|X(:)|) \leq \frac{\sqrt{n} \max(|B(:)|)}{\min(\text{svd}(A))^2}.$$

Proof of Upper Bound for $X = (A'A) \setminus B$

If A is not full rank, then $\min(\text{svd}(A)) = 0$, and if B is not equal to zero, then

$$\sqrt{n} \max(|B(:)|) / \min(\text{svd}(A))^2 = \infty \text{ and so the inequality is true.}$$

If $A'Ax = b$ and $QR = A$ is the economy-size QR decomposition of A , then $A'Ax = R'Q'QRx = R'Rx = b$.

If A is full rank then $x = R^{-1} \cdot ((R')^{-1}b)$. Let $x = X(:, j)$ be the j th column of X , and $b = B(:, j)$ be the j th column of B . Then

$$\begin{aligned}
\max(|x(:)|) &= \|x\|_\infty \\
&\leq \|x\|_2 \\
&= \|R^{-1} \cdot ((R')^{-1}b)\|_2 \\
&\leq \|R^{-1}\|_2 \|(R')^{-1}\|_2 \|b\|_2 \\
&= (1/\min(\text{svd}(A))^2) \cdot \|b\|_2 \\
&= \|b\|_2 / \min(\text{svd}(A))^2 \\
&\leq \sqrt{n} \|b\|_\infty / \min(\text{svd}(A))^2 \\
&= \sqrt{n} \max(|b(:)|) / \min(\text{svd}(A))^2.
\end{aligned}$$

Since $\max(|x(:)|) \leq \sqrt{n} \max(|b(:)|) / \min(\text{svd}(A))^2$ for all rows and columns of B and X , then

$$\max(|X(:)|) \leq \frac{\sqrt{n} \max(|B(:)|)}{\min(\text{svd}(A))^2}.$$

Lower Bound for $\min(\text{svd}(A))$

You can estimate a lower bound s of $\min(\text{svd}(A))$ for complex-valued A using the following formula,

$$s = \frac{\sigma_N}{\sqrt{2}} \sqrt{\gamma^{-1} \left(\frac{p_s \Gamma(m-n+2)^2 \Gamma(n)}{\Gamma(m+1) \Gamma(m-n+1) (m-n+1)}, m-n+1 \right)}$$

where σ_N is the standard deviation of random noise added to the elements of A , $1 - p_s$ is the probability that $s \leq \min(\text{svd}(A))$, Γ is the gamma function, and γ^{-1} is the inverse incomplete gamma function `gammaincinv`.

The proof is found in [1]. It is derived by integrating the formula in Lemma 3.4 from [3] and rearranging terms.

Since $s \leq \min(\text{svd}(A))$ with probability $1 - p_s$, then you can bound the magnitude of the elements of X without computing $\text{svd}(A)$,

$$\max(|X(:)|) \leq \frac{\sqrt{n} \max(|B(:)|)}{\min(\text{svd}(A))^2} \leq \frac{\sqrt{n} \max(|B(:)|)}{s^2} \text{ with probability } 1 - p_s.$$

You can compute s using the `fixed.complexSingularValueLowerBound` function which uses a default probability of 5 standard deviations below the mean,

$p_s = (1 + \text{erf}(-5/\sqrt{2}))/2 \approx 2.8665 \cdot 10^{-7}$, so the probability that the estimated bound for the smallest singular value s is less than the actual smallest singular value of A is $1 - p_s \approx 0.9999997$.

Example

This example runs a simulation with many random matrices and compares the analytical bounds with the actual singular values of A and the actual largest elements of $R = Q'A$, and $X = (A'A)B$.

Define System Parameters

Define the matrix attributes and system parameters for this example.

m is the number of rows in matrix A . In a problem such as beamforming or direction finding, m corresponds to the number of samples that are integrated over.

```
m = 300;
```

n is the number of columns in matrix A and rows in matrices B and X . In a least-squares problem, m is greater than n , and usually m is much larger than n . In a problem such as beamforming or direction finding, n corresponds to the number of sensors.

```
n = 10;
```

p is the number of columns in matrices B and X . It corresponds to simultaneously solving a system with p right-hand sides.

```
p = 1;
```

In this example, set the rank of matrix A to be less than the number of columns. In a problem such as beamforming or direction finding, $\text{rank}(A)$ corresponds to the number of signals impinging on the sensor array.

```
rankA = 3;
```

`precisionBits` defines the number of bits of precision required for the matrix solve. Set this value according to system requirements.

```
precisionBits = 24;
```

In this example, complex-valued matrices A and B are constructed such that the magnitude of the real and imaginary parts of their elements is less than or equal to one, so the maximum possible absolute value of any element is $|1 + 1i| = \sqrt{2}$. Your own system requirements will define what those values are. If you don't know what they are, and A and B are fixed-point inputs to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of A and B .

`max_abs_A` is an upper bound on the maximum magnitude element of A .

```
max_abs_A = sqrt(2);
```

`max_abs_B` is an upper bound on the maximum magnitude element of B .

```
max_abs_B = sqrt(2);
```

Thermal noise standard deviation is the square root of thermal noise power, which is a system parameter. A well-designed system has the quantization level lower than the thermal noise. Here, set `thermalNoiseStandardDeviation` to the equivalent of -50dB noise power.

```
thermalNoiseStandardDeviation = sqrt(10^(-50/10))
```

```
thermalNoiseStandardDeviation = 0.0032
```

The standard deviation of the noise from quantizing the real and imaginary parts of a complex signal is $2^{-\text{precisionBits}}/\sqrt{6}$ [4,5]. Use `fixed.complexQuantizationNoiseStandardDeviation` to compute this. See that it is less than `thermalNoiseStandardDeviation`.

```
quantizationNoiseStandardDeviation = fixed.complexQuantizationNoiseStandardDeviation(precisionBits)
```

```
quantizationNoiseStandardDeviation = 2.4333e-08
```

Compute Fixed-Point Types

In this example, assume that the designed system matrix A does not have full rank (there are fewer signals of interest than number of columns of matrix A), and the measured system matrix A has additive thermal noise that is larger than the quantization noise. The additive noise makes the measured matrix A have full rank.

Set $\sigma_{\text{noise}} = \sigma_{\text{thermal noise}}$.

```
noiseStandardDeviation = thermalNoiseStandardDeviation;
```

Use `fixed.complexQlessQRMatrixSolveFixedpointTypes` to compute fixed-point types.

```
T = fixed.complexQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,noiseStandardDeviation)
```

```
T = struct with fields:
    A: [0x0 embedded.fi]
    B: [0x0 embedded.fi]
    X: [0x0 embedded.fi]
```

$T.A$ is the type computed for transforming A to R in-place so that it does not overflow.

$T.A$

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 32
    FractionLength: 24
```

$T.B$ is the type computed for B so that it does not overflow.

$T.B$

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 27
    FractionLength: 24
```

$T.X$ is the type computed for the solution $X = (A'A)\backslash B$ so that there is a low probability that it overflows.

$T.X$

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
```

```
WordLength: 40
FractionLength: 24
```

Upper Bound for R

The upper bound for R is computed using the formula $\max(|R(:)|) \leq \sqrt{m} \max(|A(:)|)$, where m is the number of rows of matrix A . This upper bound is used to select a fixed-point type with the required number of bits of precision to avoid an overflow in the upper bound.

```
upperBoundR = sqrt(m)*max_abs_A
```

```
upperBoundR = 24.4949
```

Lower Bound for $\min(\text{svd}(A))$ for Complex A

A lower bound for $\min(\text{svd}(A))$ is estimated by the `fixed.complexSingularValueLowerBound` function using a probability that the estimate s is not greater than the actual smallest singular value. The default probability is 5 standard deviations below the mean. You can change this probability by specifying it as the last input parameter to the `fixed.complexSingularValueLowerBound` function.

```
estimatedSingularValueLowerBound = fixed.complexSingularValueLowerBound(m,n,noiseStandardDeviation
```

```
estimatedSingularValueLowerBound = 0.0389
```

Simulate and Compare to the Computed Bounds

The bounds are within an order of magnitude of the simulated results. This is sufficient because the number of bits translates to a logarithmic scale relative to the range of values. Being within a factor of 10 is between 3 and 4 bits. This is a good starting point for specifying a fixed-point type. If you run the simulation for more samples, then it is more likely that the simulated results will be closer to the bound. This example uses a limited number of simulations so it doesn't take too long to run. For real-world system design, you should run additional simulations.

Define the number of samples, `numSamples`, over which to run the simulation.

```
numSamples = 1e4;
```

Run the simulation.

```
[actualMaxR,singularValues,X_values] = runSimulations(m,n,p,rankA,max_abs_A,max_abs_B,numSamples
    noiseStandardDeviation,T);
```

You can see that the upper bound on R compared to the measured simulation results of the maximum value of R over all runs is within an order of magnitude.

```
upperBoundR
```

```
upperBoundR = 24.4949
```

```
max(actualMaxR)
```

```
ans = 9.4990
```

Finally, see that the estimated lower bound of $\min(\text{svd}(A))$ compared to the measured simulation results of $\min(\text{svd}(A))$ over all runs is also within an order of magnitude.

```
estimatedSingularValueLowerBound
```

```

estimatedSingularValueLowerBound = 0.0389
actualSmallestSingularValue = min(singularValues,[],'all')
actualSmallestSingularValue = 0.0443

```

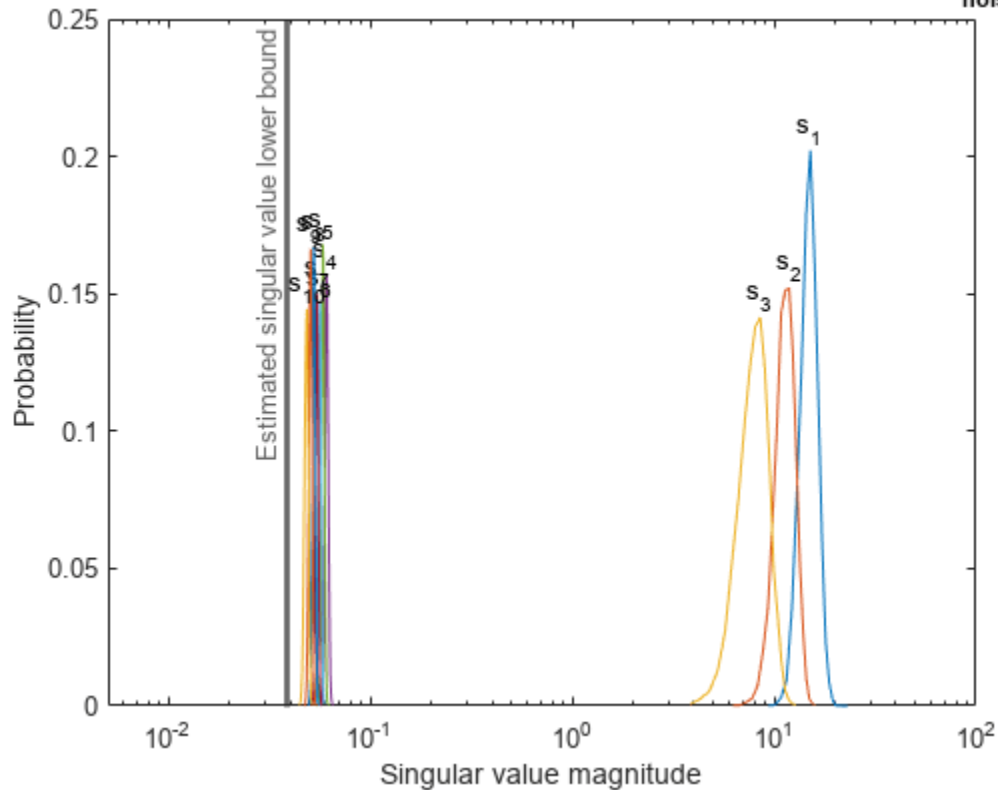
Plot the distribution of the singular values over all simulation runs. The distributions of the largest singular values correspond to the signals that determine the rank of the matrix. The distributions of the smallest singular values correspond to the noise. The derivation of the estimated bound of the smallest singular value makes use of the random nature of the noise.

```

clf
fixed.example.plot.singularValueDistribution(m,n,rankA,...
    noiseStandardDeviation,singularValues,...
    estimatedSingularValueLowerBound,"complex");

```

Singular value distributions for 300-by-10 complex matrices of rank 3 with $\sigma_{\text{noise}} = 0.1$

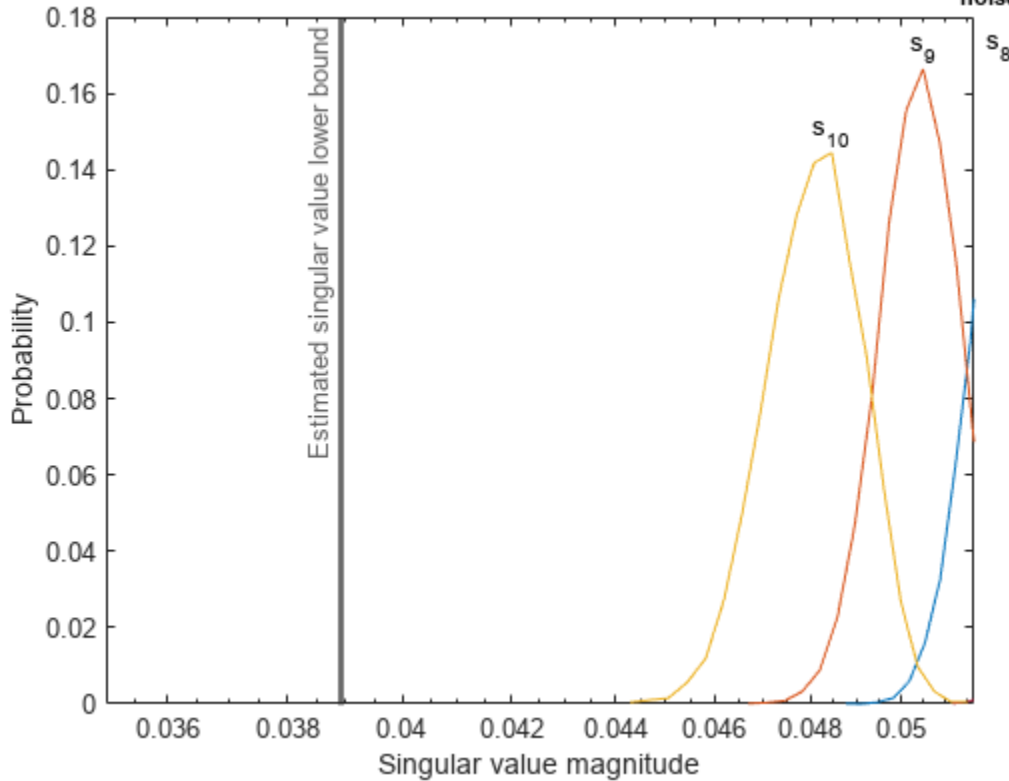


Zoom in to the smallest singular value to see that the estimated bound is close to it.

```

xlim([estimatedSingularValueLowerBound*0.9, max(singularValues(n,:))]);

```

Singular value distributions for 300-by-10 complex matrices of rank 3 with $\sigma_{\text{noise}} = 0.1$ 

Estimate the largest value of the solution, X , and compare it to the largest value of X found during the simulation runs. The estimation is within an order of magnitude of the actual value, which is sufficient for estimating a fixed-point data type, because it is between 3 and 4 bits.

This example uses a limited number of simulation runs. With additional simulation runs, the actual largest value of X will approach the estimated largest value of X .

```
estimated_largest_X = fixed.complexQlessQRMatrixSolveUpperBoundX(m,n,max_abs_B,noiseStandardDeviation)
```

```
estimated_largest_X = 9.3348e+03
```

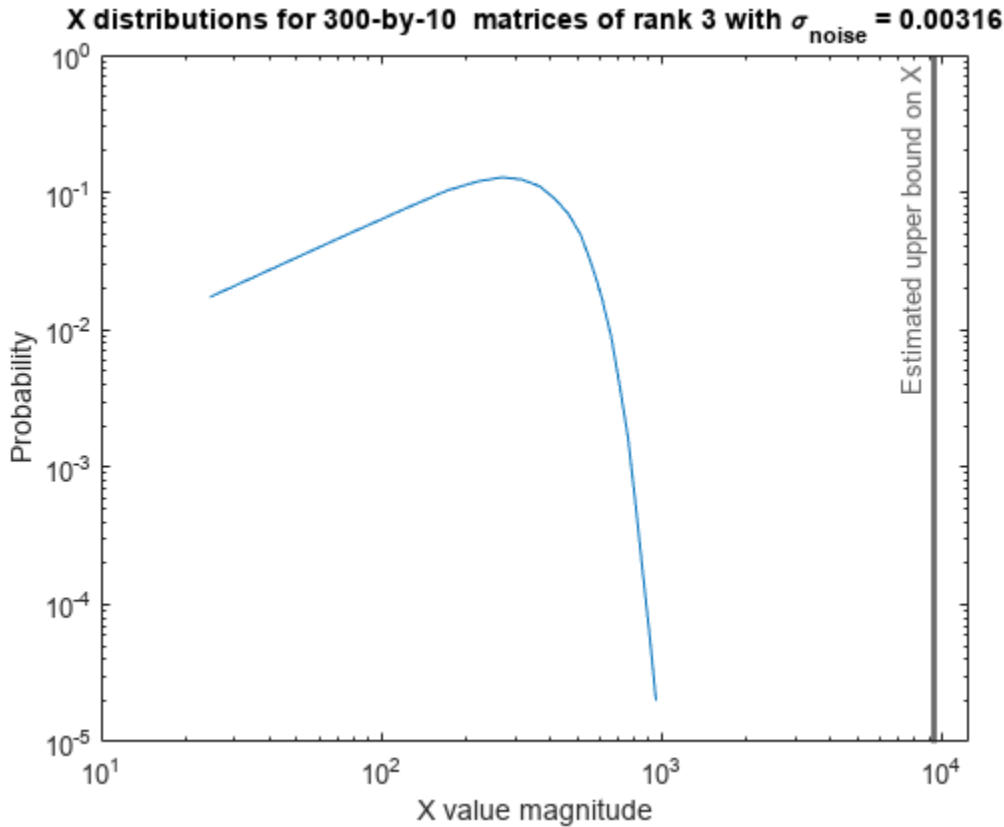
```
actual_largest_X = max(abs(X_values), [], 'all')
```

```
actual_largest_X = 977.7440
```

Plot the distribution of X values and compare it to the estimated upper bound for X .

```
clf
```

```
fixed.example.plot.xValueDistribution(m,n,rankA,noiseStandardDeviation,...
    X_values,estimated_largest_X,"complex normally distributed random");
```



Supporting Functions

The `runSimulations` function creates a series of random matrices A and B of a given size and rank, quantizes them according to the computed types, computes the QR decomposition of A , and solves the equation $A'AX = B$. It returns the maximum values of $R = Q'A$, the singular values of A , and the values of X so their distributions can be plotted and compared to the bounds.

```
function [actualMaxR,singularValues,X_values] = runSimulations(m,n,p,rankA,max_abs_A,max_abs_B,
    numSamples,noiseStandardDeviation,T)
precisionBits = T.A.FractionLength;
A_WordLength = T.A.WordLength;
B_WordLength = T.B.WordLength;
actualMaxR = zeros(1,numSamples);
singularValues = zeros(n,numSamples);
X_values = zeros(n,numSamples);
for j = 1:numSamples
    A = (max_abs_A/sqrt(2))*fixed.example.complexRandomLowRankMatrix(m,n,rankA);
    % Adding random noise makes A non-singular.
    A = A + fixed.example.complexNormalRandomArray(0,noiseStandardDeviation,m,n);
    A = quantizenumeric(A,1,A_WordLength,precisionBits);
    B = fixed.example.complexUniformRandomArray(-max_abs_B,max_abs_B,n,p);
    B = quantizenumeric(B,1,B_WordLength,precisionBits);
    [~,R] = qr(A,0);
    X = R \ (R' \ B);
    actualMaxR(j) = max(abs(R(:)));
    singularValues(:,j) = svd(A);
    X_values(:,j) = X;
end
```

end
end

References

- 1 Thomas A. Bryan and Jenna L. Warren. “Systems and Methods for Design Parameter Selection”. Patent pending. U.S. Patent Application No. 16/947,130. 2020.
- 2 “Perform QR Factorization Using CORDIC” on page 54-62. Derivation of the bound on growth when computing QR. MathWorks. 2010.
- 3 Zizhong Chen and Jack J. Dongarra. “Condition Numbers of Gaussian Random Matrices”. In: SIAM J. Matrix Anal. Appl. 27.3 (July 2005), pp. 603–620. issn: 0895-4798. doi: 10.1137/040616413. url: <https://dx.doi.org/10.1137/040616413>.
- 4 Bernard Widrow. “A Study of Rough Amplitude Quantization by Means of Nyquist Sampling Theory”. In: IRE Transactions on Circuit Theory 3.4 (Dec. 1956), pp. 266–276.
- 5 Bernard Widrow and István Kollár. Quantization Noise - Roundoff Error in Digital Computation, Signal Processing, Control, and Communications. Cambridge, UK: Cambridge University Press, 2008.
- 6 Gene H. Golub and Charles F. Van Loan. Matrix Computations. Second edition. Baltimore: Johns Hopkins University Press, 1989.

Suppress mlint warnings in this file.

```

%#ok< *NASGU>
%#ok< *ASGLU>

```

See Also

Functions

`fixed.complexQlessQRMatrixSolveFixedpointTypes` |
`fixed.complexSingularValueLowerBound` | `fixed.qlessQRMatrixSolve`

Blocks

Complex Burst Matrix Solve Using Q-less QR Decomposition | Complex Partial-Systolic Q-less QR Decomposition | Complex Partial-Systolic Q-less QR Decomposition with Forgetting Factor

Related Examples

- “Determine Fixed-Point Types for Complex Q-less QR Matrix Solve $A'AX=B$ ” on page 48-150

Determine Fixed-Point Types for Complex Q-less QR Matrix Solve $A'AX=B$

This example shows how to use the `fixed.complexQlessQRMatrixSolveFixedpointTypes` function to analytically determine fixed-point types for the solution of the complex least-squares matrix equation $A'AX = B$, where A is an m -by- n matrix with $m \geq n$, B is n -by- p , and X is n -by- p .

Fixed-point types for the solution of the matrix equation $A'AX = B$ are well-bounded if the number of rows, m , of A are much greater than the number of columns, n (i.e. $m \gg n$), and A is full rank. If A is not inherently full rank, then it can be made so by adding random noise. Random noise naturally occurs in physical systems, such as thermal noise in radar or communications systems. If $m = n$, then the dynamic range of the system can be unbounded, for example in the scalar equation $x = a/b$ and $a, b \in [-1, 1]$, then x can be arbitrarily large if b is close to 0.

Define System Parameters

Define the matrix attributes and system parameters for this example.

m is the number of rows in matrix A . In a problem such as beamforming or direction finding, m corresponds to the number of samples that are integrated over.

```
m = 300;
```

n is the number of columns in matrix A and rows in matrices B and X . In a least-squares problem, m is greater than n , and usually m is much larger than n . In a problem such as beamforming or direction finding, n corresponds to the number of sensors.

```
n = 10;
```

p is the number of columns in matrices B and X . It corresponds to simultaneously solving a system with p right-hand sides.

```
p = 1;
```

In this example, set the rank of matrix A to be less than the number of columns. In a problem such as beamforming or direction finding, $\text{rank}(A)$ corresponds to the number of signals impinging on the sensor array.

```
rankA = 3;
```

`precisionBits` defines the number of bits of precision required for the matrix solve. Set this value according to system requirements.

```
precisionBits = 24;
```

In this example, complex-valued matrices A and B are constructed such that the magnitude of the real and imaginary parts of their elements is less than or equal to one, so the maximum possible absolute value of any element is $|1 + 1i| = \sqrt{2}$. Your own system requirements will define what those values are. If you don't know what they are, and A and B are fixed-point inputs to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of A and B .

`max_abs_A` is an upper bound on the maximum magnitude element of A .

```
max_abs_A = sqrt(2);
```


`max_abs_B` is an upper bound on the maximum magnitude element of B .

```
max_abs_B = sqrt(2);
```

Thermal noise standard deviation is the square root of thermal noise power, which is a system parameter. A well-designed system has the quantization level lower than the thermal noise. Here, set `thermalNoiseStandardDeviation` to the equivalent of -50dB noise power.

```
thermalNoiseStandardDeviation = sqrt(10^(-50/10))
```

```
thermalNoiseStandardDeviation = 0.0032
```

The quantization noise standard deviation is a function of the required number of bits of precision. Use `fixed.complexQuantizationNoiseStandardDeviation` to compute this. See that it is less than `thermalNoiseStandardDeviation`.

```
quantizationNoiseStandardDeviation = fixed.complexQuantizationNoiseStandardDeviation(precisionBits)
```

```
quantizationNoiseStandardDeviation = 2.4333e-08
```

Compute Fixed-Point Types

In this example, assume that the designed system matrix A does not have full rank (there are fewer signals of interest than number of columns of matrix A), and the measured system matrix A has additive thermal noise that is larger than the quantization noise. The additive noise makes the measured matrix A have full rank.

Set $\sigma_{\text{noise}} = \sigma_{\text{thermal noise}}$.

```
noiseStandardDeviation = thermalNoiseStandardDeviation;
```

Use the `fixed.complexQlessQRMatrixSolveFixedpointTypes` function to compute fixed-point types.

```
T = fixed.complexQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,noiseStandardDeviation)
```

```
T = struct with fields:
    A: [0x0 embedded.fi]
    B: [0x0 embedded.fi]
    X: [0x0 embedded.fi]
```

`T.A` is the type computed for transforming A to $R = Q'A$ in-place so that it does not overflow.

```
T.A
```

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 32
    FractionLength: 24
```

`T.B` is the type computed for B so that it does not overflow.

```
T.B
```

```
ans =
[]

      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 27
      FractionLength: 24
```

T.X is the type computed for the solution $X = (A'A)\backslash B$ so that there is a low probability that it overflows.

```
T.X
ans =
[]

      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 40
      FractionLength: 24
```

Use the Specified Types to Solve the Matrix Equation $A'AX=B$

Create random matrices A and B such that $\text{rank}A=\text{rank}(A)$. Add random measurement noise to A which will make it become full rank.

```
rng('default');
[A,B] = fixed.example.complexRandomQlessQRMatrices(m,n,p,rankA);
A = A + fixed.example.complexNormalRandomArray(0,noiseStandardDeviation,m,n);
```

Cast the inputs to the types determined by `fixed.complexQlessQRMatrixSolveFixedpointTypes`. Quantizing to fixed-point is equivalent to adding random noise.

```
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
```

Accelerate the `fixed.qlessQRMatrixSolve` function by using `fiaccl` to generate a MATLAB executable (MEX) function.

```
fiaccl fixed.qlessQRMatrixSolve -args {A,B,T.X} -o qlessQRMatrixSolve_mex
```

Specify output type T.X and compute fixed-point $X = (A'A)\backslash B$ using the QR method.

```
X = qlessQRMatrixSolve_mex(A,B,T.X);
```

Compute the relative error to verify the accuracy of the output.

```
relative_error = norm(double(A'*A*X - B))/norm(double(B))
relative_error = 0.1054
```

Suppress `mlint` warnings in this file.

```
%#ok<*NASGU>  
%#ok<*ASGLU>
```

See Also

Functions

`fixed.complexQlessQRMatrixSolveFixedpointTypes` | `fixed.qlessQRMatrixSolve`

Blocks

Complex Burst Matrix Solve Using Q-less QR Decomposition | Complex Partial-Systolic Q-less QR Decomposition | Complex Partial-Systolic Q-less QR Decomposition with Forgetting Factor

Related Examples

- “Algorithms to Determine Fixed-Point Types for Complex Q-less QR Matrix Solve $A'AX=B$ ” on page 48-140

Algorithms to Determine Fixed-Point Types for Complex Least-Squares Matrix Solve $AX=B$

This example shows the algorithms that the `fixed.complexQRMatrixSolveFixedpointTypes` function uses to analytically determine fixed-point types for the solution of the complex least-squares matrix equation $AX = B$, where A is an m -by- n matrix with $m \geq n$, B is m -by- p , and X is n -by- p .

Overview

You can solve the fixed-point least-squares matrix equation $AX = B$ using QR decomposition. Using a sequence of orthogonal transformations, QR decomposition transforms matrix A in-place to upper triangular R , and transforms matrix B in-place to $C = Q'B$, where $QR = A$ is the economy-size QR decomposition. This reduces the equation to an upper-triangular system of equations $RX = C$. To solve for X , compute $X = R \setminus C$ through back-substitution of R into C .

You can determine appropriate fixed-point types for the least-squares matrix equation $AX = B$ by selecting the fraction length based on the number of bits of precision defined by your requirements. The `fixed.complexQRMatrixSolveFixedpointTypes` function analytically computes the following upper bounds on $R = Q'A$, $C = Q'B$, and X to determine the number of integer bits required to avoid overflow [1,2,3].

The upper bound for the magnitude of the elements of $R = Q'A$ is

$$\max(|R(:)|) \leq \sqrt{m} \max(|A(:)|).$$

The upper bound for the magnitude of the elements of $C = Q'B$ is

$$\max(|C(:)|) \leq \sqrt{m} \max(|B(:)|).$$

The upper bound for the magnitude of the elements of $X = A \setminus B$ is

$$\max(|X(:)|) \leq \frac{\sqrt{m} \max(|B(:)|)}{\min(\text{svd}(A))}.$$

Since computing `svd(A)` is more computationally expensive than solving the system of equations, the `fixed.complexQRMatrixSolveFixedpointTypes` function estimates a lower bound of `min(svd(A))`.

Fixed-point types for the solution of the matrix equation $AX = B$ are generally well-bounded if the number of rows, m , of A are much greater than the number of columns, n (i.e. $m \gg n$), and A is full rank. If A is not inherently full rank, then it can be made so by adding random noise. Random noise naturally occurs in physical systems, such as thermal noise in radar or communications systems. If $m = n$, then the dynamic range of the system can be unbounded, for example in the scalar equation $x = a/b$ and $a, b \in [-1, 1]$, then x can be arbitrarily large if b is close to 0.

Proofs of the Bounds

Properties and Definitions of Vector and Matrix Norms

The proofs of the bounds use the following properties and definitions of matrix and vector norms, where Q is an orthogonal matrix, and v is a vector of length m [6].

$$\begin{aligned} \|Av\|_2 &\leq \|A\|_2 \|v\|_2 \\ \|Q\|_2 &= 1 \\ \|v\|_\infty &= \max(|v(:)|) \\ \|v\|_\infty &\leq \|v\|_2 \leq \sqrt{m} \|v\|_\infty \end{aligned}$$

If A is an m -by- n matrix and $QR = A$ is the economy-size QR decomposition of A , where Q is orthogonal and m -by- n and R is upper-triangular and n -by- n , then the singular values of R are equal to the singular values of A . If A is nonsingular, then

$$\|R^{-1}\|_2 = \|(R')^{-1}\|_2 = \frac{1}{\min(\text{svd}(R))} = \frac{1}{\min(\text{svd}(A))}$$

Upper Bound for $R = Q'A$

The upper bound for the magnitude of the elements of R is

$$\max(|R(:)|) \leq \sqrt{m} \max(|A(:)|).$$

Proof of Upper Bound for $R = Q'A$

The j th column of R is equal to $R(:, j) = Q'A(:, j)$, so

$$\begin{aligned} \max(|R(:, j)|) &= \|R(:, j)\|_\infty \\ &\leq \|R(:, j)\|_2 \\ &= \|Q'A(:, j)\|_2 \\ &\leq \|Q'\|_2 \|A(:, j)\|_2 \\ &= \|A(:, j)\|_2 \\ &\leq \sqrt{m} \|A(:, j)\|_\infty \\ &= \sqrt{m} \max(|A(:, j)|) \\ &\leq \sqrt{m} \max(|A(:)|). \end{aligned}$$

Since $\max(|R(:, j)|) \leq \sqrt{m} \max(|A(:)|)$ for all $1 \leq j$, then

$$\max(|R(:)|) \leq \sqrt{m} \max(|A(:)|).$$

Upper Bound for $C = Q'B$

The upper bound for the magnitude of the elements of $C = Q'B$ is

$$\max(|C(:)|) \leq \sqrt{m} \max(|B(:)|).$$

Proof of Upper Bound for $C = Q'B$

The proof of the upper bound for $C = Q'B$ is the same as the proof of the upper bound for $R = Q'A$ by substituting C for R and B for A .

Upper Bound for $X = A \setminus B$

The upper bound for the magnitude of the elements of $X = A \setminus B$ is

$$\max(|X(:)|) \leq \frac{\sqrt{m}\max(|B(:)|)}{\min(\text{svd}(A))}.$$

Proof of Upper Bound for $X = A \setminus B$

If A is not full rank, then $\min(\text{svd}(A)) = 0$, and if B is not equal to zero, then $\sqrt{m}\max(|B(:)|)/\min(\text{svd}(A)) = \infty$ and so the inequality is true.

If A is full rank, then $x = R^{-1}(Q'b)$. Let $x = X(:, j)$ be the j th column of X , and $b = B(:, j)$ be the j th column of B . Then

$$\begin{aligned} \max(|x(:)|) &= \|x\|_{\infty} \\ &\leq \|x\|_2 \\ &= \|R^{-1} \cdot (Q'b)\|_2 \\ &\leq \|R^{-1}\|_2 \|Q'\|_2 \|b\|_2 \\ &= (1/\min(\text{svd}(A))) \cdot 1 \cdot \|b\|_2 \\ &= \|b\|_2 / \min(\text{svd}(A)) \\ &\leq \sqrt{m} \|b\|_{\infty} / \min(\text{svd}(A)) \\ &= \sqrt{m}\max(|b(:)|) / \min(\text{svd}(A)). \end{aligned}$$

Since $\max(|x(:)|) \leq \sqrt{m}\max(|b(:)|)/\min(\text{svd}(A))$ for all rows and columns of B and X , then

$$\max(|X(:)|) \leq \frac{\sqrt{m}\max(|B(:)|)}{\min(\text{svd}(A))}.$$

Lower Bound for $\min(\text{svd}(A))$

You can estimate a lower bound s of $\min(\text{svd}(A))$ for complex-valued A using the following formula,

$$s = \frac{\sigma_N}{\sqrt{2}} \sqrt{\gamma^{-1} \left(\frac{p_s \Gamma(m-n+2)^2 \Gamma(n)}{\Gamma(m+1) \Gamma(m-n+1) (m-n+1)}, m-n+1 \right)}$$

where σ_N is the standard deviation of random noise added to the elements of A , $1 - p_s$ is the probability that $s \leq \min(\text{svd}(A))$, Γ is the gamma function, and γ^{-1} is the inverse incomplete gamma function `gammaincinv`.

The proof is found in [1]. It is derived by integrating the formula in Lemma 3.4 from [3] and rearranging terms.

Since $s \leq \min(\text{svd}(A))$ with probability $1 - p_s$, then you can bound the magnitude of the elements of X without computing $\text{svd}(A)$,

$$\max(|X(:)|) \leq \frac{\sqrt{m}\max(|B(:)|)}{\min(\text{svd}(A))} \leq \frac{\sqrt{m}\max(|B(:)|)}{s} \text{ with probability } 1 - p_s.$$

You can compute s using the `fixed.complexSingularValueLowerBound` function which uses a default probability of 5 standard deviations below the mean,

$p_s = (1 + \text{erf}(-5/\sqrt{2}))/2 \approx 2.8665 \cdot 10^{-7}$, so the probability that the estimated bound for the smallest singular value s is less than the actual smallest singular value of A is $1 - p_s \approx 0.9999997$.

Example

This example runs a simulation with many random matrices and compares the analytical bounds with the actual singular values of A and the actual largest elements of $R = Q'A$, $C = Q'B$, and $X = A \setminus B$.

Define System Parameters

Define the matrix attributes and system parameters for this example.

m is the number of rows in matrices A and B . In a problem such as beamforming or direction finding, m corresponds to the number of samples that are integrated over.

```
m = 300;
```

n is the number of columns in matrix A and rows in matrix X . In a least-squares problem, m is greater than n , and usually m is much larger than n . In a problem such as beamforming or direction finding, n corresponds to the number of sensors.

```
n = 10;
```

p is the number of columns in matrices B and X . It corresponds to simultaneously solving a system with p right-hand sides.

```
p = 1;
```

In this example, set the rank of matrix A to be less than the number of columns. In a problem such as beamforming or direction finding, $\text{rank}(A)$ corresponds to the number of signals impinging on the sensor array.

```
rankA = 3;
```

`precisionBits` defines the number of bits of precision required for the matrix solve. Set this value according to system requirements.

```
precisionBits = 24;
```

In this example, complex-valued matrices A and B are constructed such that the magnitude of the real and imaginary parts of their elements is less than or equal to one, so the maximum possible absolute value of any element is $|1 + 1i| = \sqrt{2}$. Your own system requirements will define what those values are. If you don't know what they are, and A and B are fixed-point inputs to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of A and B .

`max_abs_A` is an upper bound on the maximum magnitude element of A .

```
max_abs_A = sqrt(2);
```

`max_abs_B` is an upper bound on the maximum magnitude element of B .

```
max_abs_B = sqrt(2);
```

Thermal noise standard deviation is the square root of thermal noise power, which is a system parameter. A well-designed system has the quantization level lower than the thermal noise. Here, set `thermalNoiseStandardDeviation` to the equivalent of -50dB noise power.

```
thermalNoiseStandardDeviation = sqrt(10^(-50/10))
```

```
thermalNoiseStandardDeviation = 0.0032
```

The standard deviation of the noise from quantizing the real and imaginary parts of a complex signal is $2^{-\text{precisionBits}}/\sqrt{6}$ [4,5]. Use the `fixed.complexQuantizationNoiseStandardDeviation` function to compute this. See that it is less than `thermalNoiseStandardDeviation`.

```
quantizationNoiseStandardDeviation = fixed.complexQuantizationNoiseStandardDeviation(precisionBits);
quantizationNoiseStandardDeviation = 2.4333e-08
```

Compute Fixed-Point Types

In this example, assume that the designed system matrix A does not have full rank (there are fewer signals of interest than number of columns of matrix A), and the measured system matrix A has additive thermal noise that is larger than the quantization noise. The additive noise makes the measured matrix A have full rank.

Set $\sigma_{\text{noise}} = \sigma_{\text{thermal noise}}$.

```
noiseStandardDeviation = thermalNoiseStandardDeviation;
```

Use `fixed.complexQRMatrixSolveFixedpointTypes` to compute fixed-point types.

```
T = fixed.complexQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,noiseStandardDeviation)
```

```
T = struct with fields:
    A: [0x0 embedded.fi]
    B: [0x0 embedded.fi]
    X: [0x0 embedded.fi]
```

$T.A$ is the type computed for transforming A to R in-place so that it does not overflow.

$T.A$

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 32
    FractionLength: 24
```

$T.B$ is the type computed for transforming B to $Q'B$ in-place so that it does not overflow.

$T.B$

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 32
    FractionLength: 24
```

$T.X$ is the type computed for the solution $X = A \setminus B$ so that there is a low probability that it overflows.

$T.X$


```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 37
    FractionLength: 24
```

Upper Bounds for R and $C=Q'B$

The upper bounds for R and $C = Q'B$ are computed using the following formulas, where m is the number of rows of matrices A and B .

$$\max(|R(:)|) \leq \sqrt{m} \max(|A(:)|)$$

$$\max(|C(:)|) \leq \sqrt{m} \max(|B(:)|)$$

These upper bounds are used to select a fixed-point type with the required number of bits of precision to avoid overflows.

```
upperBoundR = sqrt(m)*max_abs_A
```

```
upperBoundR = 24.4949
```

```
upperBoundQB = sqrt(m)*max_abs_B
```

```
upperBoundQB = 24.4949
```

Lower Bound for $\min(\text{svd}(A))$ for Complex A

A lower bound for $\min(\text{svd}(A))$ is estimated by the `fixed.complexSingularValueLowerBound` function using a probability that the estimate s is not greater than the actual smallest singular value. The default probability is 5 standard deviations below the mean. You can change this probability by specifying it as the last input parameter to the `fixed.complexSingularValueLowerBound` function.

```
estimatedSingularValueLowerBound = fixed.complexSingularValueLowerBound(m,n,noiseStandardDeviation,probability)
```

```
estimatedSingularValueLowerBound = 0.0389
```

Simulate and Compare to the Computed Bounds

The bounds are within an order of magnitude of the simulated results. This is sufficient because the number of bits translates to a logarithmic scale relative to the range of values. Being within a factor of 10 is between 3 and 4 bits. This is a good starting point for specifying a fixed-point type. If you run the simulation for more samples, then it is more likely that the simulated results will be closer to the bound. This example uses a limited number of simulations so it doesn't take too long to run. For real-world system design, you should run additional simulations.

Define the number of samples, `numSamples`, over which to run the simulation.

```
numSamples = 1e4;
```

Run the simulation.

```
[actualMaxR,actualMaxQB,singularValues,X_values] = runSimulations(m,n,p,rankA,max_abs_A,max_abs_B,rankB,rankC,numSamples,noiseStandardDeviation,T);
```

You can see that the upper bound on R compared to the measured simulation results of the maximum value of R over all runs is within an order of magnitude.

```
upperBoundR
upperBoundR = 24.4949
max(actualMaxR)
ans = 9.6720
```

You can see that the upper bound on $C = QB$ compared to the measured simulation results of the maximum value of $C = QB$ over all runs is also within an order of magnitude.

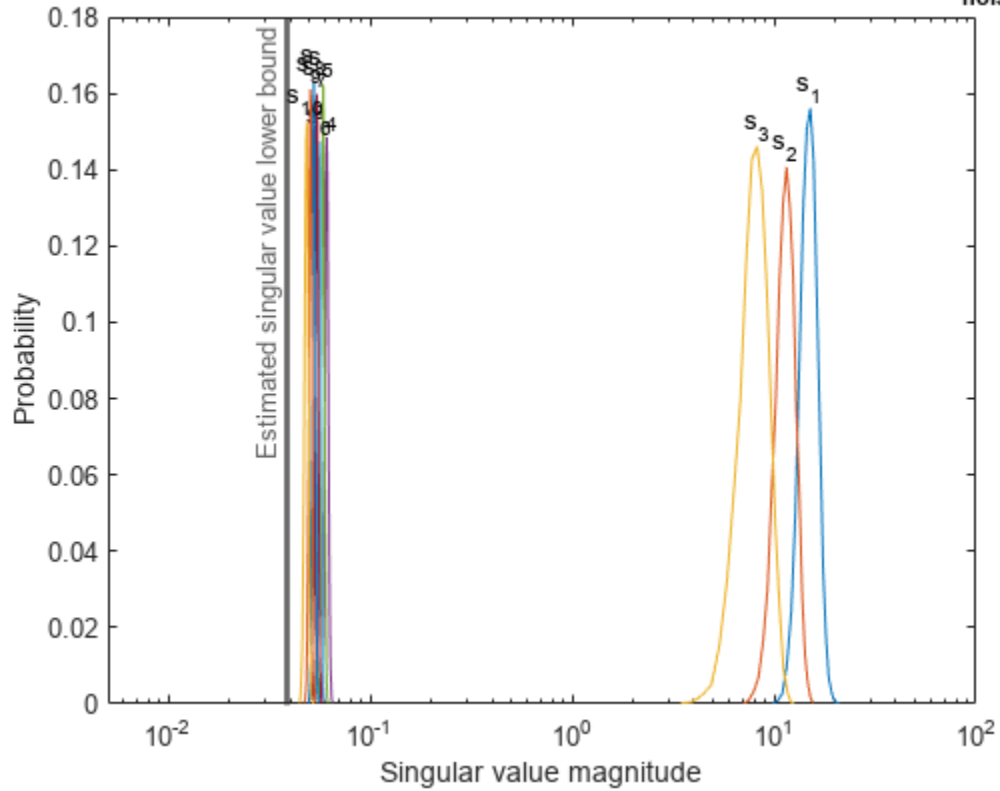
```
upperBoundQB
upperBoundQB = 24.4949
max(actualMaxQB)
ans = 4.4764
```

Finally, see that the estimated lower bound of $\min(\text{svd}(A))$ compared to the measured simulation results of $\min(\text{svd}(A))$ over all runs is also within an order of magnitude.

```
estimatedSingularValueLowerBound
estimatedSingularValueLowerBound = 0.0389
actualSmallestSingularValue = min(singularValues,[], 'all')
actualSmallestSingularValue = 0.0443
```

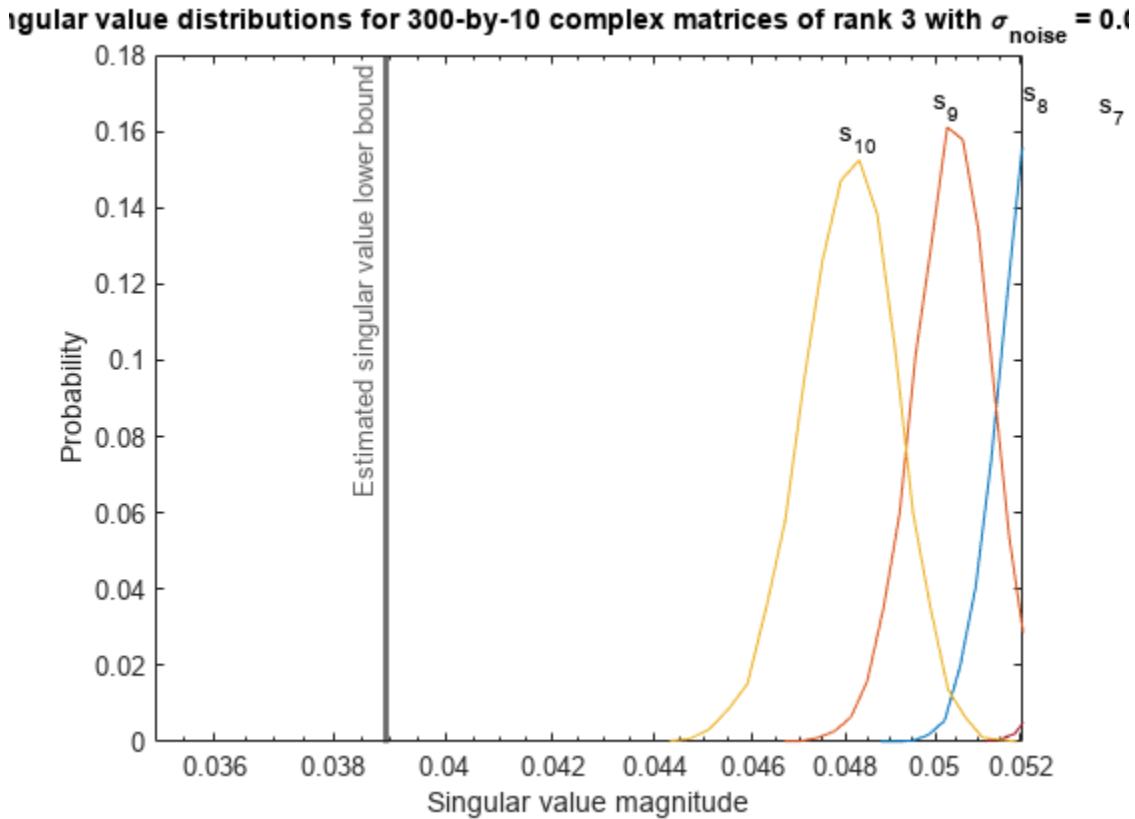
Plot the distribution of the singular values over all simulation runs. The distributions of the largest singular values correspond to the signals that determine the rank of the matrix. The distributions of the smallest singular values correspond to the noise. The derivation of the estimated bound of the smallest singular value makes use of the random nature of the noise.

```
clf
fixed.example.plot.singularValueDistribution(m,n,rankA,noiseStandardDeviation,...
    singularValues,estimatedSingularValueLowerBound,"complex");
```

Singular value distributions for 300-by-10 complex matrices of rank 3 with $\sigma_{\text{noise}} = 0.1$ 

Zoom in to the smallest singular value to see that the estimated bound is close to it.

```
xlim([estimatedSingularValueLowerBound*0.9, max(singularValues(n,:))]);
```



Estimate the largest value of the solution, X , and compare it to the largest value of X found during the simulation runs. The estimation is within an order of magnitude of the actual value, which is sufficient for estimating a fixed-point data type, because it is between 3 and 4 bits.

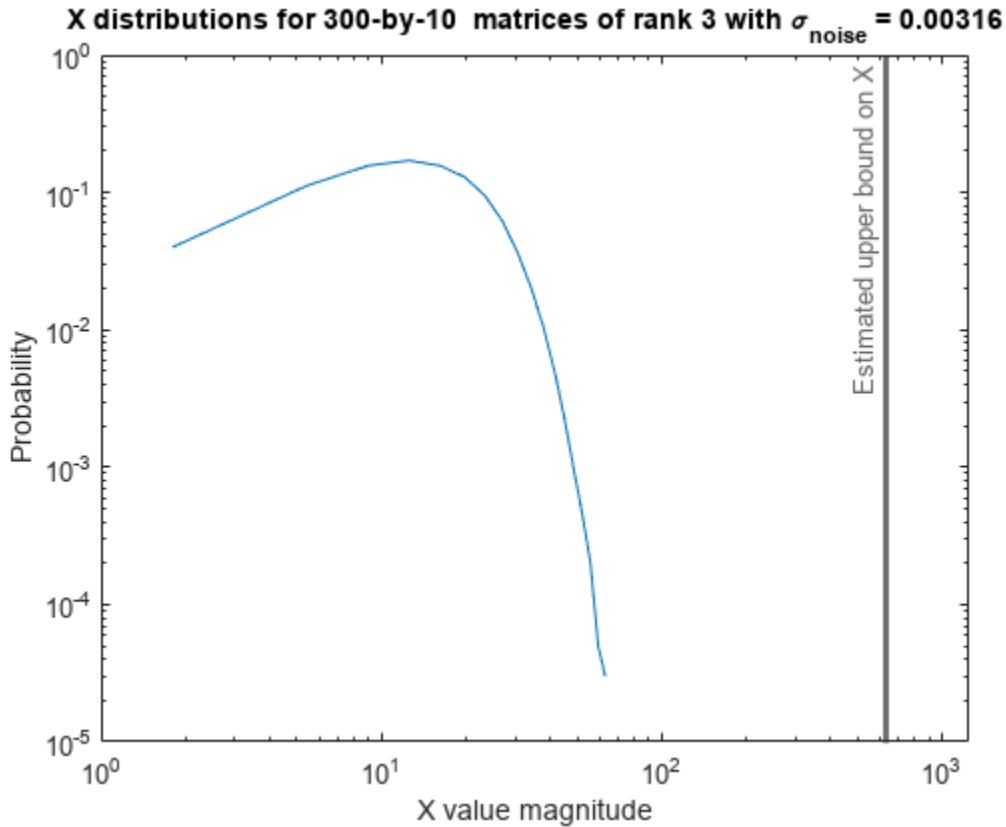
This example uses a limited number of simulation runs. With additional simulation runs, the actual largest value of X will approach the estimated largest value of X .

```
estimated_largest_X = fixed.complexMatrixSolveUpperBoundX(m,n,max_abs_B,noiseStandardDeviation)
estimated_largest_X = 629.3194
```

```
actual_largest_X = max(abs(X_values),[],'all')
actual_largest_X = 70.2644
```

Plot the distribution of X values and compare it to the estimated upper bound for X .

```
clf
fixed.example.plot.xValueDistribution(m,n,rankA,noiseStandardDeviation,...
    X_values,estimated_largest_X,"complex normally distributed random");
```



Supporting Functions

The `runSimulations` function creates a series of random matrices A and B of a given size and rank, quantizes them according to the computed types, computes the QR decomposition of A , and solves the equation $AX = B$. It returns the maximum values of $R = Q'A$ and $C = Q'B$, the singular values of A , and the values of X so their distributions can be plotted and compared to the bounds.

```
function [actualMaxR,actualMaxQB,singularValues,X_values] = runSimulations(m,n,p,rankA,max_abs_A,
    numSamples,noiseStandardDeviation,T)
precisionBits = T.A.FractionLength;
A_WordLength = T.A.WordLength;
B_WordLength = T.B.WordLength;
actualMaxR = zeros(1,numSamples);
actualMaxQB = zeros(1,numSamples);
singularValues = zeros(n,numSamples);
X_values = zeros(n,numSamples);
for j = 1:numSamples
    A = (max_abs_A/sqrt(2))*fixed.example.complexRandomLowRankMatrix(m,n,rankA);
    % Adding normally distributed random noise makes A non-singular.
    A = A + fixed.example.complexNormalRandomArray(0,noiseStandardDeviation,m,n);
    A = quantiznumeric(A,1,A_WordLength,precisionBits);
    B = fixed.example.complexUniformRandomArray(-max_abs_B,max_abs_B,m,p);
    B = quantiznumeric(B,1,B_WordLength,precisionBits);
    [Q,R] = qr(A,0);
    C = Q'*B;
    X = R\C;
    actualMaxR(j) = max(abs(R(:)));
end
```

```
        actualMaxQB(j) = max(abs(C(:)));  
        singularValues(:,j) = svd(A);  
        X_values(:,j) = X;  
    end  
end
```

References

- 1 Thomas A. Bryan and Jenna L. Warren. “Systems and Methods for Design Parameter Selection”. Patent pending. U.S. Patent Application No. 16/947,130. 2020.
- 2 Perform QR Factorization Using CORDIC. Derivation of the bound on growth when computing QR. MathWorks. 2010.
- 3 Zizhong Chen and Jack J. Dongarra. “Condition Numbers of Gaussian Random Matrices”. In: SIAM J. Matrix Anal. Appl. 27.3 (July 2005), pp. 603–620. issn: 0895-4798. doi: 10.1137/040616413. url: <https://dx.doi.org/10.1137/040616413>.
- 4 Bernard Widrow. “A Study of Rough Amplitude Quantization by Means of Nyquist Sampling Theory”. In: IRE Transactions on Circuit Theory 3.4 (Dec. 1956), pp. 266–276.
- 5 Bernard Widrow and István Kollár. Quantization Noise - Roundoff Error in Digital Computation, Signal Processing, Control, and Communications. Cambridge, UK: Cambridge University Press, 2008.
- 6 Gene H. Golub and Charles F. Van Loan. Matrix Computations. Second edition. Baltimore: Johns Hopkins University Press, 1989.

Suppress mlint warnings in this file.

```
 %#ok< *NASGU>  
 %#ok< *ASGLU>
```

See Also

Functions

`fixed.complexQRMatrixSolveFixedpointTypes` |
`fixed.complexSingularValueLowerBound` | `fixed.qrMatrixSolve`

Blocks

`Complex Burst Matrix Solve Using QR Decomposition` | `Complex Partial-Systolic Matrix Solve Using QR Decomposition`

Related Examples

- “Determine Fixed-Point Types for Complex Least-Squares Matrix Solve AX=B” on page 48-165

Determine Fixed-Point Types for Complex Least-Squares Matrix Solve AX=B

This example shows how to use the `fixed.complexQRMatrixSolveFixedpointTypes` function to analytically determine fixed-point types for the solution of the complex least-squares matrix equation $AX = B$, where A is an m -by- n matrix with $m \geq n$, B is m -by- p , and X is n -by- p .

Fixed-point types for the solution of the matrix equation $AX = B$ are well-bounded if the number of rows, m , of A are much greater than the number of columns, n (i.e. $m \gg n$), and A is full rank. If A is not inherently full rank, then it can be made so by adding random noise. Random noise naturally occurs in physical systems, such as thermal noise in radar or communications systems. If $m = n$, then the dynamic range of the system can be unbounded, for example in the scalar equation $x = a/b$ and $a, b \in [-1, 1]$, then x can be arbitrarily large if b is close to 0.

Define System Parameters

Define the matrix attributes and system parameters for this example.

m is the number of rows in matrices A and B . In a problem such as beamforming or direction finding, m corresponds to the number of samples that are integrated over.

```
m = 300;
```

n is the number of columns in matrix A and rows in matrix X . In a least-squares problem, m is greater than n , and usually m is much larger than n . In a problem such as beamforming or direction finding, n corresponds to the number of sensors.

```
n = 10;
```

p is the number of columns in matrices B and X . It corresponds to simultaneously solving a system with p right-hand sides.

```
p = 1;
```

In this example, set the rank of matrix A to be less than the number of columns. In a problem such as beamforming or direction finding, $\text{rank}(A)$ corresponds to the number of signals impinging on the sensor array.

```
rankA = 3;
```

`precisionBits` defines the number of bits of precision required for the matrix solve. Set this value according to system requirements.

```
precisionBits = 24;
```

In this example, complex-valued matrices A and B are constructed such that the magnitude of the real and imaginary parts of their elements is less than or equal to one, so the maximum possible absolute value of any element is $|1 + 1i| = \sqrt{2}$. Your own system requirements will define what those values are. If you don't know what they are, and A and B are fixed-point inputs to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of A and B .

`max_abs_A` is an upper bound on the maximum magnitude element of A .

```
max_abs_A = sqrt(2);
```

`max_abs_B` is an upper bound on the maximum magnitude element of `B`.

```
max_abs_B = sqrt(2);
```

Thermal noise standard deviation is the square root of thermal noise power, which is a system parameter. A well-designed system has the quantization level lower than the thermal noise. Here, set `thermalNoiseStandardDeviation` to the equivalent of -50dB noise power.

```
thermalNoiseStandardDeviation = sqrt(10^(-50/10))
```

```
thermalNoiseStandardDeviation = 0.0032
```

The quantization noise standard deviation is a function of the required number of bits of precision. Use `fixed.complexQuantizationNoiseStandardDeviation` to compute this. See that it is less than `thermalNoiseStandardDeviation`.

```
quantizationNoiseStandardDeviation = fixed.complexQuantizationNoiseStandardDeviation(precisionBits)
```

```
quantizationNoiseStandardDeviation = 2.4333e-08
```

Compute Fixed-Point Types

In this example, assume that the designed system matrix `A` does not have full rank (there are fewer signals of interest than number of columns of matrix `A`), and the measured system matrix `A` has additive thermal noise that is larger than the quantization noise. The additive noise makes the measured matrix `A` have full rank.

Set $\sigma_{\text{noise}} = \sigma_{\text{thermal noise}}$.

```
noiseStandardDeviation = thermalNoiseStandardDeviation;
```

Use `fixed.complexQRMatrixSolveFixedpointTypes` to compute fixed-point types.

```
T = fixed.complexQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,noiseStandardDeviation)
```

```
T = struct with fields:
    A: [0x0 embedded.fi]
    B: [0x0 embedded.fi]
    X: [0x0 embedded.fi]
```

`T.A` is the type computed for transforming `A` to $R = QA$ in-place so that it does not overflow.

```
T.A
```

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 32
    FractionLength: 24
```

`T.B` is the type computed for transforming `B` to $C = QB$ in-place so that it does not overflow.

```
T.B
```



```
ans =
[]
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 32
    FractionLength: 24
```

T.X is the type computed for the solution $X = A \setminus B$ so that there is a low probability that it overflows.

```
T.X
ans =
[]
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 37
    FractionLength: 24
```

Use the Specified Types to Solve the Matrix Equation $AX=B$

Create random matrices A and B such that B is in the range of A, and $\text{rankA}=\text{rank}(A)$. Add random measurement noise to A which will make it become full rank, but it will also affect the solution so that B is only close to the range of A.

```
rng('default');
[A,B] = fixed.example.complexRandomLeastSquaresMatrices(m,n,p,rankA);
A = A + fixed.example.complexNormalRandomArray(0,noiseStandardDeviation,m,n);
```

Cast the inputs to the types determined by `fixed.complexQRMatrixSolveFixedpointTypes`. Quantizing to fixed-point is equivalent to adding random noise.

```
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
```

Accelerate the `fixed.qrMatrixSolve` function by using `fiaccl` to generate a MATLAB executable (MEX) function.

```
fiaccl fixed.qrMatrixSolve -args {A,B,T.X} -o qrComplexMatrixSolve_mex
```

Specify the output type T.X and compute fixed-point $X = A \setminus B$ using the QR method.

```
X = qrComplexMatrixSolve_mex(A,B,T.X);
```

Compute the relative error to verify the accuracy of the output.

```
relative_error = norm(double(A*X - B))/norm(double(B))
```

```
relative_error = 0.0056
```

Suppress `mlint` warnings in this file.

```
 %#ok<*NASGU>  
 %#ok<*ASGLU>
```

See Also

Functions

`fixed.complexQRMatrixSolveFixedpointTypes` | `fixed.qrMatrixSolve`

Blocks

Complex Burst Matrix Solve Using QR Decomposition | Complex Partial-Systolic Matrix Solve Using QR Decomposition

Related Examples

- “Algorithms to Determine Fixed-Point Types for Complex Least-Squares Matrix Solve $AX=B$ ” on page 48-154

Algorithms to Determine Fixed-Point Types for Real Q-less QR Matrix Solve $A'AX=B$

This example shows the algorithms that the `fixed.realQlessQRMatrixSolveFixedpointTypes` function uses to analytically determine fixed-point types for the solution of the real matrix equation $A'AX = B$, where A is an m -by- n matrix with $m > n$, B is n -by- p , and X is n -by- p .

Overview

You can solve the fixed-point matrix equation $A'AX = B$ using QR decomposition. Using a sequence of orthogonal transformations, QR decomposition transforms matrix A in-place to upper triangular R , where $QR = A$ is the economy-size QR decomposition. This reduces the equation to an upper-triangular system of equations $R'RX = B$. To solve for X , compute $X = R \setminus (R' \setminus B)$ through forward- and backward-substitution of R into B .

You can determine appropriate fixed-point types for the matrix equation $A'AX = B$ by selecting the fraction length based on the number of bits of precision defined by your requirements. The `fixed.realQlessQRMatrixSolveFixedpointTypes` function analytically computes the following upper bounds on R , and X to determine the number of integer bits required to avoid overflow [1,2,3].

The upper bound for the magnitude of the elements of $R = Q'A$ is

$$\max(|R(:)|) \leq \sqrt{m} \max(|A(:)|).$$

The upper bound for the magnitude of the elements of $X = (A'A) \setminus B$ is

$$\max(|X(:)|) \leq \frac{\sqrt{n} \max(|B(:)|)}{\min(\text{svd}(A))^2}.$$

Since computing `svd(A)` is more computationally expensive than solving the system of equations, the `fixed.realQlessQRMatrixSolveFixedpointTypes` function estimates a lower bound of `min(svd(A))`.

Fixed-point types for the solution of the matrix equation $(A'A)X = B$ are generally well-bounded if the number of rows, m , of A are much greater than the number of columns, n (i.e. $m \gg n$), and A is full rank. If A is not inherently full rank, then it can be made so by adding random noise. Random noise naturally occurs in physical systems, such as thermal noise in radar or communications systems. If $m = n$, then the dynamic range of the system can be unbounded, for example in the scalar equation $x = a^2/b$ and $a, b \in [-1, 1]$, then x can be arbitrarily large if b is close to 0.

Proofs of the Bounds

Properties and Definitions of Vector and Matrix Norms

The proofs of the bounds use the following properties and definitions of matrix and vector norms, where Q is an orthogonal matrix, and v is a vector of length m [6].

$$\|Av\|_2 \leq \|A\|_2 \|v\|_2$$

$$\|Q\|_2 = 1$$

$$\|v\|_\infty = \max(|v(:)|)$$

$$\|v\|_\infty \leq \|v\|_2 \leq \sqrt{m} \|v\|_\infty$$

If A is an m -by- n matrix and $QR = A$ is the economy-size QR decomposition of A , where Q is orthogonal and m -by- n and R is upper-triangular and n -by- n , then the singular values of R are equal to the singular values of A . If A is nonsingular, then

$$\|R^{-1}\|_2 = \|(R')^{-1}\|_2 = \frac{1}{\min(\text{svd}(R))} = \frac{1}{\min(\text{svd}(A))}$$

Upper Bound for $R = Q'A$

The upper bound for the magnitude of the elements of R is

$$\max(|R(\cdot)|) \leq \sqrt{m} \max(|A(\cdot)|).$$

Proof of Upper Bound for $R = Q'A$

The j th column of R is equal to $R(:, j) = Q'A(:, j)$, so

$$\begin{aligned} \max(|R(:, j)|) &= \|R(:, j)\|_\infty \\ &\leq \|R(:, j)\|_2 \\ &= \|Q'A(:, j)\|_2 \\ &\leq \|Q'\|_2 \|A(:, j)\|_2 \\ &= \|A(:, j)\|_2 \\ &\leq \sqrt{m} \|A(:, j)\|_\infty \\ &= \sqrt{m} \max(|A(:, j)|) \\ &\leq \sqrt{m} \max(|A(\cdot)|). \end{aligned}$$

Since $\max(|R(:, j)|) \leq \sqrt{m} \max(|A(\cdot)|)$ for all $1 \leq j$, then

$$\max(|R(\cdot)|) \leq \sqrt{m} \max(|A(\cdot)|).$$

Upper Bound for $X = (A'A) \setminus B$

The upper bound for the magnitude of the elements of $X = (A'A) \setminus B$ is

$$\max(|X(\cdot)|) \leq \frac{\sqrt{n} \max(|B(\cdot)|)}{\min(\text{svd}(A))^2}.$$

Proof of Upper Bound for $X = (A'A) \setminus B$

If A is not full rank, then $\min(\text{svd}(A)) = 0$, and if B is not equal to zero, then

$$\sqrt{n} \max(|B(\cdot)|) / \min(\text{svd}(A))^2 = \infty \text{ and so the inequality is true.}$$

If $A'Ax = b$ and $QR = A$ is the economy-size QR decomposition of A , then $A'Ax = R'Q'QRx = R'Rx = b$.

If A is full rank then $x = R^{-1} \cdot ((R')^{-1}b)$. Let $x = X(:, j)$ be the j th column of X , and $b = B(:, j)$ be the j th column of B . Then

$$\begin{aligned}
 \max(|x(:)|) &= \|x\|_\infty \\
 &\leq \|x\|_2 \\
 &= \|R^{-1} \cdot ((R')^{-1}b)\|_2 \\
 &\leq \|R^{-1}\|_2 \|(R')^{-1}\|_2 \|b\|_2 \\
 &= (1/\min(\text{svd}(A))^2) \cdot \|b\|_2 \\
 &= \|b\|_2 / \min(\text{svd}(A))^2 \\
 &\leq \sqrt{n} \|b\|_\infty / \min(\text{svd}(A))^2 \\
 &= \sqrt{n} \max(|b(:)|) / \min(\text{svd}(A))^2.
 \end{aligned}$$

Since $\max(|x(:)|) \leq \sqrt{n} \max(|b(:)|) / \min(\text{svd}(A))^2$ for all rows and columns of B and X , then

$$\max(|X(:)|) \leq \frac{\sqrt{n} \max(|B(:)|)}{\min(\text{svd}(A))^2}.$$

Lower Bound for min(svd(A))

You can estimate a lower bound s of $\min(\text{svd}(A))$ for real-valued A using the following formula,

$$s = \sigma_N \sqrt{2\gamma^{-1} \left(\frac{p_s \Gamma(m-n+1) \Gamma(n/2)}{2^{m-n} \Gamma(\frac{m+1}{2}) \Gamma(\frac{m-n+1}{2})}, \frac{m-n+1}{2} \right)}$$

where σ_N is the standard deviation of random noise added to the elements of A , $1 - p_s$ is the probability that $s \leq \min(\text{svd}(A))$, Γ is the gamma function, and γ^{-1} is the inverse incomplete gamma function `gammaincinv`.

The proof is found in [1]. It is derived by integrating the formula in Lemma 3.3 from [3] and rearranging terms.

Since $s \leq \min(\text{svd}(A))$ with probability $1 - p_s$, then you can bound the magnitude of the elements of X without computing $\text{svd}(A)$,

$$\max(|X(:)|) \leq \frac{\sqrt{n} \max(|B(:)|)}{\min(\text{svd}(A))^2} \leq \frac{\sqrt{n} \max(|B(:)|)}{s^2} \text{ with probability } 1 - p_s.$$

You can compute s using the `fixed.realSingularValueLowerBound` function which uses a default probability of 5 standard deviations below the mean,

$p_s = (1 + \text{erf}(-5/\sqrt{2}))/2 \approx 2.8665 \cdot 10^{-7}$, so the probability that the estimated bound for the smallest singular value s is less than the actual smallest singular value of A is $1 - p_s \approx 0.9999997$.

Example

This example runs a simulation with many random matrices and compares the analytical bounds with the actual singular values of A and the actual largest elements of $R = Q'A$, and $X = (A'A) \setminus B$.

Define System Parameters

Define the matrix attributes and system parameters for this example.

`m` is the number of rows in matrix `A`. In a problem such as beamforming or direction finding, `m` corresponds to the number of samples that are integrated over.

```
m = 300;
```

`n` is the number of columns in matrix `A` and rows in matrices `B` and `X`. In a least-squares problem, `m` is greater than `n`, and usually `m` is much larger than `n`. In a problem such as beamforming or direction finding, `n` corresponds to the number of sensors.

```
n = 10;
```

`p` is the number of columns in matrices `B` and `X`. It corresponds to simultaneously solving a system with `p` right-hand sides.

```
p = 1;
```

In this example, set the rank of matrix `A` to be less than the number of columns. In a problem such as beamforming or direction finding, `rank(A)` corresponds to the number of signals impinging on the sensor array.

```
rankA = 3;
```

`precisionBits` defines the number of bits of precision required for the matrix solve. Set this value according to system requirements.

```
precisionBits = 24;
```

In this example, real-valued matrices `A` and `B` are constructed such that the magnitude of their elements is less than or equal to one. Your own system requirements will define what those values are. If you don't know what they are, and `A` and `B` are fixed-point inputs to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of `A` and `B`.

`max_abs_A` is an upper bound on the maximum magnitude element of `A`.

```
max_abs_A = 1;
```

`max_abs_B` is an upper bound on the maximum magnitude element of `B`.

```
max_abs_B = 1;
```

Thermal noise standard deviation is the square root of thermal noise power, which is a system parameter. A well-designed system has the quantization level lower than the thermal noise. Here, set `thermalNoiseStandardDeviation` to the equivalent of -50 dB noise power.

```
thermalNoiseStandardDeviation = sqrt(10^(-50/10))
```

```
thermalNoiseStandardDeviation = 0.0032
```

The standard deviation of the noise from quantizing a real signal is $2^{-\text{precisionBits}}/\sqrt{12}$ [4,5]. Use `fixed.realQuantizationNoiseStandardDeviation` to compute this. See that it is less than `thermalNoiseStandardDeviation`.

```
quantizationNoiseStandardDeviation = fixed.realQuantizationNoiseStandardDeviation(precisionBits)
```

```
quantizationNoiseStandardDeviation = 1.7206e-08
```

Compute Fixed-Point Types

In this example, assume that the designed system matrix A does not have full rank (there are fewer signals of interest than number of columns of matrix A), and the measured system matrix A has additive thermal noise that is larger than the quantization noise. The additive noise makes the measured matrix A have full rank.

Set $\sigma_{\text{noise}} = \sigma_{\text{thermal noise}}$.

```
noiseStandardDeviation = thermalNoiseStandardDeviation;
```

Use `fixed.realQlessQRMatrixSolveFixedpointTypes` to compute fixed-point types.

```
T = fixed.realQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,noiseStandardDeviation)
```

```
T = struct with fields:
    A: [0x0 embedded.fi]
    B: [0x0 embedded.fi]
    X: [0x0 embedded.fi]
```

`T.A` is the type computed for transforming A to R in-place so that it does not overflow.

`T.A`

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 31
    FractionLength: 24
```

`T.B` is the type computed for B so that it does not overflow.

`T.B`

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 27
    FractionLength: 24
```

`T.X` is the type computed for the solution $X = (A'A)\backslash B$ so that there is a low probability that it overflows.

`T.X`

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
```

```
WordLength: 40
FractionLength: 24
```

Upper Bound for R

The upper bound for R is computed using the formula $\max(|R(:)|) \leq \sqrt{m} \max(|A(:)|)$, where m is the number of rows of matrix A . This upper bound is used to select a fixed-point type with the required number of bits of precision to avoid an overflow in the upper bound.

```
upperBoundR = sqrt(m)*max_abs_A
```

```
upperBoundR = 17.3205
```

Lower Bound for min(svd(A)) for Real A

A lower bound for $\min(\text{svd}(A))$ is estimated by the `fixed.realSingularValueLowerBound` function using a probability that the estimate s is not greater than the actual smallest singular value. The default probability is 5 standard deviations below the mean. You can change this probability by specifying it as the last input parameter to the `fixed.realSingularValueLowerBound` function.

```
estimatedSingularValueLowerBound = fixed.realSingularValueLowerBound(m,n,noiseStandardDeviation)
```

```
estimatedSingularValueLowerBound = 0.0371
```

Simulate and Compare to the Computed Bounds

The bounds are within an order of magnitude of the simulated results. This is sufficient because the number of bits translates to a logarithmic scale relative to the range of values. Being within a factor of 10 is between 3 and 4 bits. This is a good starting point for specifying a fixed-point type. If you run the simulation for more samples, then it is more likely that the simulated results will be closer to the bound. This example uses a limited number of simulations so it doesn't take too long to run. For real-world system design, you should run additional simulations.

Define the number of samples, `numSamples`, over which to run the simulation.

```
numSamples = 1e4;
```

Run the simulation.

```
[actualMaxR,singularValues,X_values] = runSimulations(m,n,p,rankA,max_abs_A,max_abs_B,numSamples
    noiseStandardDeviation,T);
```

You can see that the upper bound on R compared to the measured simulation results of the maximum value of R over all runs is within an order of magnitude.

```
upperBoundR
```

```
upperBoundR = 17.3205
```

```
max(actualMaxR)
```

```
ans = 8.1682
```

Finally, see that the estimated lower bound of $\min(\text{svd}(A))$ compared to the measured simulation results of $\min(\text{svd}(A))$ over all runs is also within an order of magnitude.

```
estimatedSingularValueLowerBound
```

```
estimatedSingularValueLowerBound = 0.0371
```

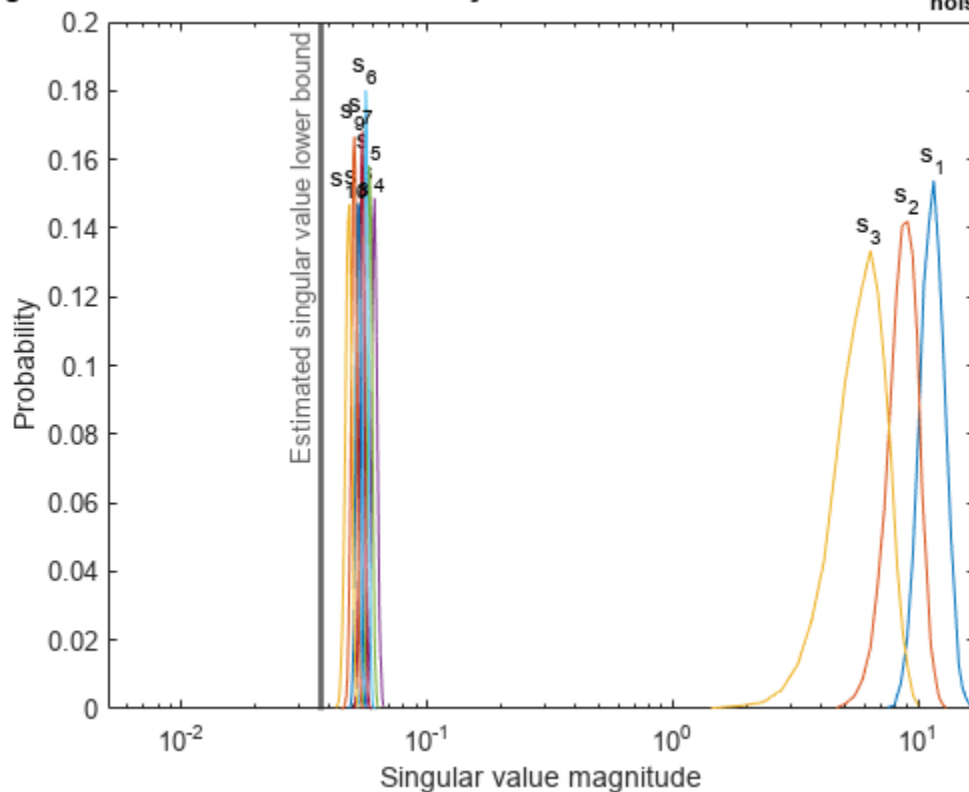


```
actualSmallestSingularValue = min(singularValues,[], 'all')
actualSmallestSingularValue = 0.0421
```

Plot the distribution of the singular values over all simulation runs. The distributions of the largest singular values correspond to the signals that determine the rank of the matrix. The distributions of the smallest singular values correspond to the noise. The derivation of the estimated bound of the smallest singular value makes use of the random nature of the noise.

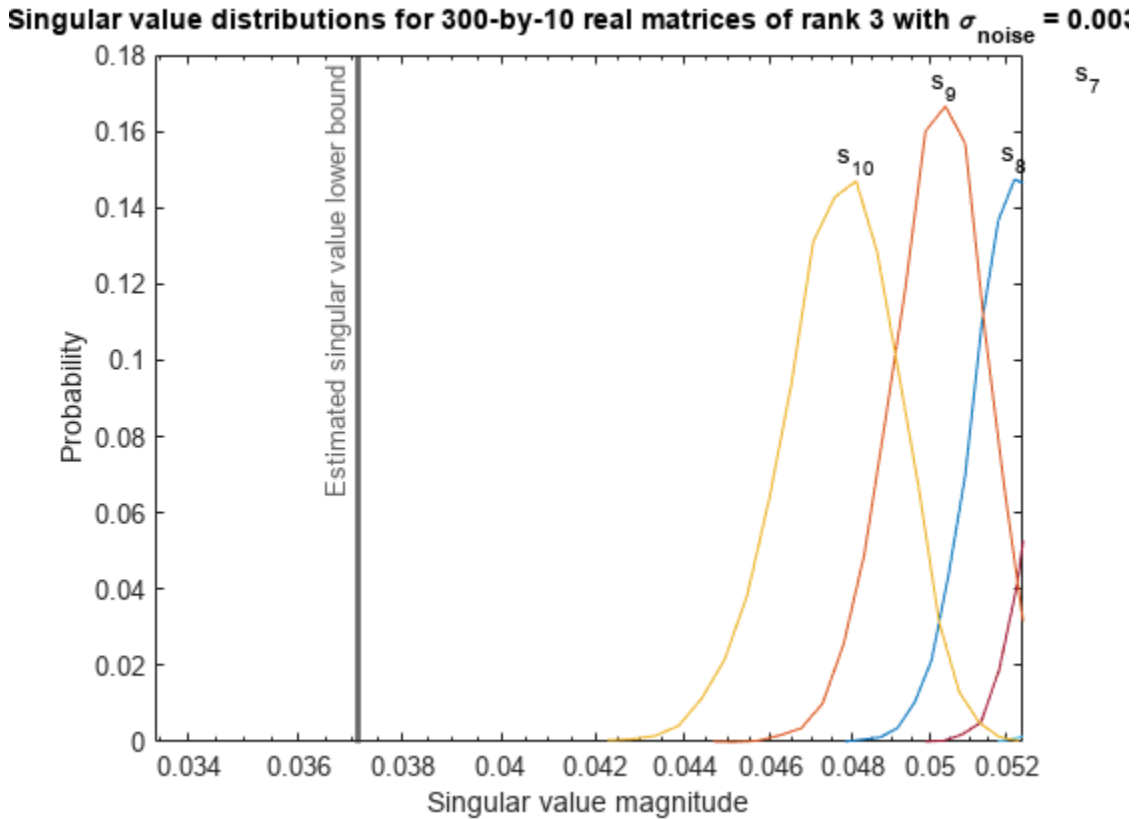
```
clf
fixed.example.plot.singularValueDistribution(m,n,rankA,...
    noiseStandardDeviation,singularValues,...
    estimatedSingularValueLowerBound,"real");
```

Singular value distributions for 300-by-10 real matrices of rank 3 with $\sigma_{\text{noise}} = 0.00$:



Zoom in to the smallest singular value to see that the estimated bound is close to it.

```
xlim([estimatedSingularValueLowerBound*0.9, max(singularValues(n,:))]);
```



Estimate the largest value of the solution, X , and compare it to the largest value of X found during the simulation runs. The estimation is within an order of magnitude of the actual value, which is sufficient for estimating a fixed-point data type, because it is between 3 and 4 bits.

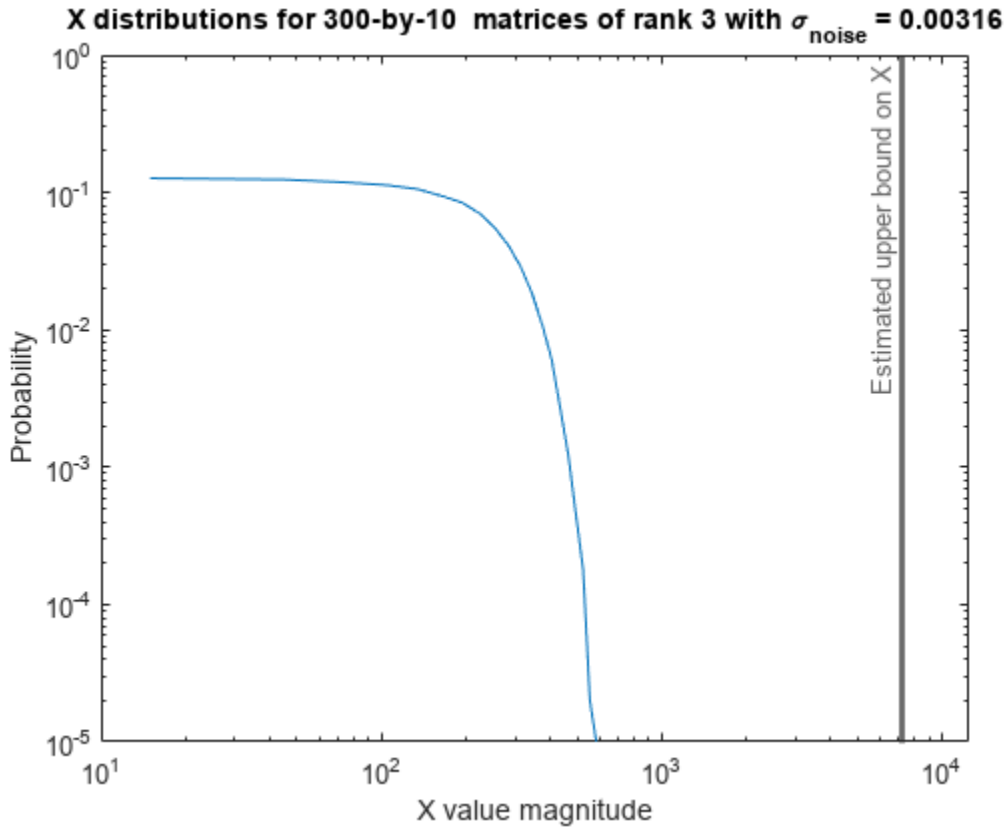
This example uses a limited number of simulation runs. With additional simulation runs, the actual largest value of X will approach the estimated largest value of X .

```
estimated_largest_X = fixed.realQlessQRMatrixSolveUpperBoundX(m,n,max_abs_B,noiseStandardDeviation)
estimated_largest_X = 7.2565e+03
```

```
actual_largest_X = max(abs(X_values),[],'all')
actual_largest_X = 582.6761
```

Plot the distribution of X values and compare it to the estimated upper bound for X .

```
clf
fixed.example.plot.xValueDistribution(m,n,rankA,noiseStandardDeviation,...
    X_values,estimated_largest_X,"real normally distributed random");
```



Supporting Functions

The `runSimulations` function creates a series of random matrices A and B of a given size and rank, quantizes them according to the computed types, computes the QR decomposition of A , and solves the equation $A'AX = B$. It returns the maximum values of $R = Q'A$, the singular values of A , and the values of X so their distributions can be plotted and compared to the bounds.

```
function [actualMaxR,singularValues,X_values] = runSimulations(m,n,p,rankA,max_abs_A,max_abs_B, .
    numSamples,noiseStandardDeviation,T)
precisionBits = T.A.FractionLength;
A_WordLength = T.A.WordLength;
B_WordLength = T.B.WordLength;
actualMaxR = zeros(1,numSamples);
singularValues = zeros(n,numSamples);
X_values = zeros(n,numSamples);
for j = 1:numSamples
    A = max_abs_A*fixed.example.realRandomLowRankMatrix(m,n,rankA);
    % Adding random noise makes A non-singular.
    A = A + fixed.example.realNormalRandomArray(0,noiseStandardDeviation,m,n);
    A = quantizenumeric(A,1,A_WordLength,precisionBits);
    B = fixed.example.realUniformRandomArray(-max_abs_B,max_abs_B,n,p);
    B = quantizenumeric(B,1,B_WordLength,precisionBits);
    [~,R] = qr(A,0);
    X = R\(R'\B);
    actualMaxR(j) = max(abs(R(:)));
    singularValues(:,j) = svd(A);
    X_values(:,j) = X;
```

```
end  
end
```

References

- 1 Thomas A. Bryan and Jenna L. Warren. “Systems and Methods for Design Parameter Selection”. Patent pending. U.S. Patent Application No. 16/947,130. 2020.
- 2 Perform QR Factorization Using CORDIC. Derivation of the bound on growth when computing QR. MathWorks. 2010.
- 3 Zizhong Chen and Jack J. Dongarra. “Condition Numbers of Gaussian Random Matrices”. In: SIAM J. Matrix Anal. Appl. 27.3 (July 2005), pp. 603–620. issn: 0895-4798. doi: 10.1137/040616413. url: <https://dx.doi.org/10.1137/040616413>.
- 4 Bernard Widrow. “A Study of Rough Amplitude Quantization by Means of Nyquist Sampling Theory”. In: IRE Transactions on Circuit Theory 3.4 (Dec. 1956), pp. 266–276.
- 5 Bernard Widrow and István Kollár. Quantization Noise - Roundoff Error in Digital Computation, Signal Processing, Control, and Communications. Cambridge, UK: Cambridge University Press, 2008.
- 6 Gene H. Golub and Charles F. Van Loan. Matrix Computations. Second edition. Baltimore: Johns Hopkins University Press, 1989.

Suppress mlint warnings in this file.

```
 %#ok< *NASGU>  
 %#ok< *ASGLU>
```

See Also

Functions

`fixed.realQlessQRMatrixSolveFixedpointTypes` | `fixed.realSingularValueLowerBound` | `fixed.qlessQRMatrixSolve`

Blocks

Real Burst Matrix Solve Using Q-less QR Decomposition | Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition | Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Forgetting Factor

Related Examples

- “Determine Fixed-Point Types for Real Q-less QR Matrix Solve $A'AX=B$ ” on page 48-179

Determine Fixed-Point Types for Real Q-less QR Matrix Solve $A'AX=B$

This example shows how to use the `fixed.realQlessQRMatrixSolveFixedpointTypes` function to analytically determine fixed-point types for the solution of the real least-squares matrix equation $A'AX = B$, where A is an m -by- n matrix with $m \geq n$, B is n -by- p , and X is n -by- p .

Fixed-point types for the solution of the matrix equation $A'AX = B$ are well-bounded if the number of rows, m , of A are much greater than the number of columns, n (i.e. $m \gg n$), and A is full rank. If A is not inherently full rank, then it can be made so by adding random noise. Random noise naturally occurs in physical systems, such as thermal noise in radar or communications systems. If $m = n$, then the dynamic range of the system can be unbounded, for example in the scalar equation $x = a/b$ and $a, b \in [-1, 1]$, then x can be arbitrarily large if b is close to 0.

Define System Parameters

Define the matrix attributes and system parameters for this example.

m is the number of rows in matrix A . In a problem such as beamforming or direction finding, m corresponds to the number of samples that are integrated over.

```
m = 300;
```

n is the number of columns in matrix A and rows in matrices B and X . In a least-squares problem, m is greater than n , and usually m is much larger than n . In a problem such as beamforming or direction finding, n corresponds to the number of sensors.

```
n = 10;
```

p is the number of columns in matrices B and X . It corresponds to simultaneously solving a system with p right-hand sides.

```
p = 1;
```

In this example, set the rank of matrix A to be less than the number of columns. In a problem such as beamforming or direction finding, $\text{rank}(A)$ corresponds to the number of signals impinging on the sensor array.

```
rankA = 3;
```

`precisionBits` defines the number of bits of precision required for the matrix solve. Set this value according to system requirements.

```
precisionBits = 24;
```

In this example, real-valued matrices A and B are constructed such that the magnitude of their elements is less than or equal to one. Your own system requirements will define what those values are. If you don't know what they are, and A and B are fixed-point inputs to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of A and B .

`max_abs_A` is an upper bound on the maximum magnitude element of A .

```
max_abs_A = 1;
```

`max_abs_B` is an upper bound on the maximum magnitude element of B .

```
max_abs_B = 1;
```

Thermal noise standard deviation is the square root of thermal noise power, which is a system parameter. A well-designed system has the quantization level lower than the thermal noise. Here, set `thermalNoiseStandardDeviation` to the equivalent of -50dB noise power.

```
thermalNoiseStandardDeviation = sqrt(10^(-50/10))
```

```
thermalNoiseStandardDeviation = 0.0032
```

The quantization noise standard deviation is a function of the required number of bits of precision. Use `fixed.realQuantizationNoiseStandardDeviation` to compute this. See that it is less than `thermalNoiseStandardDeviation`.

```
quantizationNoiseStandardDeviation = fixed.realQuantizationNoiseStandardDeviation(precisionBits)
```

```
quantizationNoiseStandardDeviation = 1.7206e-08
```

Compute Fixed-Point Types

In this example, assume that the designed system matrix A does not have full rank (there are fewer signals of interest than number of columns of matrix A), and the measured system matrix A has additive thermal noise that is larger than the quantization noise. The additive noise makes the measured matrix A have full rank.

Set $\sigma_{\text{noise}} = \sigma_{\text{thermal noise}}$.

```
noiseStandardDeviation = thermalNoiseStandardDeviation;
```

Use `fixed.realQlessQRMatrixSolveFixedpointTypes` to compute fixed-point types.

```
T = fixed.realQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,noiseStandardDeviation)
```

```
T = struct with fields:
    A: [0x0 embedded.fi]
    B: [0x0 embedded.fi]
    X: [0x0 embedded.fi]
```

`T.A` is the type computed for transforming A to $R = Q'A$ in-place so that it does not overflow.

```
T.A
```

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 31
    FractionLength: 24
```

`T.B` is the type computed for B so that it does not overflow.

```
T.B
```

```
ans =
```

```

[]

        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 27
        FractionLength: 24

```

T.X is the type computed for the solution $X = (A'A)\backslash B$ so that there is a low probability that it overflows.

T.X

ans =

```

[]

        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 40
        FractionLength: 24

```

Use the Specified Types to Solve the Matrix Equation $A'AX=B$

Create random matrices A and B such that $\text{rank}A=\text{rank}(A)$. Add random measurement noise to A which will make it become full rank.

```

rng('default');
[A,B] = fixed.example.realRandomQlessQRMatrices(m,n,p,rankA);
A = A + fixed.example.realNormalRandomArray(0,noiseStandardDeviation,m,n);

```

Cast the inputs to the types determined by `fixed.realQlessQRMatrixSolveFixedpointTypes`. Quantizing to fixed-point is equivalent to adding random noise [4,5].

```

A = cast(A,'like',T.A);
B = cast(B,'like',T.B);

```

Accelerate the `fixed.qlessQRMatrixSolve` function by using `fiaccel` to generate a MATLAB executable (MEX) function.

```

fiaccel fixed.qlessQRMatrixSolve -args {A,B,T.X} -o qlessQRMatrixSolve_mex

```

Specify output type T.X and compute fixed-point $X = (A'A)\backslash B$ using the QR method.

```

X = qlessQRMatrixSolve_mex(A,B,T.X);

```

Compute the relative error to verify the accuracy of the output.

```

relative_error = norm(double(A'*A*X - B))/norm(double(B))

```

```

relative_error = 0.0561

```

Suppress mlint warnings in this file.

```
 %#ok< *NASGU>  
 %#ok< *ASGLU>
```

See Also

Functions

`fixed.realQlessQRMatrixSolveFixedpointTypes` | `fixed.qlessQRMatrixSolve`

Blocks

Real Burst Matrix Solve Using Q-less QR Decomposition | Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition | Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Forgetting Factor

Related Examples

- “Algorithms to Determine Fixed-Point Types for Real Q-less QR Matrix Solve $A'AX=B$ ” on page 48-169

Algorithms to Determine Fixed-Point Types for Real Least-Squares Matrix Solve $AX=B$

This example shows the algorithms that the `fixed.realQRMatrixSolveFixedpointTypes` function uses to analytically determine fixed-point types for the solution of the real least-squares matrix equation $AX = B$, where A is an m -by- n matrix with $m \geq n$, B is m -by- p , and X is n -by- p .

Overview

You can solve the fixed-point least-squares matrix equation $AX = B$ using QR decomposition. Using a sequence of orthogonal transformations, QR decomposition transforms matrix A in-place to upper triangular R , and transforms matrix B in-place to $C = QB$, where $QR = A$ is the economy-size QR decomposition. This reduces the equation to an upper-triangular system of equations $RX = C$. To solve for X , compute $X = R \setminus C$ through back-substitution of R into C .

You can determine appropriate fixed-point types for the least-squares matrix equation $AX = B$ by selecting the fraction length based on the number of bits of precision defined by your requirements. The `fixed.realQRMatrixSolveFixedpointTypes` function analytically computes the following upper bounds on R , $C = QB$, and X to determine the number of integer bits required to avoid overflow [1,2,3].

The upper bound for the magnitude of the elements of R is

$$\max(|R(:)|) \leq \sqrt{m} \max(|A(:)|).$$

The upper bound for the magnitude of the elements of $C = QB$ is

$$\max(|C(:)|) \leq \sqrt{m} \max(|B(:)|).$$

The upper bound for the magnitude of the elements of $X = A \setminus B$ is

$$\max(|X(:)|) \leq \frac{\sqrt{m} \max(|B(:)|)}{\min(\text{svd}(A))}.$$

Since computing `svd(A)` is more computationally expensive than solving the system of equations, the `fixed.realQRMatrixSolveFixedpointTypes` function estimates a lower bound of `min(svd(A))`.

Fixed-point types for the solution of the matrix equation $AX = B$ are generally well-bounded if the number of rows, m , of A are much greater than the number of columns, n (i.e. $m \gg n$), and A is full rank. If A is not inherently full rank, then it can be made so by adding random noise. Random noise naturally occurs in physical systems, such as thermal noise in radar or communications systems. If $m = n$, then the dynamic range of the system can be unbounded, for example in the scalar equation $x = a/b$ and $a, b \in [-1, 1]$, then x can be arbitrarily large if b is close to 0.

Proofs of the Bounds

Properties and Definitions of Vector and Matrix Norms

The proofs of the bounds use the following properties and definitions of matrix and vector norms, where Q is an orthogonal matrix, and v is a vector of length m [6].

$$\begin{aligned}
\|Av\|_2 &\leq \|A\|_2 \|v\|_2 \\
\|Q\|_2 &= 1 \\
\|v\|_\infty &= \max(|v(:)|) \\
\|v\|_\infty &\leq \|v\|_2 \leq \sqrt{m} \|v\|_\infty
\end{aligned}$$

If A is an m -by- n matrix and $QR = A$ is the economy-size QR decomposition of A , where Q is orthogonal and m -by- n and R is upper-triangular and n -by- n , then the singular values of R are equal to the singular values of A . If A is nonsingular, then

$$\|R^{-1}\|_2 = \|(R')^{-1}\|_2 = \frac{1}{\min(\text{svd}(R))} = \frac{1}{\min(\text{svd}(A))}$$

Upper Bound for $R = Q'A$

The upper bound for the magnitude of the elements of R is

$$\max(|R(:)|) \leq \sqrt{m} \max(|A(:)|).$$

Proof of Upper Bound for $R = Q'A$

The j th column of R is equal to $R(:, j) = Q'A(:, j)$, so

$$\begin{aligned}
\max(|R(:, j)|) &= \|R(:, j)\|_\infty \\
&\leq \|R(:, j)\|_2 \\
&= \|Q'A(:, j)\|_2 \\
&\leq \|Q'\|_2 \|A(:, j)\|_2 \\
&= \|A(:, j)\|_2 \\
&\leq \sqrt{m} \|A(:, j)\|_\infty \\
&= \sqrt{m} \max(|A(:, j)|) \\
&\leq \sqrt{m} \max(|A(:)|).
\end{aligned}$$

Since $\max(|R(:, j)|) \leq \sqrt{m} \max(|A(:)|)$ for all $1 \leq j$, then

$$\max(|R(:)|) \leq \sqrt{m} \max(|A(:)|).$$

Upper Bound for $C = Q'B$

The upper bound for the magnitude of the elements of $C = Q'B$ is

$$\max(|C(:)|) \leq \sqrt{m} \max(|B(:)|).$$

Proof of Upper Bound for $C = Q'B$

The proof of the upper bound for $C = Q'B$ is the same as the proof of the upper bound for $R = Q'A$ by substituting C for R and B for A .

Upper Bound for $X = A \setminus B$

The upper bound for the magnitude of the elements of $X = A \setminus B$ is

$$\max(|X(\cdot)|) \leq \frac{\sqrt{m}\max(|B(\cdot)|)}{\min(\text{svd}(A))}.$$

Proof of Upper Bound for $X = A \setminus B$

If A is not full rank, then $\min(\text{svd}(A)) = 0$, and if B is not equal to zero, then $\sqrt{m}\max(|B(\cdot)|)/\min(\text{svd}(A)) = \infty$ and so the inequality is true.

If A is full rank, then $x = R^{-1}(Q'b)$. Let $x = X(\cdot, j)$ be the j th column of X , and $b = B(\cdot, j)$ be the j th column of B . Then

$$\begin{aligned} \max(|x(\cdot)|) &= \|x\|_{\infty} \\ &\leq \|x\|_2 \\ &= \|R^{-1} \cdot (Q'b)\|_2 \\ &\leq \|R^{-1}\|_2 \|Q'\|_2 \|b\|_2 \\ &= (1/\min(\text{svd}(A))) \cdot 1 \cdot \|b\|_2 \\ &= \|b\|_2 / \min(\text{svd}(A)) \\ &\leq \sqrt{m} \|b\|_{\infty} / \min(\text{svd}(A)) \\ &= \sqrt{m}\max(|b(\cdot)|) / \min(\text{svd}(A)). \end{aligned}$$

Since $\max(|x(\cdot)|) \leq \sqrt{m}\max(|b(\cdot)|) / \min(\text{svd}(A))$ for all rows and columns of B and X , then

$$\max(|X(\cdot)|) \leq \frac{\sqrt{m}\max(|B(\cdot)|)}{\min(\text{svd}(A))}.$$

Lower Bound for $\min(\text{svd}(A))$

You can estimate a lower bound s of $\min(\text{svd}(A))$ for real-valued A using the following formula,

$$s = \sigma_N \sqrt{2\gamma^{-1} \left(\frac{p_s \Gamma(m-n+1) \Gamma(n/2)}{2^{m-n} \Gamma(\frac{m+1}{2}) \Gamma(\frac{m-n+1}{2})}, \frac{m-n+1}{2} \right)}$$

where σ_N is the standard deviation of random noise added to the elements of A , $1 - p_s$ is the probability that $s \leq \min(\text{svd}(A))$, Γ is the gamma function, and γ^{-1} is the inverse incomplete gamma function `gammaincinv`.

The proof is found in [1]. It is derived by integrating the formula in Lemma 3.3 from [3] and rearranging terms.

Since $s \leq \min(\text{svd}(A))$ with probability $1 - p_s$, then you can bound the magnitude of the elements of X without computing $\text{svd}(A)$,

$$\max(|X(\cdot)|) \leq \frac{\sqrt{m}\max(|B(\cdot)|)}{\min(\text{svd}(A))} \leq \frac{\sqrt{m}\max(|B(\cdot)|)}{s} \text{ with probability } 1 - p_s.$$

You can compute s using the `fixed.realSingularValueLowerBound` function which uses a default probability of 5 standard deviations below the mean $p_s = (1 + \text{erf}(-5/\sqrt{2}))/2 \approx 2.8665 \cdot 10^{-7}$,

so the probability that the estimated bound for the smallest singular value s is less than the actual smallest singular value of A is $1 - p_s \approx 0.9999997$.

Example

This example runs a simulation with many random matrices and compares the analytical bounds with the actual singular values of A and the actual largest elements of $R = Q'A$, $C = Q'B$, and $X = A \setminus B$.

Define System Parameters

Define the matrix attributes and system parameters for this example.

m is the number of rows in matrices A and B . In a problem such as beamforming or direction finding, m corresponds to the number of samples that are integrated over.

```
m = 300;
```

n is the number of columns in matrix A and rows in matrix X . In a least-squares problem, m is greater than n , and usually m is much larger than n . In a problem such as beamforming or direction finding, n corresponds to the number of sensors.

```
n = 10;
```

p is the number of columns in matrices B and X . It corresponds to simultaneously solving a system with p right-hand sides.

```
p = 1;
```

In this example, set the rank of matrix A to be less than the number of columns. In a problem such as beamforming or direction finding, $\text{rank}(A)$ corresponds to the number of signals impinging on the sensor array.

```
rankA = 3;
```

`precisionBits` defines the number of bits of precision required for the matrix solve. Set this value according to system requirements.

```
precisionBits = 24;
```

In this example, real-valued matrices A and B are constructed such that the magnitude of their elements is less than or equal to one. Your own system requirements will define what those values are. If you don't know what they are, and A and B are fixed-point inputs to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of A and B .

`max_abs_A` is an upper bound on the maximum magnitude element of A .

```
max_abs_A = 1;
```

`max_abs_B` is an upper bound on the maximum magnitude element of B .

```
max_abs_B = 1;
```

Thermal noise standard deviation is the square root of thermal noise power, which is a system parameter. A well-designed system has the quantization level lower than the thermal noise. Here, set `thermalNoiseStandardDeviation` to the equivalent of -50dB noise power.

```
thermalNoiseStandardDeviation = sqrt(10^(-50/10))
```

```
thermalNoiseStandardDeviation = 0.0032
```

The standard deviation of the noise from quantizing the elements of a real signal is $2^{-\text{precisionBits}}/\sqrt{12}$ [4,5]. Use the `fixed.realQuantizationNoiseStandardDeviation` function to compute this. See that it is less than `thermalNoiseStandardDeviation`.

```
quantizationNoiseStandardDeviation = fixed.realQuantizationNoiseStandardDeviation(precisionBits)
quantizationNoiseStandardDeviation = 1.7206e-08
```

Compute Fixed-Point Types

In this example, assume that the designed system matrix A does not have full rank (there are fewer signals of interest than number of columns of matrix A), and the measured system matrix A has additive thermal noise that is larger than the quantization noise. The additive noise makes the measured matrix A have full rank.

Set $\sigma_{\text{noise}} = \sigma_{\text{thermal noise}}$.

```
noiseStandardDeviation = thermalNoiseStandardDeviation;
```

Use `fixed.realQRMatrixSolveFixedpointTypes` to compute fixed-point types.

```
T = fixed.realQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,noiseStandardDeviation)
```

```
T = struct with fields:
    A: [0x0 embedded.fi]
    B: [0x0 embedded.fi]
    X: [0x0 embedded.fi]
```

T.A is the type computed for transforming A to R in-place so that it does not overflow.

T.A

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 31
    FractionLength: 24
```

T.B is the type computed for transforming B to $Q'B$ in-place so that it does not overflow.

T.B

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 31
    FractionLength: 24
```

T.X is the type computed for the solution $X = A \setminus B$ so that there is a low probability that it overflows.

T.X

ans =

[]

```

        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 36
        FractionLength: 24

```

Upper Bounds for R and C=Q'B

The upper bounds for R and $C = Q'B$ are computed using the following formulas, where m is the number of rows of matrices A and B .

$$\max(|R(:)|) \leq \sqrt{m} \max(|A(:)|)$$

$$\max(|C(:)|) \leq \sqrt{m} \max(|B(:)|)$$

These upper bounds are used to select a fixed-point type with the required number of bits of precision to avoid overflows.

```
upperBoundR = sqrt(m)*max_abs_A
```

```
upperBoundR = 17.3205
```

```
upperBoundQB = sqrt(m)*max_abs_B
```

```
upperBoundQB = 17.3205
```

Lower Bound for min(svd(A)) for Real A

A lower bound for $\min(\text{svd}(A))$ is estimated by the `fixed.realSingularValueLowerBound` function using a probability that the estimate s is not greater than the actual smallest singular value. The default probability is 5 standard deviations below the mean. You can change this probability by specifying it as the last input parameter to the `fixed.realSingularValueLowerBound` function.

```
estimatedSingularValueLowerBound = fixed.realSingularValueLowerBound(m,n,noiseStandardDeviation)
```

```
estimatedSingularValueLowerBound = 0.0371
```

Simulate and Compare to the Computed Bounds

The bounds are within an order of magnitude of the simulated results. This is sufficient because the number of bits translates to a logarithmic scale relative to the range of values. Being within a factor of 10 is between 3 and 4 bits. This is a good starting point for specifying a fixed-point type. If you run the simulation for more samples, then it is more likely that the simulated results will be closer to the bound. This example uses a limited number of simulations so it doesn't take too long to run. For real-world system design, you should run additional simulations.

Define the number of samples, `numSamples`, over which to run the simulation.

```
numSamples = 1e4;
```

Run the simulation.

```
[actualMaxR,actualMaxQB,singularValues,X_values] = runSimulations(m,n,p,rankA,max_abs_A,max_abs_B,
    numSamples,noiseStandardDeviation,T);
```

You can see that the upper bound on R compared to the measured simulation results of the maximum value of R over all runs is within an order of magnitude.

```
upperBoundR
```

```
upperBoundR = 17.3205
```

```
max(actualMaxR)
```

```
ans = 8.3029
```

You can see that the upper bound on $C = QB$ compared to the measured simulation results of the maximum value of $C = QB$ over all runs is also within an order of magnitude.

```
upperBoundQB
```

```
upperBoundQB = 17.3205
```

```
max(actualMaxQB)
```

```
ans = 2.5707
```

Finally, see that the estimated lower bound of $\min(\text{svd}(A))$ compared to the measured simulation results of $\min(\text{svd}(A))$ over all runs is also within an order of magnitude.

```
estimatedSingularValueLowerBound
```

```
estimatedSingularValueLowerBound = 0.0371
```

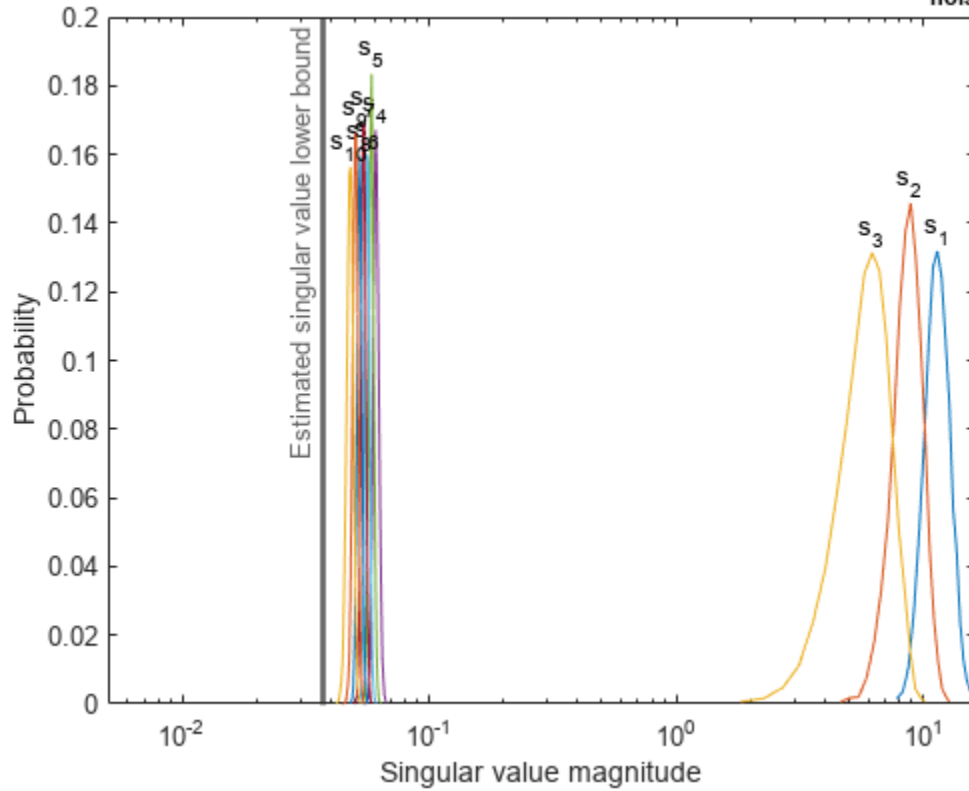
```
actualSmallestSingularValue = min(singularValues,[],'all')
```

```
actualSmallestSingularValue = 0.0420
```

Plot the distribution of the singular values over all simulation runs. The distributions of the largest singular values correspond to the signals that determine the rank of the matrix. The distributions of the smallest singular values correspond to the noise. The derivation of the estimated bound of the smallest singular value makes use of the random nature of the noise.

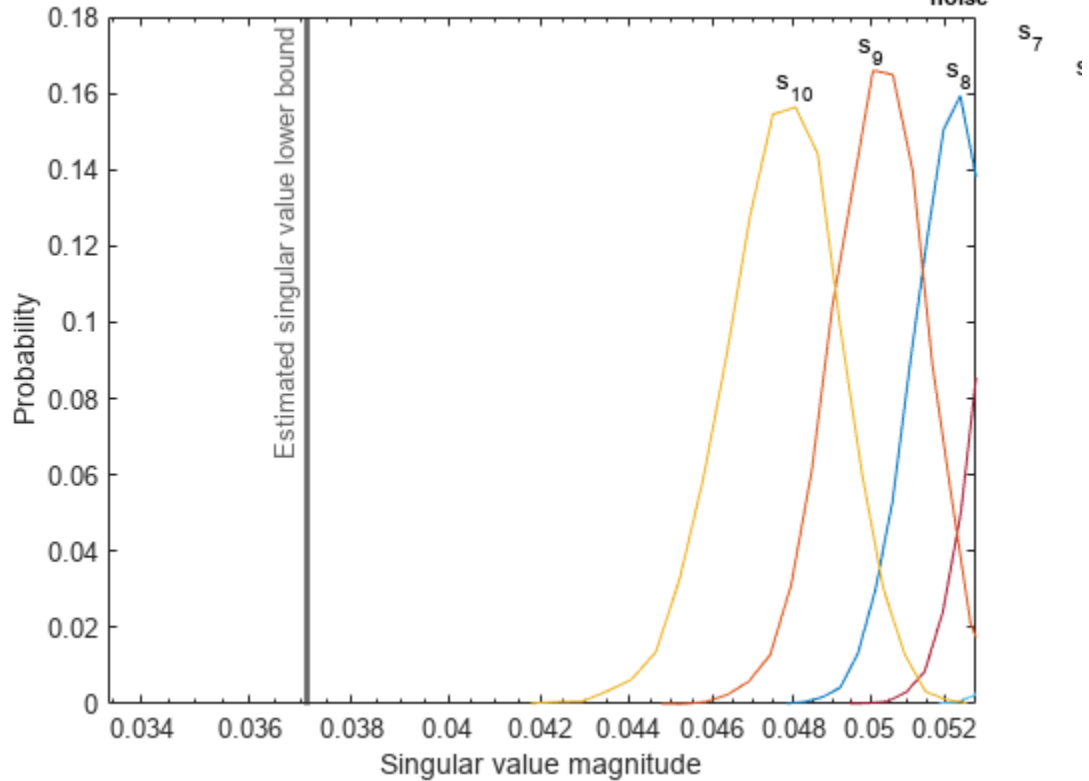
```
clf
```

```
fixed.example.plot.singularValueDistribution(m,n,rankA,noiseStandardDeviation,...
    singularValues,estimatedSingularValueLowerBound,"real");
```

Singular value distributions for 300-by-10 real matrices of rank 3 with $\sigma_{\text{noise}} = 0.001$:

Zoom in to smallest singular value to see that the estimated bound is close to it.

```
xlim([estimatedSingularValueLowerBound*0.9, max(singularValues(n,:))]);
```


Singular value distributions for 300-by-10 real matrices of rank 3 with $\sigma_{\text{noise}} = 0.001$ 

Estimate the largest value of the solution, X , and compare it to the largest value of X found during the simulation runs. The estimation is within an order of magnitude of the actual value, which is sufficient for estimating a fixed-point data type, because it is between 3 and 4 bits.

This example uses a limited number of simulation runs. With additional simulation runs, the actual largest value of X will approach the estimated largest value of X .

```
estimated_largest_X = fixed.realMatrixSolveUpperBoundX(m,n,max_abs_B,noiseStandardDeviation)
```

```
estimated_largest_X = 466.5772
```

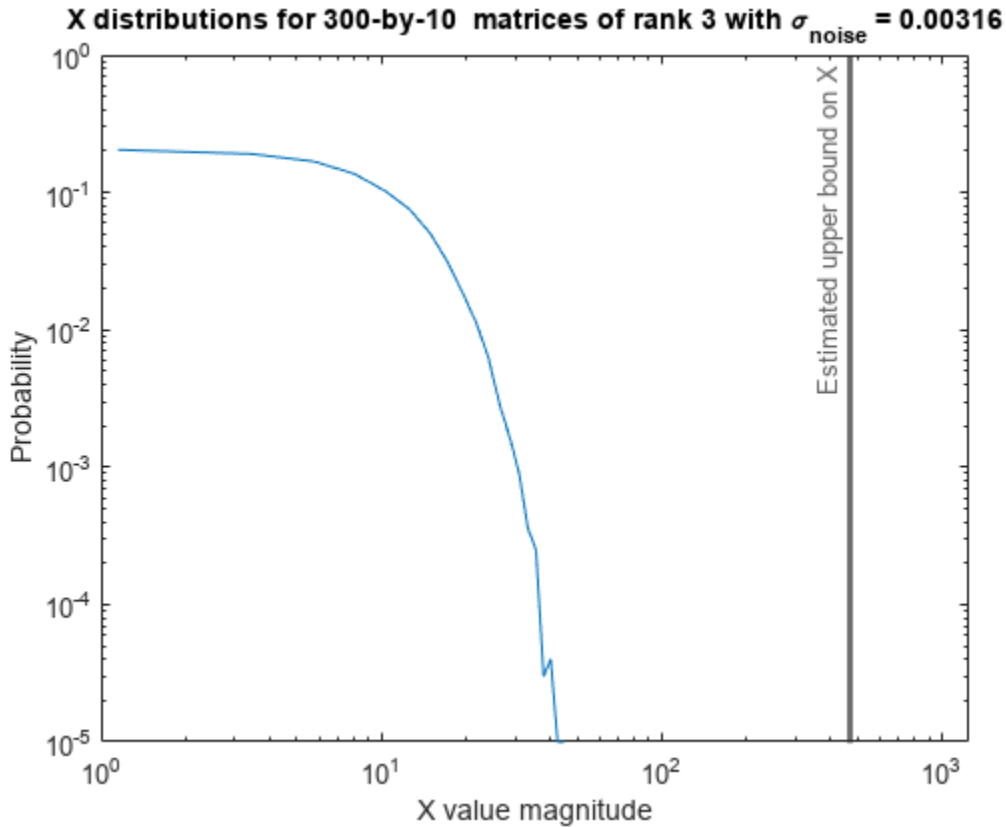
```
actual_largest_X = max(abs(X_values), [], 'all')
```

```
actual_largest_X = 44.8056
```

Plot the distribution of X values and compare it to the estimated upper bound for X .

```
clf
```

```
fixed.example.plot.xValueDistribution(m,n,rankA,noiseStandardDeviation,...  
X_values,estimated_largest_X,"real normally distributed random");
```



Supporting Functions

The `runSimulations` function creates a series of random matrices A and B of a given size and rank, quantizes them according to the computed types, computes the QR decomposition of A , and solves the equation $AX = B$. It returns the maximum values of $R = Q'A$ and $C = Q'B$, the singular values of A , and the values of X so their distributions can be plotted and compared to the bounds.

```
function [actualMaxR,actualMaxQB,singularValues,X_values] = runSimulations(m,n,p,rankA,max_abs_A,
    numSamples,noiseStandardDeviation,T)
precisionBits = T.A.FractionLength;
A_WordLength = T.A.WordLength;
B_WordLength = T.B.WordLength;
actualMaxR = zeros(1,numSamples);
actualMaxQB = zeros(1,numSamples);
singularValues = zeros(n,numSamples);
X_values = zeros(n,numSamples);
for j = 1:numSamples
    A = max_abs_A*fixed.example.realRandomLowRankMatrix(m,n,rankA);
    % Adding normally distributed random noise makes A non-singular.
    A = A + fixed.example.realNormalRandomArray(0,noiseStandardDeviation,m,n);
    A = quantiznumeric(A,1,A_WordLength,precisionBits);
    B = fixed.example.realUniformRandomArray(-max_abs_B,max_abs_B,m,p);
    B = quantiznumeric(B,1,B_WordLength,precisionBits);
    [Q,R] = qr(A,0);
    C = Q'*B;
    X = R\C;
    actualMaxR(j) = max(abs(R(:)));
end
```

```

    actualMaxQB(j) = max(abs(C(:)));
    singularValues(:,j) = svd(A);
    X_values(:,j) = X;
end
end

```

References

- 1 Thomas A. Bryan and Jenna L. Warren. “Systems and Methods for Design Parameter Selection”. Patent pending. U.S. Patent Application No. 16/947,130. 2020.
- 2 Perform QR Factorization Using CORDIC. Derivation of the bound on growth when computing QR. MathWorks. 2010.
- 3 Zizhong Chen and Jack J. Dongarra. “Condition Numbers of Gaussian Random Matrices”. In: SIAM J. Matrix Anal. Appl. 27.3 (July 2005), pp. 603–620. issn: 0895-4798. doi: 10.1137/040616413. url: <https://dx.doi.org/10.1137/040616413>.
- 4 Bernard Widrow. “A Study of Rough Amplitude Quantization by Means of Nyquist Sampling Theory”. In: IRE Transactions on Circuit Theory 3.4 (Dec. 1956), pp. 266–276.
- 5 Bernard Widrow and István Kollár. Quantization Noise - Roundoff Error in Digital Computation, Signal Processing, Control, and Communications. Cambridge, UK: Cambridge University Press, 2008.
- 6 Gene H. Golub and Charles F. Van Loan. Matrix Computations. Second edition. Baltimore: Johns Hopkins University Press, 1989.

Suppress mlint warnings in this file.

```

%#ok< *NASGU>
%#ok< *ASGLU>

```

See Also

Functions

`fixed.realQRMatrixSolveFixedpointTypes` | `fixed.realSingularValueLowerBound` | `fixed.qrMatrixSolve`

Blocks

Real Burst Matrix Solve Using QR Decomposition | Real Partial-Systolic Matrix Solve Using QR Decomposition

Related Examples

- “Determine Fixed-Point Types for Real Least-Squares Matrix Solve $AX=B$ ” on page 48-194

Determine Fixed-Point Types for Real Least-Squares Matrix Solve $AX=B$

This example shows how to use the `fixed.realQRMatrixSolveFixedpointTypes` function to analytically determine fixed-point types for the solution of the real least-squares matrix equation $AX = B$, where A is an m -by- n matrix with $m \geq n$, B is m -by- p , and X is n -by- p .

Fixed-point types for the solution of the matrix equation $AX = B$ are well-bounded if the number of rows, m , of A are much greater than the number of columns, n (i.e. $m \gg n$), and A is full rank. If A is not inherently full rank, then it can be made so by adding random noise. Random noise naturally occurs in physical systems, such as thermal noise in radar or communications systems. If $m = n$, then the dynamic range of the system can be unbounded, for example in the scalar equation $x = a/b$ and $a, b \in [-1, 1]$, then x can be arbitrarily large if b is close to 0.

Define System Parameters

Define the matrix attributes and system parameters for this example.

m is the number of rows in matrices A and B . In a problem such as beamforming or direction finding, m corresponds to the number of samples that are integrated over.

```
m = 300;
```

n is the number of columns in matrix A and rows in matrix X . In a least-squares problem, m is greater than n , and usually m is much larger than n . In a problem such as beamforming or direction finding, n corresponds to the number of sensors.

```
n = 10;
```

p is the number of columns in matrices B and X . It corresponds to simultaneously solving a system with p right-hand sides.

```
p = 1;
```

In this example, set the rank of matrix A to be less than the number of columns. In a problem such as beamforming or direction finding, $\text{rank}(A)$ corresponds to the number of signals impinging on the sensor array.

```
rankA = 3;
```

`precisionBits` defines the number of bits of precision required for the matrix solve. Set this value according to system requirements.

```
precisionBits = 24;
```

In this example, real-valued matrices A and B are constructed such that the magnitude of their elements is less than or equal to one. Your own system requirements will define what those values are. If you don't know what they are, and A and B are fixed-point inputs to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of A and B .

`max_abs_A` is an upper bound on the maximum magnitude element of A .

```
max_abs_A = 1;
```

`max_abs_B` is an upper bound on the maximum magnitude element of B .

```
max_abs_B = 1;
```

Thermal noise standard deviation is the square root of thermal noise power, which is a system parameter. A well-designed system has the quantization level lower than the thermal noise. Here, set `thermalNoiseStandardDeviation` to the equivalent of -50dB noise power.

```
thermalNoiseStandardDeviation = sqrt(10^(-50/10))
```

```
thermalNoiseStandardDeviation = 0.0032
```

The quantization noise standard deviation is a function of the required number of bits of precision. Use `fixed.realQuantizationNoiseStandardDeviation` to compute this. See that it is less than `thermalNoiseStandardDeviation`.

```
quantizationNoiseStandardDeviation = fixed.realQuantizationNoiseStandardDeviation(precisionBits)
```

```
quantizationNoiseStandardDeviation = 1.7206e-08
```

Compute Fixed-Point Types

In this example, assume that the designed system matrix A does not have full rank (there are fewer signals of interest than number of columns of matrix A), and the measured system matrix A has additive thermal noise that is larger than the quantization noise. The additive noise makes the measured matrix A have full rank.

Set $\sigma_{\text{noise}} = \sigma_{\text{thermal noise}}$.

```
noiseStandardDeviation = thermalNoiseStandardDeviation;
```

Use `fixed.realQRMatrixSolveFixedpointTypes` to compute fixed-point types.

```
T = fixed.realQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,noiseStandardDeviation)
```

```
T = struct with fields:
    A: [0x0 embedded.fi]
    B: [0x0 embedded.fi]
    X: [0x0 embedded.fi]
```

$T.A$ is the type computed for transforming A to $R = Q'A$ in-place so that it does not overflow.

```
T.A
```

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 31
    FractionLength: 24
```

$T.B$ is the type computed for transforming B to $C = QB$ in-place so that it does not overflow.

```
T.B
```

```
ans =
```

```

[]

        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 31
        FractionLength: 24

```

T.X is the type computed for the solution $X = A \setminus B$ so that there is a low probability that it overflows.

T.X

ans =

```

[]

        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 36
        FractionLength: 24

```

Use the Specified Types to Solve the Matrix Equation $AX=B$

Create random matrices A and B such that B is in the range of A, and $\text{rank}A=\text{rank}(A)$. Add random measurement noise to A which will make it become full rank, but it will also affect the solution so that B is only close to the range of A.

```

rng('default');
[A,B] = fixed.example.realRandomLeastSquaresMatrices(m,n,p,rankA);
A = A + fixed.example.realNormalRandomArray(0,noiseStandardDeviation,m,n);

```

Cast the inputs to the types determined by `fixed.realQRMatrixSolveFixedpointTypes`. Quantizing to fixed-point is equivalent to adding random noise [4,5].

```

A = cast(A,'like',T.A);
B = cast(B,'like',T.B);

```

Accelerate the `fixed.qrMatrixSolve` function by using `fiaccl` to generate a MATLAB executable (MEX) function.

```

fiaccl fixed.qrMatrixSolve -args {A,B,T.X} -o qrRealMatrixSolve_mex

```

Specify output type T.X and compute fixed-point $X = A \setminus B$ using the QR method.

```

X = qrRealMatrixSolve_mex(A,B,T.X);

```

Compute the relative error to verify the accuracy of the output.

```

relative_error = norm(double(A*X - B))/norm(double(B))
relative_error = 0.0063

```

Suppress `mlint` warnings in this file.

`%#ok<*NASGU>`
`%#ok<*ASGLU>`

See Also

Functions

`fixed.realQRMatrixSolveFixedpointTypes` | `fixed.qrMatrixSolve`

Blocks

Real Burst Matrix Solve Using QR Decomposition | Real Partial-Systolic Matrix Solve Using QR Decomposition

Related Examples

- “Algorithms to Determine Fixed-Point Types for Real Least-Squares Matrix Solve $AX=B$ ” on page 48-183

Compute Forgetting Factor Required for Streaming Input Data

This example shows how to use the `fixed.forgettingFactor` and `fixed.forgettingFactorInverse` functions.

The growth in the QR decomposition can be seen by looking at the magnitude of the first element $R(1, 1)$ of the upper-triangular factor R , which is equal to the Euclidean norm of the first column of matrix A ,

$$|R(1, 1)| = \|A(:, 1)\|_2.$$

To see this, create matrix A as a column of ones of length n and compute R of the economy-size QR decomposition.

```
n = 1e4;
A = ones(n, 1);
```

$$\text{Then } |R(1, 1)| = \|A(:, 1)\|_2 = \sqrt{\sum_{i=1}^n 1^2} = \sqrt{n}.$$

```
R = fixed.qlessQR(A)
```

```
R = 100.0000
```

```
norm(A)
```

```
ans = 100
```

```
sqrt(n)
```

```
ans = 100
```

The diagonal elements of the upper-triangular factor R of the QR decomposition may be positive, negative, or zero, but `fixed.qlessQR` and `fixed.qrAB` always return the diagonal elements of R as non-negative.

In a real-time application, such as when data is streaming continuously from a radar array, you can update the QR decomposition with an exponential forgetting factor α where $0 < \alpha < 1$. Use the `fixed.forgettingFactor` function to compute a forgetting factor α that acts as if the matrix were being integrated over m rows to maintain a gain of about \sqrt{m} . The relationship between α and m is $\alpha = e^{-1/(2m)}$.

```
m = 16;
alpha = fixed.forgettingFactor(m);
R_alpha = fixed.qlessQR(A, alpha)
```

```
R_alpha = 3.9377
```

```
sqrt(m)
```

```
ans = 4
```

If you are working with a system and have been given a forgetting factor α , and want to know the effective number of rows m that you are integrating over, then you can use the

`fixed.forgettingFactorInverse` function. The relationship between m and α is $m = \frac{-1}{2\log(\alpha)}$.


```
fixed.forgettingFactorInverse(alpha)
```

```
ans = 16
```

See Also

Functions

[fixed.forgettingFactor](#) | [fixed.forgettingFactorInverse](#) | [fixed.qlessQR](#) | [fixed.qlessQRMatrixSolve](#) | [fixed.qlessQRUpdate](#)

Blocks

[Real Partial-Systolic Q-less QR Decomposition with Forgetting Factor](#) | [Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Forgetting Factor](#)

Estimate Standard Deviation of Quantization Noise of Complex-Valued Signal

Quantizing a complex signal to p bits of precision can be modeled as a linear system that adds normally distributed noise with a standard deviation of $\zeta_{\text{noise}} = \frac{2^{-p}}{\sqrt{6}}$ [1,2].

Compute the theoretical quantization noise standard deviation with p bits of precision using the `fixed.complexQuantizationNoiseStandardDeviation` function.

```
p = 14;
theoreticalQuantizationNoiseStandardDeviation = fixed.complexQuantizationNoiseStandardDeviation(p)
```

The returned value is $\zeta_{\text{noise}} = \frac{2^{-p}}{\sqrt{6}}$.

Create a complex signal with n samples.

```
rng('default');
n = 1e6;
x = complex(rand(1,n), rand(1,n));
```

Quantize the signal with p bits of precision.

```
wordLength = 16;
x_quantized = quantizenumeric(x,1,wordLength,p);
```

Compute the quantization noise by taking the difference between the quantized signal and the original signal.

```
quantizationNoise = x_quantized - x;
```

Compute the measured quantization noise standard deviation.

```
measuredQuantizationNoiseStandardDeviation = std(quantizationNoise)
measuredQuantizationNoiseStandardDeviation = 2.4902e-05
```

Compare the actual quantization noise standard deviation to the theoretical and see that they are close for large values of n .

```
theoreticalQuantizationNoiseStandardDeviation
theoreticalQuantizationNoiseStandardDeviation = 2.4917e-05
```

References

- 1 Bernard Widrow. "A Study of Rough Amplitude Quantization by Means of Nyquist Sampling Theory". In: IRE Transactions on Circuit Theory 3.4 (Dec. 1956), pp. 266-276.

- 2 Bernard Widrow and István Kollár. Quantization Noise - Roundoff Error in Digital Computation, Signal Processing, Control, and Communications. Cambridge, UK: Cambridge University Press, 2008.

See Also

`fixed.complexQuantizationNoiseStandardDeviation`

Estimate Standard Deviation of Quantization Noise of Real-Valued Signal

Quantizing a real signal to p bits of precision can be modeled as a linear system that adds normally distributed noise with a standard deviation of $\zeta_{\text{noise}} = \frac{2^{-p}}{\sqrt{12}}$ [1,2].

Compute the theoretical quantization noise standard deviation with p bits of precision using the `fixed.realQuantizationNoiseStandardDeviation` function.

```
p = 14;
theoreticalQuantizationNoiseStandardDeviation = fixed.realQuantizationNoiseStandardDeviation(p);
```

The returned value is $\zeta_{\text{noise}} = \frac{2^{-p}}{\sqrt{12}}$.

Create a real signal with n samples.

```
rng('default');
n = 1e6;
x = rand(1,n);
```

Quantize the signal with p bits of precision.

```
wordLength = 16;
x_quantized = quantizenumeric(x,1,wordLength,p);
```

Compute the quantization noise by taking the difference between the quantized signal and the original signal.

```
quantizationNoise = x_quantized - x;
```

Compute the measured quantization noise standard deviation.

```
measuredQuantizationNoiseStandardDeviation = std(quantizationNoise)
measuredQuantizationNoiseStandardDeviation = 1.7607e-05
```

Compare the actual quantization noise standard deviation to the theoretical and see that they are close for large values of n .

```
theoreticalQuantizationNoiseStandardDeviation
theoreticalQuantizationNoiseStandardDeviation = 1.7619e-05
```

References

- 1 Bernard Widrow. "A Study of Rough Amplitude Quantization by Means of Nyquist Sampling Theory". In: IRE Transactions on Circuit Theory 3.4 (Dec. 1956), pp. 266-276.

- 2 Bernard Widrow and István Kollár. Quantization Noise - Roundoff Error in Digital Computation, Signal Processing, Control, and Communications. Cambridge, UK: Cambridge University Press, 2008.

See Also

`fixed.realQuantizationNoiseStandardDeviation`

Implement Hardware-Efficient Real Partial-Systolic Q-less QR with Forgetting Factor

This example shows how to use the hardware-efficient Real Partial-Systolic Q-less QR Decomposition with Forgetting Factor block.

Q-less QR Decomposition with Forgetting Factor

The Real Partial-Systolic Q-less QR Decomposition with Forgetting Factor block implements the following recursion to compute the upper-triangular factor R of continuously streaming n-by-1 row vectors A(k,:) using forgetting factor α . It is as if matrix A is infinitely tall. The forgetting factor in the range $0 < \alpha < 1$ keeps it from integrating without bound.

$$\begin{aligned}
 R_0 &= \text{zeros}(n, n) \\
 [\sim, R_1] &= \text{qr} \left(\begin{bmatrix} R_0 \\ A(1, :) \end{bmatrix}, 0 \right) \\
 R_1 &= \alpha R_1 \\
 \\
 [\sim, R_2] &= \text{qr} \left(\begin{bmatrix} R_1 \\ A(2, :) \end{bmatrix}, 0 \right) \\
 R_2 &= \alpha R_2 \\
 \\
 &\vdots \\
 \\
 [\sim, R_k] &= \text{qr} \left(\begin{bmatrix} R_{k-1} \\ A(k, :) \end{bmatrix}, 0 \right) \\
 R_k &= \alpha R_k \\
 \\
 &\vdots
 \end{aligned}$$

AMBA AXI Handshaking Process

The Data Handler subsystem in this model takes real matrix A as input. It sends rows of A to the QR Decomposition block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set delay for the feeding in rows of A in the Data Handler to emulate the processing time of the upstream block. `validOut` signal of the Data Handler remain high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Define System Parameters

n is the length of the row vectors A(k,:) and the number of rows and columns in R.

n = 5;

m is the effective number of rows of A to integrate over.

```
m = 100;
```

Use the `fixed.forgettingFactor` function to compute the forgetting factor as a function of the number of rows that you are integrating over.

```
forgettingFactor = fixed.forgettingFactor(m)
```

```
forgettingFactor =
```

```
    0.9950
```

`precisionBits` defines the number of bits of precision required for the QR Decomposition. Set this value according to system requirements.

```
precisionBits = 24;
```

In this example, real-valued matrix `A` is constructed such that the magnitude of its elements are less than or equal to one. Your own system requirements will define what those values are. If you don't know what they are, and `A` is a fixed-point input to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of `A`.

`max_abs_A` is an upper bound on the maximum magnitude element of `A`.

```
max_abs_A = 1;
```

Select Fixed-Point Types

Use the `fixed.qlessqrFixedpointTypes` function to compute fixed-point types.

```
T = fixed.qlessqrFixedpointTypes(m,max_abs_A,precisionBits)
```

```
T =
```

```
  struct with fields:
```

```
    A: [0x0 embedded.fi]
```

`T.A` is the fixed-point type computed for transforming `A` to `R` in-place so that it does not overflow.

```
T.A
```

```
ans =
```

```
 []
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 31
    FractionLength: 24
```

Define Simulation Parameters

Create random matrix `A` to contain a specified number of inputs.

`numInputs` is the number of input rows $A(k,:)$ for this example.

```
numInputs = 500;
rng('default')
A = fixed.example.realUniformRandomArray(-1,1,numInputs,n);
```

Cast the inputs to the types determined by `fixed.qlessqrFixedpointTypes`.

```
A = cast(A,'like',T.A);
```

Cast the forgetting factor to a fixed-point type with the same word length as A and best-precision scaling.

```
forgettingFactor = fi(forgettingFactor,1,T.A.WordLength);
```

Set delay of clock cycles between feeding in rows of A

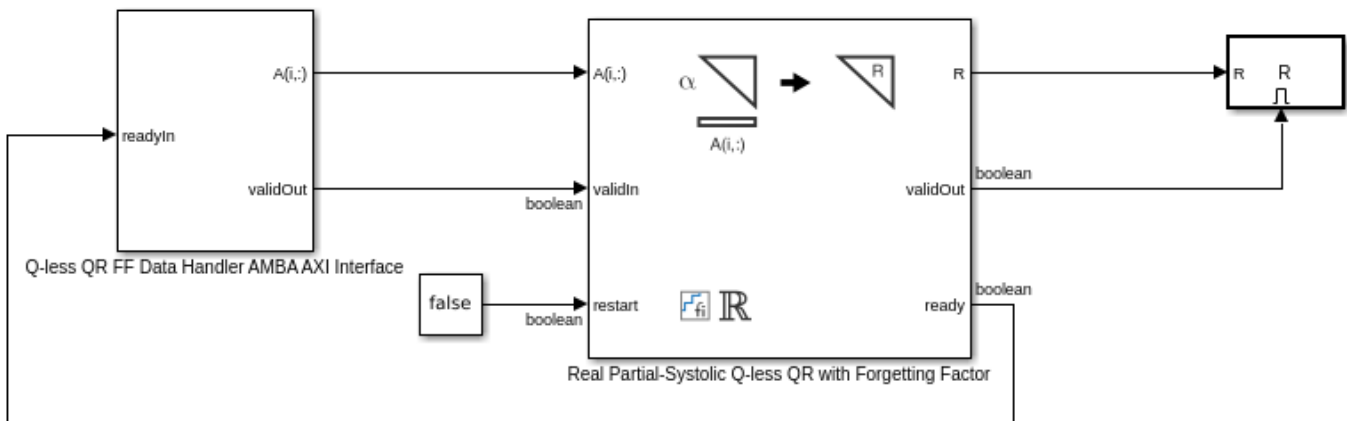
```
rowDelay = 1;
```

Select a stop time for the simulation that is long enough to process all the inputs from A .

```
stopTime = 2*numInputs*T.A.WordLength;
```

Open the Model

```
model = 'RealPartialSystolicQlessQRForgettingFactorModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace.

```
fixed.example.setModelWorkspace(model,'A',A,'n',n,...
    'forgettingFactor',forgettingFactor,...
    'regularizationParameter',0,...
    'rowDelay',rowDelay,...
    'stopTime',stopTime);
```


Simulate the Model

```
out = sim(model);
```

Verify the Accuracy of the Output

Define matrix A_k as follows

$$A_k = \begin{bmatrix} \alpha^k & & & & \\ & \alpha^{k-1} & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & \alpha \end{bmatrix} A(1:k, :).$$

Then using the formula for the computation of the k th output R_k , and the fact that $[Q, R] = \text{qr}(A, 0) \Rightarrow A'A = R'Q'QR = R'R$, you can show that

$$A'_k A_k = R'_k R_k.$$

So to verify the output, the difference between $A'_k A_k$ and $R'_k R_k$ should be small.

Choose the last output of the simulation.

```
R = double(out.R(:, :, end))
```

```
R =
```

```
    5.2717    0.1389   -0.2004    0.4152   -0.0565
         0    5.5798    0.0561   -0.1367   -0.5387
         0         0    5.4077   -0.0043   -0.1051
         0         0         0    5.6626   -0.1625
         0         0         0         0    5.4957
```

Verify that R is upper triangular.

```
isequal(R, triu(R))
```

```
ans =
```

```
logical
```

```
1
```

Verify that the diagonal is greater than or equal to zero.

```
diag(R)
```

```
ans =
```

```
    5.2717
    5.5798
```

```
5.4077
5.6626
5.4957
```

Synchronize the last output R with the input by finding the number of inputs that produced it.

```
A = double(A);
alpha = double(forgettingFactor);
relative_errors = nan(1,n);
for k = 1:numInputs
    A_k = alpha.^(k:-1:1)' .* A(1:k,:);
    relative_errors(k) = norm(A_k'*A_k - R'*R)/norm(A_k'*A_k);
end
```

k is the number of inputs A(k,:) that produced the last R.

```
k = find(relative_errors==min(relative_errors),1,'last')
```

```
k =
    500
```

Verify that

$$A_k' A_k = R_k' R_k$$

with a small relative error.

```
A_k = alpha.^(k:-1:1)' .* A(1:k,:);
relative_error = norm(A_k'*A_k - R'*R)/norm(A_k'*A_k)
```

```
relative_error =
    8.3674e-06
```

Suppress mlint warnings in this file.

```
 %#ok<*NOPTS>
```

See Also

Real Partial-Systolic Q-less QR Decomposition with Forgetting Factor

Implement Hardware-Efficient Complex Partial-Systolic Q-less QR with Forgetting Factor

This example shows how to use the hardware-efficient Complex Partial-Systolic Q-less QR Decomposition with Forgetting Factor block.

Q-less QR Decomposition with Forgetting Factor

The Complex Partial-Systolic Q-less QR Decomposition with Forgetting Factor block implements the following recursion to compute the upper-triangular factor R of continuously streaming n-by-1 row vectors A(k,:) using forgetting factor α . It is as if matrix A is infinitely tall. The forgetting factor in the range $0 < \alpha < 1$ keeps it from integrating without bound.

$$\begin{aligned}
 R_0 &= \text{zeros}(n, n) \\
 [\sim, R_1] &= \text{qr} \left(\begin{bmatrix} R_0 \\ A(1, :) \end{bmatrix}, 0 \right) \\
 R_1 &= \alpha R_1 \\
 \\
 [\sim, R_2] &= \text{qr} \left(\begin{bmatrix} R_1 \\ A(2, :) \end{bmatrix}, 0 \right) \\
 R_2 &= \alpha R_2 \\
 \\
 &\vdots \\
 \\
 [\sim, R_k] &= \text{qr} \left(\begin{bmatrix} R_{k-1} \\ A(k, :) \end{bmatrix}, 0 \right) \\
 R_k &= \alpha R_k \\
 \\
 &\vdots
 \end{aligned}$$

AMBA AXI Handshaking Process

The Data Handler subsystem in this model takes complex matrix A as input. It sends rows of A to the QR Decomposition block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set delay for the feeding in rows of A in the Data Handler to emulate the processing time of the upstream block. `validOut` signal of the Data Handler remain high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Define System Parameters

n is the length of the row vectors A(k,:) and the number of rows and columns in R.

n = 5;

m is the effective number of rows of A to integrate over.

```
m = 100;
```

Use the `fixed.forgettingFactor` function to compute the forgetting factor as a function of the number of rows that you are integrating over.

```
forgettingFactor = fixed.forgettingFactor(m)
```

```
forgettingFactor =
```

```
    0.9950
```

`precisionBits` defines the number of bits of precision required for the QR Decomposition. Set this value according to system requirements.

```
precisionBits = 24;
```

In this example, complex-valued matrix *A* is constructed such that the magnitude of the real and imaginary parts of its elements is less than or equal to one, so the maximum possible absolute value of any element is $|1 + 1i| = \sqrt{2}$. Your own system requirements will define what those values are. If you don't know what they are, and *A* is a fixed-point input to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of *A*.

`max_abs_A` is an upper bound on the maximum magnitude element of *A*.

```
max_abs_A = sqrt(2);
```

Select Fixed-Point Types

Use the `fixed.qlessqrFixedpointTypes` function to compute fixed-point types.

```
T = fixed.qlessqrFixedpointTypes(m,max_abs_A,precisionBits)
```

```
T =
```

```
  struct with fields:
```

```
    A: [0x0 embedded.fi]
```

`T.A` is the fixed-point type computed for transforming *A* to *R* in-place so that it does not overflow.

```
T.A
```

```
ans =
```

```
 []
```

```
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 31  
    FractionLength: 24
```

Define Simulation Parameters

Create random matrix *A* to contain a specified number of inputs.

numInputs is the number of input rows A(k,:) for this example.

```
numInputs = 500;
rng('default')
A = fixed.example.complexUniformRandomArray(-1,1,numInputs,n);
```

Cast the inputs to the types determined by fixed.qlessqrFixedpointTypes.

```
A = cast(A,'like',T.A);
```

Cast the forgetting factor to a fixed-point type with the same word length as A and best-precision scaling.

```
forgettingFactor = fi(forgettingFactor,1,T.A.WordLength);
```

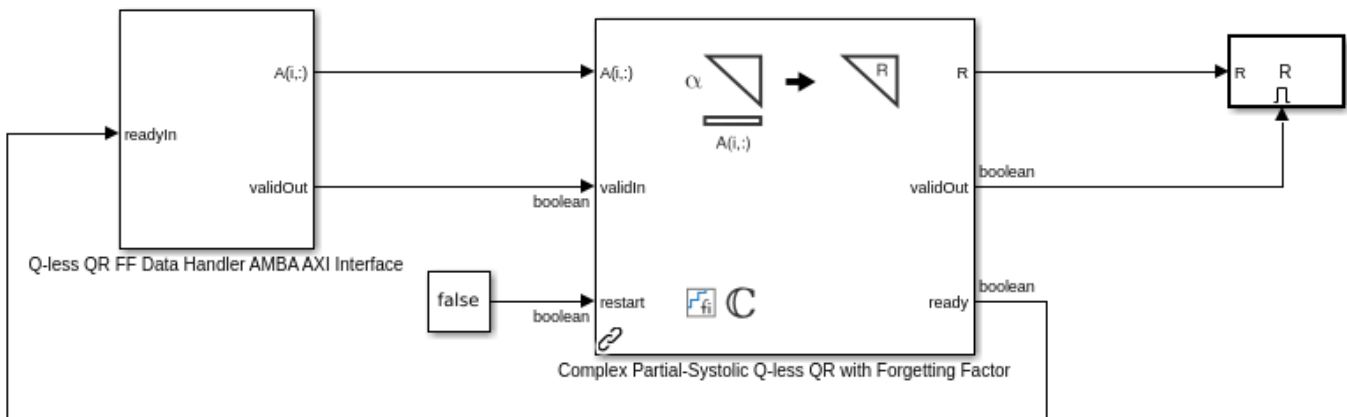
```
rowDelay = 1; % Delay of clock cycles between feeding in rows of A
```

Select a stop time for the simulation that is long enough to process all the inputs from A.

```
stopTime = 4*numInputs*T.A.WordLength;
```

Open the Model

```
model = 'ComplexPartialSystolicQlessQRForgettingFactorModel';
open_system(model);
```



Copyright 2022 The MathWorks, Inc.

Set Variables in the Model Workspace

Use the helper function setModelWorkspace to add the variables defined above to the model workspace.

```
fixed.example.setModelWorkspace(model,'A',A,'n',n,...
    'forgettingFactor',forgettingFactor,...
    'regularizationParameter',0,...
    'rowDelay',rowDelay,...
    'stopTime',stopTime);
```

Simulate the Model

```
out = sim(model);
```

Verify the Accuracy of the Output

Define matrix A_k as follows

$$A_k = \begin{bmatrix} \alpha^k & & & \\ & \alpha^{k-1} & & \\ & & \ddots & \\ & & & \alpha \end{bmatrix} A(1:k, :).$$

Then using the formula for the computation of the k th output R_k , and the fact that $[Q, R] = \text{qr}(A, 0) \Rightarrow A'A = R'Q'QR = R'R$, you can show that

$$A'_k A_k = R'_k R_k.$$

So to verify the output, the difference between $A'_k A_k$ and $R'_k R_k$ should be small.

Choose the last output of the simulation.

```
R = double(out.R(:, :, end))
```

```
R =
```

```
Columns 1 through 4
```

```
7.8025 + 0.0000i    0.3158 + 0.0965i   -0.0992 + 0.2743i    0.7696 + 0.2507i
0.0000 + 0.0000i    7.8339 + 0.0000i   -0.3554 - 0.4953i   -0.4352 - 0.3383i
0.0000 + 0.0000i    0.0000 + 0.0000i    8.1951 + 0.0000i   -0.2289 - 0.0269i
0.0000 + 0.0000i    0.0000 + 0.0000i    0.0000 + 0.0000i    8.0860 + 0.0000i
0.0000 + 0.0000i    0.0000 + 0.0000i    0.0000 + 0.0000i    0.0000 + 0.0000i
```

```
Column 5
```

```
-0.2899 - 0.2469i
-0.4371 - 0.7667i
-0.3997 + 0.4916i
0.1868 + 0.0171i
7.9768 + 0.0000i
```

Verify that R is upper triangular.

```
isequal(R, triu(R))
```

```
ans =
```

```
logical
```

```
1
```

Verify that the diagonal is greater than or equal to zero.

```
diag(R)
```

```
ans =
    7.8025
    7.8339
    8.1951
    8.0860
    7.9768
```

Synchronize the last output R with the input by finding the number of inputs that produced it.

```
A = double(A);
alpha = double(forgettingFactor);
relative_errors = nan(1,n);
for k = 1:numInputs
    A_k = alpha.^(k:-1:1)' .* A(1:k,:);
    relative_errors(k) = norm(A_k'*A_k - R'*R)/norm(A_k'*A_k);
end
```

k is the number of inputs A(k,:) that produced the last R.

```
k = find(relative_errors==min(relative_errors),1,'last')
```

```
k =
    500
```

Verify that

$$A_k' A_k = R_k' R_k$$

with a small relative error.

```
A_k = alpha.^(k:-1:1)' .* A(1:k,:);
relative_error = norm(A_k'*A_k - R'*R)/norm(A_k'*A_k)
```

```
relative_error =
    6.8663e-06
```

Suppress mlint warnings in this file.

```
 %#ok<*NOPTS>
```

See Also

Complex Partial-Systolic Q-less QR Decomposition with Forgetting Factor

Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Forgetting Factor

This example shows how to use the hardware-efficient Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Forgetting Factor block.

Q-less QR Decomposition with Forgetting Factor

The Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Forgetting Factor block implements the following recursion to compute the upper-triangular factor R of continuously streaming n-by-1 row vectors $A(k,:)$ using forgetting factor α . It is as if matrix A is infinitely tall. The forgetting factor in the range $0 < \alpha < 1$ keeps it from integrating without bound.

$$\begin{aligned}
 R_0 &= \text{zeros}(n, n) \\
 [\sim, R_1] &= \text{qr} \left(\begin{bmatrix} R_0 \\ A(1,:) \end{bmatrix}, 0 \right) \\
 R_1 &= \alpha R_1 \\
 \\
 [\sim, R_2] &= \text{qr} \left(\begin{bmatrix} R_1 \\ A(2,:) \end{bmatrix}, 0 \right) \\
 R_2 &= \alpha R_2 \\
 \\
 &\vdots \\
 \\
 [\sim, R_k] &= \text{qr} \left(\begin{bmatrix} R_{k-1} \\ A(k,:) \end{bmatrix}, 0 \right) \\
 R_k &= \alpha R_k \\
 \\
 &\vdots
 \end{aligned}$$

Forward and Backward Substitution

When an upper triangular factor is ready, then forward and backward substitution are computed with the current input B to produce output X.

$$X = R_k \setminus (R'_k \setminus B)$$

AMBA AXI Handshaking Process

The Data Handler subsystem in this model takes real matrices A and B as inputs. It sends rows of A and full matrix of B to the QR Decomposition block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set delays for the feeding in rows of A and the feeding in of B matrices in the Data Handler to emulate the processing time of the upstream block. `validInA` and `validInB` remain high when `aDelay` and `bDelay` are set to 0 because this indicates the Data Handler always has data available.

Asynchronous Matrix Solver

This block operates asynchronously. First, Q-less QR decomposition is performed on the input A matrix and the resulting R matrix is put into a buffer. Then, the Forward Backward Substitute block uses the input B matrix and the buffered R matrix to compute $R'RX = B$. Because the R and B matrices are stored separately in buffers, the upstream Q-less QR decomposition block and the downstream Forward Backward Substitute block can run independently. The Forward Backward Substitute block starts processing when the first R and B matrices are available. Then it runs continuously using the latest buffered R and B matrices, regardless of the status of the Q-less QR Decomposition block. For example, if the upstream block stops providing A and B matrices, the Forward Backward Substitute block continues to generate the same output using the last pair of R and B matrices.

Define System Parameters

n is the length of the row vectors $A(k,:)$, the number of rows in B, and the number of rows and columns in R.

```
n = 5;
```

p is the number of columns in B.

```
p = 1;
```

m is the effective number of rows of A to integrate over.

```
m = 100;
```

Use the `fixed.forgettingFactor` function to compute the forgetting factor as a function of the number of rows that you are integrating over.

```
forgettingFactor = fixed.forgettingFactor(m)
```

```
forgettingFactor =
```

```
    0.9950
```

`precisionBits` defines the number of bits of precision required for the QR decomposition. Set this value according to system requirements.

```
precisionBits = 24;
```

In this example, real-valued matrices A and B are constructed such that the magnitude of their elements are less than or equal to one. Your own system requirements will define what those values are. If you don't know what they are, and A and B are fixed-point input to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of A and B.

`max_abs_A` is an upper bound on the maximum magnitude element of A.

```
max_abs_A = 1;
```

`max_abs_B` is an upper bound on the maximum magnitude element of B.

```
max_abs_B = 1;
```

Select Fixed-Point Types

Use the `fixed.realQlessQRMatrixSolveFixedpointTypes` function to compute fixed-point types.

```
T = fixed.realQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,precisionBits);
```

T.A is the fixed-point type computed for transforming A to R in-place so that it does not overflow.

T.A

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 31  
    FractionLength: 24
```

T.B is the type computed for B so that it does not overflow.

T.B

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 27  
    FractionLength: 24
```

T.X is the type computed for the output X so that there is a low probability of overflow.

T.X

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 76  
    FractionLength: 24
```

Define Simulation Parameters

Create random matrix A to contain a specified number of inputs, and n-by-p random matrix B.

`numInputs` is the number of input rows `A(k,:)` for this example.

```
numInputs = 500;  
rng('default')  
[A,B] = fixed.example.realRandomQlessQRMatrices(numInputs,n,p);
```

Cast the inputs to the types determined by `fixed.realQlessQRMatrixSolveFixedpointTypes`

```
A = cast(A, 'like', T.A);
B = cast(B, 'like', T.B);
```

Use the `fixed.extractNumericType` function to extract a `numericType` object to use as an input parameter to the block.

```
OutputType = fixed.extractNumericType(T.X)
```

OutputType =

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 76
FractionLength: 24
```

Cast the forgetting factor to a fixed-point type with the same word length as A and best-precision scaling.

```
forgettingFactor = fi(forgettingFactor,1,T.A.WordLength);
```

Set delay for feeding in rows of A.

```
aDelay = 1;
```

Set delay for feeding in B matrices.

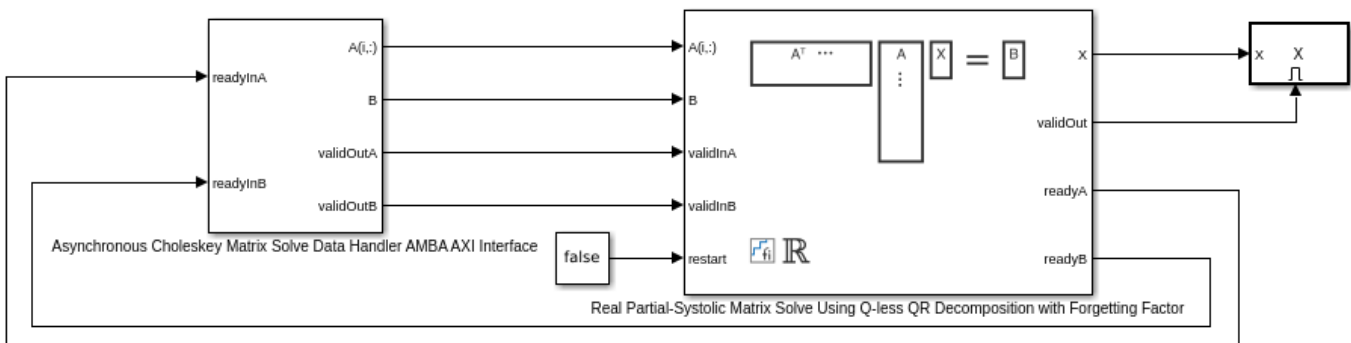
```
bDelay = 1;
```

Select a stop time for the simulation that is long enough to process all the inputs from A.

```
stopTime = (2*numInputs + n)*T.A.WordLength;
```

Open the Model

```
model = 'RealPartialSystolicMatrixSolveQlessQRForgettingFactorModel';
open_system(model);
```



Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace.

```
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'n', n, 'p', p, ...
    'forgettingFactor', forgettingFactor, 'OutputType', OutputType, ...
    'regularizationParameter', 0, ...
    'aDelay', aDelay, 'bDelay', bDelay, ...
    'stopTime', stopTime);
```

Simulate the Model

```
out = sim(model);
```

Verify the Accuracy of the Output

Define matrix A_k as follows

$$A_k = \begin{bmatrix} \alpha^k & & & \\ & \alpha^{k-1} & & \\ & & \ddots & \\ & & & \alpha \end{bmatrix} A(1:k, :).$$

Then using the formula for the computation of the k th output R_k , and the fact that $[Q, R] = \text{qr}(A, 0) \Rightarrow A'A = R'Q'QR = R'R$, you can show that

$$A_k' A_k X = R_k' R_k X = B.$$

So to verify the output, the difference between $A_k' A_k X$ and B should be small.

Choose the last output of the simulation.

```
X = double(out.X(:, :, end));
```

Synchronize the last output X with the input by finding the number of inputs that produced it.

```
A = double(A);
B = double(B);
alpha = double(forgettingFactor);
relative_errors = nan(1, numInputs);
for k = 1:numInputs
    A_k = alpha.^(k:-1:1)' .* A(1:k, :);
    relative_errors(k) = norm(A_k'*A_k*X - B)/norm(B);
end
```

k is the number of inputs $A(k, :)$ that produced the last X .

```
k = find(relative_errors==min(relative_errors))
```

```
k =
```

```
500
```

Verify that

$$A_k' A_k X = B$$

with a small relative error.

```
A_k = alpha.^(k:-1:1)' .* A(1:k,:);  
relative_error = norm(A_k'*A_k*X - B)/norm(B)
```

```
relative_error =
```

```
8.7165e-04
```

Suppress mlint warnings in this file.

```
 %#ok<*NOPTS>
```

See Also

Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Forgetting Factor

Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Forgetting Factor

This example shows how to use the hardware-efficient Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Forgetting Factor block.

Q-less QR Decomposition with Forgetting Factor

The Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Forgetting Factor block implements the following recursion to compute the upper-triangular factor R of continuously streaming n-by-1 row vectors A(k,:) using forgetting factor α . It is as if matrix A is infinitely tall. The forgetting factor in the range $0 < \alpha < 1$ keeps it from integrating without bound.

$$\begin{aligned}
 R_0 &= \text{zeros}(n, n) \\
 [\sim, R_1] &= \text{qr} \left(\begin{bmatrix} R_0 \\ A(1, :) \end{bmatrix}, 0 \right) \\
 R_1 &= \alpha R_1 \\
 \\
 [\sim, R_2] &= \text{qr} \left(\begin{bmatrix} R_1 \\ A(2, :) \end{bmatrix}, 0 \right) \\
 R_2 &= \alpha R_2 \\
 \\
 &\vdots \\
 \\
 [\sim, R_k] &= \text{qr} \left(\begin{bmatrix} R_{k-1} \\ A(k, :) \end{bmatrix}, 0 \right) \\
 R_k &= \alpha R_k \\
 \\
 &\vdots
 \end{aligned}$$

Forward and Backward Substitution

When an upper triangular factor is ready, then forward and backward substitution are computed with the current input B to produce output X.

$$X = R_k \setminus (R'_k \setminus B)$$

AMBA AXI Handshaking Process

The Data Handler subsystem in this model takes real matrices A and B as inputs. It sends rows of A and full matrix of B to the QR Decomposition block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set delays for the feeding in rows of A and the feeding in of B matrices in the Data Handler to emulate the processing time of the upstream block. `validInA` and `validInB` remain high when `aDelay` and `bDelay` are set to 0 because this indicates the Data Handler always has data available.

Asynchronous Matrix Solver

This block operates asynchronously. First, Q-less QR decomposition is performed on the input A matrix and the resulting R matrix is put into a buffer. Then, the Forward Backward Substitute block uses the input B matrix and the buffered R matrix to compute $R'RX = B$. Because the R and B matrices are stored separately in buffers, the upstream Q-less QR decomposition block and the downstream Forward Backward Substitute block can run independently. The Forward Backward Substitute block starts processing when the first R and B matrices are available. Then it runs continuously using the latest buffered R and B matrices, regardless of the status of the Q-less QR Decomposition block. For example, if the upstream block stops providing A and B matrices, the Forward Backward Substitute block continues to generate the same output using the last pair of R and B matrices.

Define System Parameters

`n` is the length of the row vectors $A(k,:)$, the number of rows in B, and the number of rows and columns in R.

```
n = 5;
```

`p` is the number of columns in B.

```
p = 1;
```

`m` is the effective number of rows of A to integrate over.

```
m = 100;
```

Use the `fixed.forgettingFactor` function to compute the forgetting factor as a function of the number of rows that you are integrating over.

```
forgettingFactor = fixed.forgettingFactor(m)
```

```
forgettingFactor =
```

```
    0.9950
```

`precisionBits` defines the number of bits of precision required for the QR decomposition. Set this value according to system requirements.

```
precisionBits = 24;
```

In this example, complex-valued matrices A and B are constructed such that the magnitude of the real and imaginary parts of their elements is less than or equal to one, so the maximum possible absolute value of any element is $|1 + 1i| = \sqrt{2}$. Your own system requirements will define what those values are. If you don't know what they are, and A and B are fixed-point input to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of A and B.

`max_abs_A` is an upper bound on the maximum magnitude element of A.

```
max_abs_A = sqrt(2);
```

`max_abs_B` is an upper bound on the maximum magnitude element of B.

```
max_abs_B = sqrt(2);
```

Select Fixed-Point Types

Use the `fixed.complexQlessQRMatrixSolveFixedpointTypes` function to compute fixed-point types.

```
T = fixed.complexQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,precisionBits);
```

T.A is the fixed-point type computed for transforming A to R in-place so that it does not overflow.

T.A

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 31  
    FractionLength: 24
```

T.B is the type computed for B so that it does not overflow.

T.B

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 27  
    FractionLength: 24
```

T.X is the type computed for the output X so that there is a low probability of overflow.

T.X

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 75  
    FractionLength: 24
```

Define Simulation Parameters

Create random matrix A to contain a specified number of inputs, and n-by-p random matrix B.

`numInputs` is the number of input rows `A(k,:)` for this example.

```
numInputs = 500;  
rng('default')  
[A,B] = fixed.example.complexRandomQlessQRMatrices(numInputs,n,p);
```


Cast the inputs to the types determined by `fixed.complexQlessQRMatrixSolveFixedpointTypes`.

```
A = cast(A, 'like', T.A);
B = cast(B, 'like', T.B);
```

Use the `fixed.extractNumericType` function to extract a `numericType` object to use as an input parameter to the block.

```
OutputType = fixed.extractNumericType(T.X)
```

```
OutputType =
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 75
    FractionLength: 24
```

Cast the forgetting factor to a fixed-point type with the same word length as A and best-precision scaling.

```
forgettingFactor = fi(forgettingFactor, 1, T.A.WordLength);
```

Set delay for feeding in rows of A.

```
aDelay = 1;
```

Set delay for feeding in B matrices.

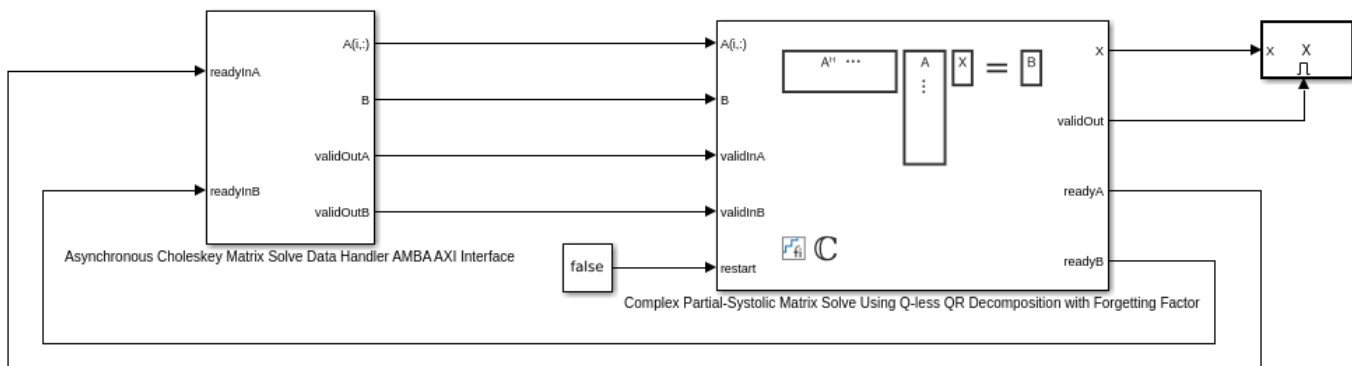
```
bDelay = 1;
```

Select a stop time for the simulation that is long enough to process all the inputs from A.

```
stopTime = 2*(2*numInputs + n)*T.A.WordLength;
```

Open the Model

```
model = 'ComplexPartialSystolicSolveQlessQRForgettingFactorModel';
open_system(model);
```



Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace.

```
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'n', n, 'p', p, ...
    'forgettingFactor', forgettingFactor, 'OutputType', OutputType, ...
    'regularizationParameter', 0, ...
    'aDelay', aDelay, 'bDelay', bDelay, ...
    'stopTime', stopTime);
```

Simulate the Model

```
out = sim(model);
```

Verify the Accuracy of the Output

Define matrix A_k as follows

$$A_k = \begin{bmatrix} \alpha^k & & & \\ & \alpha^{k-1} & & \\ & & \ddots & \\ & & & \alpha \end{bmatrix} A(1:k, :).$$

Then using the formula for the computation of the k th output R_k , and the fact that $[Q, R] = \text{qr}(A, 0) \Rightarrow A'A = R'Q'QR = R'R$, you can show that

$$A_k' A_k X = R_k' R_k X = B.$$

So to verify the output, the difference between $A_k' A_k X$ and B should be small.

Choose the last output of the simulation.

```
X = double(out.X(:, :, end));
```

Synchronize the last output X with the input by finding the number of inputs that produced it.

```
A = double(A);
B = double(B);
alpha = double(forgettingFactor);
relative_errors = nan(1, numInputs);
for k = 1:numInputs
    A_k = alpha.^(k:-1:1)' .* A(1:k, :);
    relative_errors(k) = norm(A_k'*A_k*X - B)/norm(B);
end
```

k is the number of inputs $A(k, :)$ that produced the last X .

```
k = find(relative_errors==min(relative_errors))
```

```
k =
```

```
500
```

Verify that

$$A_k' A_k X = B$$

with a small relative error.

```
A_k = alpha.^(k:-1:1)' .* A(1:k,:);  
relative_error = norm(A_k'*A_k*X - B)/norm(B)
```

```
relative_error =
```

```
4.1749e-05
```

Suppress mlint warnings in this file.

```
 %#ok<*NOPTS>
```

See Also

Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Forgetting Factor

Implement Hardware-Efficient Complex Burst Matrix Solve Using QR Decomposition with Tikhonov Regularization

This example shows how to use the Complex Burst Matrix Solve Using QR Decomposition block to solve the regularized least-squares matrix equation

$$\begin{bmatrix} \lambda I_n \\ A \end{bmatrix} X = \begin{bmatrix} 0_{n,p} \\ B \end{bmatrix},$$

where A is an m-by-n matrix with $m \geq n$, B is m-by-p, X is n-by-p, $I_n = \text{eye}(n)$, $0_{n,p} = \text{zeros}(n,p)$, and λ is a regularization parameter.

The least-squares solution is

$$X_{ls} = (\lambda^2 I_n + A^H A)^{-1} A^H B$$

but is computed without squares or inverses.

Define Matrix Dimensions

Specify the number of rows in matrices A and B, the number of columns in matrix A, and the number of columns in matrix B.

```
m = 100; % Number of rows in matrices A and B
n = 10;  % Number of columns in matrix A
p = 1;  % Number of columns in matrix B
```

Define Tikhonov Regularization Parameter

Small, positive values of the regularization parameter can improve the conditioning of the problem and reduce the variance of the estimates. While biased, the reduced variance of the estimate often results in a smaller mean squared error when compared to least-squares estimates.

```
regularizationParameter = 0.01;
```

Generate Random Least-Squares Matrices

For this example, use the helper function `complexRandomLeastSquaresMatrices` to generate random matrices A and B for the least-squares problem $AX=B$. The matrices are generated such that the real and imaginary parts of the elements of A and B are between -1 and +1, and A has rank r.

```
rng('default')
r = 3; % Rank of A
[A,B] = fixed.example.complexRandomLeastSquaresMatrices(m,n,p,r);
```

Select Fixed-Point Data Types

Use the helper function `complexQRMatrixSolveFixedpointTypes` to select fixed-point data types for input matrices A and B, and output X such that there is a low probability of overflow during the computation.

The real and imaginary parts of the elements of A and B are between -1 and 1, so the maximum possible absolute value of any element is $\sqrt{2}$.

```

max_abs_A = sqrt(2); % Upper bound on max(abs(A(:)))
max_abs_B = sqrt(2); % Upper bound on max(abs(B(:)))
precisionBits = 32; % Number of bits of precision
T = fixed.complexQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,[],[],regularizationParameter);
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);

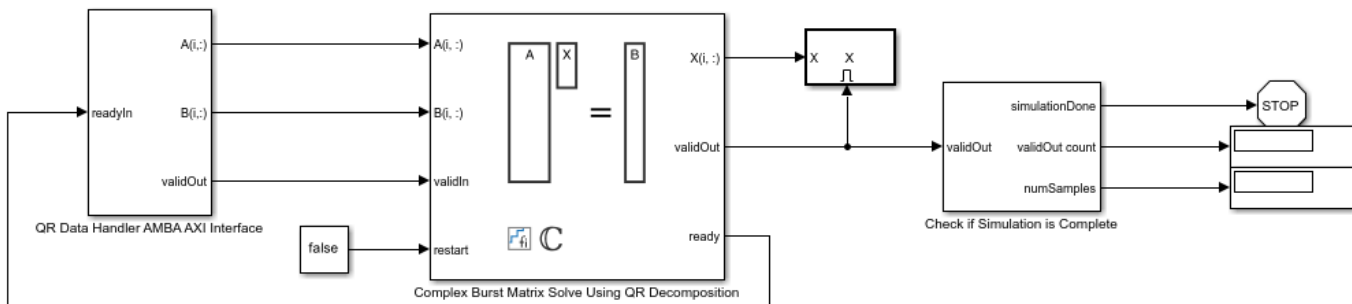
```

Open the Model

```

model = 'ComplexBurstQRMatrixSolveModel';
open_system(model);

```



Copyright 2021 The MathWorks, Inc.

The Data Handler subsystem in this model takes complex matrices A and B as inputs. It sends rows of A and B to QR block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set a delay between the feeding in rows of A and B in the Data Handler to emulate the processing time of the upstream block. `validIn` remains high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Complex Burst Matrix Solve Using QR Decomposition block.

```

numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A and B
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
    'regularizationParameter', regularizationParameter, ...
    'numSamples', numSamples, 'rowDelay', rowDelay, 'OutputType', OutputType);

```

Simulate the Model

```

out = sim(model);

```

Construct the Solution from the Output Data

The Complex Burst Matrix Solve Using QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of X are output in the order they are computed, last row first, so you must reconstruct the data to interpret the results. To

reconstruct the matrix X from the output data, use the helper function `matrixSolveModelOutputToArray`.

```
X = fixed.example.matrixSolveModelOutputToArray(out.X,n,p,numSamples);
```

Verify the Accuracy of the Output

Verify that the relative error between the fixed-point output and builtin MATLAB in double-precision floating-point is small.

$$X_{\text{double}} = \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \setminus \begin{bmatrix} 0_{n,p} \\ B \end{bmatrix}$$

```
A_lambda = double([regularizationParameter*eye(n);A]);  
B_0 = [zeros(n,p);double(B)];  
X_double = A_lambda\B_0;  
relativeError = norm(X_double - double(X))/norm(X_double)
```

```
relativeError =
```

```
1.3921e-05
```

Suppress `mlint` warnings in this file.

```
##ok<*NOPTS>  
##ok<*NASGU>  
##ok<*ASGLU>
```

See Also

Complex Burst Matrix Solve Using QR Decomposition

Implement Hardware-Efficient Complex Burst Matrix Solve Using Q-less QR Decomposition with Tikhonov Regularization

This example shows how to use the Complex Burst Matrix Solve Using QR Decomposition block to solve the regularized least-squares matrix equation

$$\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}^H \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} X = (\lambda^2 I_n + A^H A) X = B$$

where A is an m -by- n matrix with $m \geq n$, B is n -by- p , X is n -by- p , $I_n = \text{eye}(n)$, and λ is a regularization parameter.

Define System Parameters

Define the matrix attributes and system parameters for this example.

m is the number of rows in matrix A . In a problem such as beamforming or direction finding, m corresponds to the number of samples that are integrated over.

```
m = 100;
```

n is the number of columns in matrix A and rows in matrices B and X . In a least-squares problem, m is greater than n , and usually m is much larger than n . In a problem such as beamforming or direction finding, n corresponds to the number of sensors.

```
n = 10;
```

p is the number of columns in matrices B and X . It corresponds to simultaneously solving a system with p right-hand sides.

```
p = 1;
```

In this example, set the rank of matrix A to be less than the number of columns. In a problem such as beamforming or direction finding, $\text{rank}(A)$ corresponds to the number of signals impinging on the sensor array.

```
rankA = 3;
```

`precisionBits` defines the number of bits of precision required for the matrix solve. Set this value according to system requirements.

```
precisionBits = 32;
```

Small, positive values of the regularization parameter can improve the conditioning of the problem and reduce the variance of the estimates. While biased, the reduced variance of the estimate often results in a smaller mean squared error when compared to least-squares estimates.

```
regularizationParameter = 0.01;
```

In this example, complex-valued matrices A and B are constructed such that the magnitude of the real and imaginary parts of their elements is less than or equal to one, so the maximum possible absolute value of any element is $|1 + 1i| = \sqrt{2}$. Your own system requirements will define what those values

are. If you don't know what they are, and A and B are fixed-point inputs to the system, then you can use the upperbound function to determine the upper bounds of the fixed-point types of A and B.

max_abs_A is an upper bound on the maximum magnitude element of A.

```
max_abs_A = sqrt(2);
```

max_abs_B is an upper bound on the maximum magnitude element of B.

```
max_abs_B = sqrt(2);
```

Thermal noise standard deviation is the square root of thermal noise power, which is a system parameter. A well-designed system has the quantization level lower than the thermal noise. Here, set thermalNoiseStandardDeviation to the equivalent of -50 dB noise power.

```
thermalNoiseStandardDeviation = sqrt(10^(-50/10))
```

```
thermalNoiseStandardDeviation =
```

```
0.0032
```

The quantization noise standard deviation is a function of the required number of bits of precision. Use fixed.complexQuantizationNoiseStandardDeviation to compute this. See that it is less than thermalNoiseStandardDeviation.

```
quantizationNoiseStandardDeviation = fixed.complexQuantizationNoiseStandardDeviation(precisionBits)
```

```
quantizationNoiseStandardDeviation =
```

```
9.5053e-11
```

Compute Fixed-Point Types

In this example, assume that the designed system matrix A does not have full rank (there are fewer signals of interest than number of columns of matrix A), and the measured system matrix A has additive thermal noise that is larger than the quantization noise. The additive noise makes the measured matrix A have full rank.

Set $\sigma_{\text{noise}} = \sigma_{\text{thermal noise}}$.

```
noiseStandardDeviation = thermalNoiseStandardDeviation;
```

Use the fixed.complexQlessQRMatrixSolveFixedpointTypes function to compute fixed-point types.

```
T = fixed.complexQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,noiseStandardDeviation,[],regularizationParameter)
```

```
T =
```

```
struct with fields:
```

```
A: [0x0 embedded.fi]
```



```
B: [0x0 embedded.fi]
X: [0x0 embedded.fi]
```

T.A is the type computed for transforming $\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}$ to $R = Q^H \begin{bmatrix} \lambda I_n \\ A \end{bmatrix}$ in-place so that it does not overflow.

T.A

ans =

[]

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 39
FractionLength: 32
```

T.B is the type computed for B so that it does not overflow.

T.B

ans =

[]

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 35
FractionLength: 32
```

T.X is the type computed for the solution $X = \left(\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}^H \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \right) \setminus B$ so that there is a low probability that it overflows.

T.X

ans =

[]

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 50
FractionLength: 32
```

Use the Specified Types to Solve the Matrix Equation

Create random matrices A and B such that rankA=rank(A). Add random measurement noise to A which will make it become full rank.

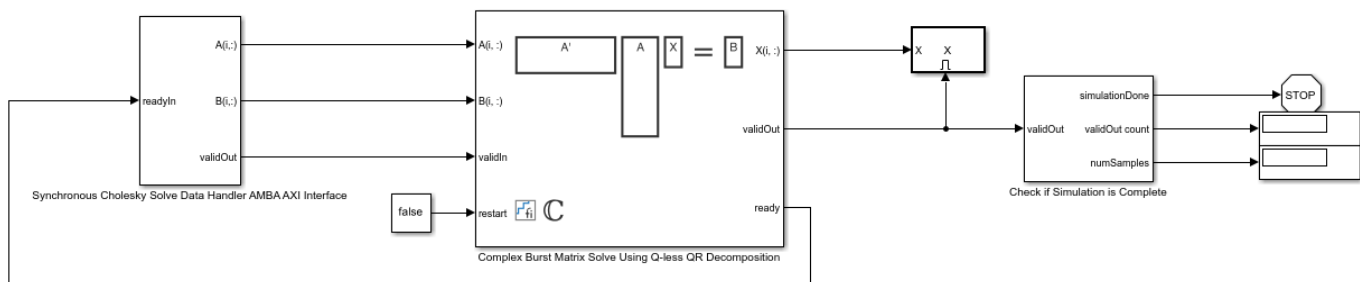
```
rng('default');
[A,B] = fixed.example.complexRandomQlessQRMatrices(m,n,p,rankA);
A = A + fixed.example.complexNormalRandomArray(0,noiseStandardDeviation,m,n);
```

Cast the inputs to the types determined by `fixed.complexQlessQRMatrixSolveFixedpointTypes`. Quantizing to fixed-point is equivalent to adding random noise.

```
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);
```

Open the Model

```
model = 'ComplexBurstQlessQRMatrixSolveModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Complex Burst Matrix Solve Using QR Decomposition block.

```
numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A and B
fixed.example.setModelWorkspace(model,'A',A,'B',B,'m',m,'n',n,'p',p,...
    'regularizationParameter',regularizationParameter,...
    'numSamples',numSamples,'rowDelay',rowDelay,'OutputType',OutputType);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Complex Burst Matrix Solve Using Q-less QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of `X` are output in the order they are computed, last row first, so you must reconstruct the data to interpret the results. To reconstruct the matrix `X` from the output data, use the helper function `matrixSolveModelOutputToArray`.

```
X = fixed.example.matrixSolveModelOutputToArray(out.X,n,p);
```

Verify the Accuracy of the Output

Verify that the relative error between the fixed-point output and builtin MATLAB in double-precision floating-point is small.

$$X_{\text{double}} = \left(\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}^H \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \right) \backslash B$$

```
A_lambda = double([regularizationParameter*eye(n);A]);
X_double = (A_lambda'*A_lambda)\double(B);
relativeError = norm(X_double - double(X))/norm(X_double)
```

```
relativeError =
```

```
5.3708e-06
```

Suppress mlint warnings in this file.

```
 %#ok<*NASGU>
 %#ok<*ASGLU>
 %#ok<*NOPTS>
```

See Also

Complex Burst Matrix Solve Using Q-less QR Decomposition

Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using QR Decomposition with Tikhonov Regularization

This example shows how to use the Complex Partial-Systolic Matrix Solve Using QR Decomposition block to solve the regularized least-squares matrix equation

$$\begin{bmatrix} \lambda I_n \\ A \end{bmatrix} X = \begin{bmatrix} 0_{n,p} \\ B \end{bmatrix},$$

where A is an m-by-n matrix with $m \geq n$, B is m-by-p, X is n-by-p, $I_n = \text{eye}(n)$, $0_{n,p} = \text{zeros}(n,p)$, and λ is a regularization parameter.

The least-squares solution is

$$X_{ls} = (\lambda^2 I_n + A^H A)^{-1} A^H B$$

but is computed without squares or inverses.

Define Matrix Dimensions

Specify the number of rows in matrices A and B, the number of columns in matrix A, and the number of columns in matrix B.

```
m = 100; % Number of rows in matrices A and B
n = 10;  % Number of columns in matrix A
p = 1;  % Number of columns in matrix B
```

Define Tikhonov Regularization Parameter

Small, positive values of the regularization parameter can improve the conditioning of the problem and reduce the variance of the estimates. While biased, the reduced variance of the estimate often results in a smaller mean squared error when compared to least-squares estimates.

```
regularizationParameter = 0.01;
```

Generate Random Least-Squares Matrices

For this example, use the helper function `complexRandomLeastSquaresMatrices` to generate random matrices A and B for the least-squares problem $AX=B$. The matrices are generated such that the real and imaginary parts of the elements of A and B are between -1 and +1, and A has rank r.

```
rng('default')
r = 3; % Rank of A
[A,B] = fixed.example.complexRandomLeastSquaresMatrices(m,n,p,r);
```

Select Fixed-Point Data Types

Use the helper function `complexQRMatrixSolveFixedpointTypes` to select fixed-point data types for input matrices A and B, and output X such that there is a low probability of overflow during the computation.

The real and imaginary parts of the elements of A and B are between -1 and 1, so the maximum possible absolute value of any element is $\sqrt{2}$.

```

max_abs_A = sqrt(2); % Upper bound on max(abs(A(:)))
max_abs_B = sqrt(2); % Upper bound on max(abs(B(:)))
precisionBits = 32; % Number of bits of precision
T = fixed.complexQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,[],[],regularizationParameter);
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);

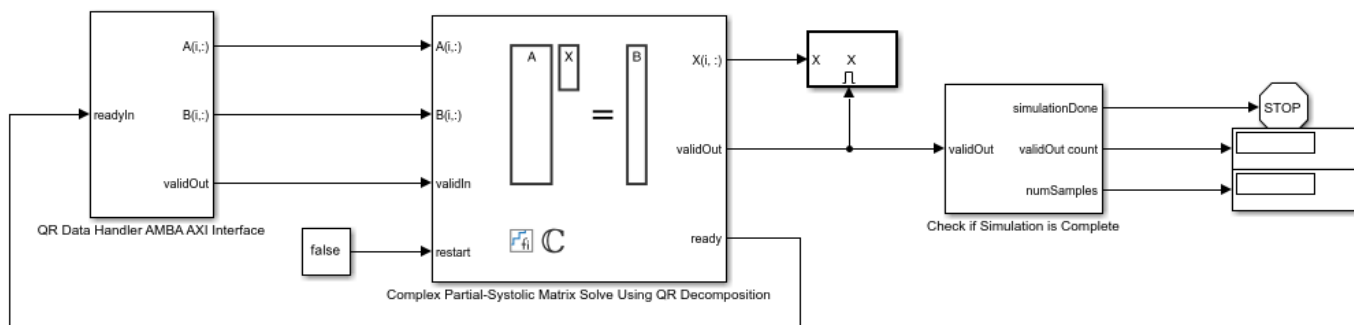
```

Open the Model

```

model = 'ComplexPartialSystolicQRMatrixSolveModel';
open_system(model);

```



Copyright 2021 The MathWorks, Inc.

The Data Handler subsystem in this model takes complex matrices A and B as inputs. It sends rows of A and B to QR block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set a delay between the feeding in rows of A and B in the Data Handler to emulate the processing time of the upstream block. `validIn` remains high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Complex Partial-Systolic Matrix Solve Using QR Decomposition block.

```

numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A and B
fixed.example.setModelWorkspace(model,'A',A,'B',B,'m',m,'n',n,'p',p,...
    'regularizationParameter',regularizationParameter,...
    'numSamples',numSamples,'rowDelay',rowDelay,'OutputType',OutputType);

```

Simulate the Model

```

out = sim(model);

```

Construct the Solution from the Output Data

The Complex Partial-Systolic Matrix Solve Using QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of X are output in the

order they are computed, last row first, so you must reconstruct the data to interpret the results. To reconstruct the matrix X from the output data, use the helper function `matrixSolveModelOutputToArray`.

```
X = fixed.example.matrixSolveModelOutputToArray(out.X,n,p,numSamples);
```

Verify the Accuracy of the Output

Verify that the relative error between the fixed-point output and builtin MATLAB in double-precision floating-point is small.

$$X_{\text{double}} = \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \setminus \begin{bmatrix} 0_{n,p} \\ B \end{bmatrix}$$

```
A_lambda = double([regularizationParameter*eye(n);A]);  
B_0 = [zeros(n,p);double(B)];  
X_double = A_lambda\B_0;  
relativeError = norm(X_double - double(X))/norm(X_double)
```

```
relativeError =
```

```
1.3921e-05
```

Suppress `mLint` warnings in this file.

```
##ok<*NOPTS>  
##ok<*NASGU>  
##ok<*ASGLU>
```

See Also

Complex Partial-Systolic Matrix Solve Using QR Decomposition

Implement Hardware-Efficient Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Tikhonov Regularization

This example shows how to use the Complex Partial-Systolic Matrix Solve Using QR Decomposition block to solve the regularized least-squares matrix equation

$$\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}^H \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} X = (\lambda^2 I_n + A^H A) X = B$$

where A is an m -by- n matrix with $m \geq n$, B is n -by- p , X is n -by- p , $I_n = \text{eye}(n)$, and λ is a regularization parameter.

Define System Parameters

Define the matrix attributes and system parameters for this example.

m is the number of rows in matrix A . In a problem such as beamforming or direction finding, m corresponds to the number of samples that are integrated over.

`m = 100;`

n is the number of columns in matrix A and rows in matrices B and X . In a least-squares problem, m is greater than n , and usually m is much larger than n . In a problem such as beamforming or direction finding, n corresponds to the number of sensors.

`n = 10;`

p is the number of columns in matrices B and X . It corresponds to simultaneously solving a system with p right-hand sides.

`p = 1;`

In this example, set the rank of matrix A to be less than the number of columns. In a problem such as beamforming or direction finding, $\text{rank}(A)$ corresponds to the number of signals impinging on the sensor array.

`rankA = 3;`

`precisionBits` defines the number of bits of precision required for the matrix solve. Set this value according to system requirements.

`precisionBits = 32;`

Small, positive values of the regularization parameter can improve the conditioning of the problem and reduce the variance of the estimates. While biased, the reduced variance of the estimate often results in a smaller mean squared error when compared to least-squares estimates.

`regularizationParameter = 0.01;`

In this example, complex-valued matrices A and B are constructed such that the magnitude of the real and imaginary parts of their elements is less than or equal to one, so the maximum possible absolute

value of any element is $|1 + 1i| = \sqrt{2}$. Your own system requirements will define what those values are. If you don't know what they are, and A and B are fixed-point inputs to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of A and B.

`max_abs_A` is an upper bound on the maximum magnitude element of A.

```
max_abs_A = sqrt(2);
```

`max_abs_B` is an upper bound on the maximum magnitude element of B.

```
max_abs_B = sqrt(2);
```

Thermal noise standard deviation is the square root of thermal noise power, which is a system parameter. A well-designed system has the quantization level lower than the thermal noise. Here, set `thermalNoiseStandardDeviation` to the equivalent of -50 dB noise power.

```
thermalNoiseStandardDeviation = sqrt(10^(-50/10))
```

```
thermalNoiseStandardDeviation =  
    0.0032
```

The quantization noise standard deviation is a function of the required number of bits of precision. Use `fixed.complexQuantizationNoiseStandardDeviation` to compute this. See that it is less than `thermalNoiseStandardDeviation`.

```
quantizationNoiseStandardDeviation = fixed.complexQuantizationNoiseStandardDeviation(precisionBits)
```

```
quantizationNoiseStandardDeviation =  
    9.5053e-11
```

Compute Fixed-Point Types

In this example, assume that the designed system matrix **A** does not have full rank (there are fewer signals of interest than number of columns of matrix **A**), and the measured system matrix **A** has additive thermal noise that is larger than the quantization noise. The additive noise makes the measured matrix **A** have full rank.

Set $\sigma_{\text{noise}} = \sigma_{\text{thermal noise}}$.

```
noiseStandardDeviation = thermalNoiseStandardDeviation;
```

Use the `fixed.complexQlessQRMatrixSolveFixedpointTypes` function to compute fixed-point types.

```
T = fixed.complexQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...  
    precisionBits,noiseStandardDeviation,[],regularizationParameter)
```

```
T =
```

```
struct with fields:
```



```
A: [0x0 embedded.fi]
B: [0x0 embedded.fi]
X: [0x0 embedded.fi]
```

T.A is the type computed for transforming $\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}$ to $R = Q^H \begin{bmatrix} \lambda I_n \\ A \end{bmatrix}$ in-place so that it does not overflow.

T.A

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 39
    FractionLength: 32
```

T.B is the type computed for B so that it does not overflow.

T.B

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 35
    FractionLength: 32
```

T.X is the type computed for the solution $X = \left(\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}^H \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \right) \setminus B$ so that there is a low probability that it overflows.

T.X

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 50
    FractionLength: 32
```

Use the Specified Types to Solve the Matrix Equation

Create random matrices A and B such that rankA=rank(A). Add random measurement noise to A which will make it become full rank.

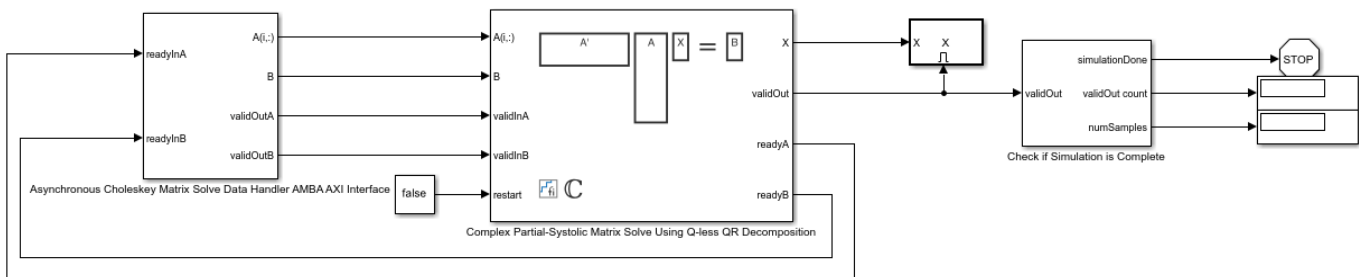
```
rng('default');
[A,B] = fixed.example.complexRandomQlessQRMatrices(m,n,p,rankA);
A = A + fixed.example.complexNormalRandomArray(0,noiseStandardDeviation,m,n);
```

Cast the inputs to the types determined by `fixed.complexQlessQRMatrixSolveFixedpointTypes`. Quantizing to fixed-point is equivalent to adding random noise.

```
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);
```

Open the Model

```
model = 'ComplexPartialSystolicQlessQRMatrixSolveModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Complex Burst Matrix Solve Using QR Decomposition block.

```
numOutputs = 1; % Number of recorded outputs
aDelay = 1; % Delay of clock cycles between feeding in rows of A
bDelay = 1; % Delay of clock cycles between feeding in B matrices
fixed.example.setModelWorkspace(model,'A',A,'B',B,'m',m,'n',n,'p',p,...
    'regularizationParameter',regularizationParameter,...
    'aDelay',aDelay,'bDelay',bDelay,...
    'numOutputs',numOutputs,'OutputType',OutputType);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition block outputs matrix `X` at each time step. When a valid result matrix is output, the block sets `validOut` to true.

```
X = out.X;
```

Verify the Accuracy of the Output

Verify that the relative error between the fixed-point output and builtin MATLAB in double-precision floating-point is small.

$$X_{\text{double}} = \left(\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}^H \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \right) \backslash B$$

```
A_lambda = double([regularizationParameter*eye(n);A]);  
X_double = (A_lambda'*A_lambda)\double(B);  
relativeError = norm(X_double - double(X))/norm(X_double)
```

```
relativeError =
```

```
5.1465e-06
```

Suppress mLint warnings in this file.

```
 %#ok< *NASGU>  
 %#ok< *ASGLU>  
 %#ok< *NOPTS>
```

See Also

Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition

Implement Hardware-Efficient Real Burst Matrix Solve Using QR Decomposition with Tikhonov Regularization

This example shows how to use the Real Burst Matrix Solve Using QR Decomposition block to solve the regularized least-squares matrix equation

$$\begin{bmatrix} \lambda I_n \\ A \end{bmatrix} X = \begin{bmatrix} 0_{n,p} \\ B \end{bmatrix},$$

where A is an m-by-n matrix with $m \geq n$, B is m-by-p, X is n-by-p, $I_n = \text{eye}(n)$, $0_{n,p} = \text{zeros}(n,p)$, and λ is a regularization parameter.

The least-squares solution is

$$X_{ls} = (\lambda^2 I_n + A^T A)^{-1} A^T B$$

but is computed without squares or inverses.

Define Matrix Dimensions

Specify the number of rows in matrices A and B, the number of columns in matrix A, and the number of columns in matrix B.

```
m = 100; % Number of rows in matrices A and B
n = 10;  % Number of columns in matrix A
p = 1;  % Number of columns in matrix B
```

Define Tikhonov Regularization Parameter

Small, positive values of the regularization parameter can improve the conditioning of the problem and reduce the variance of the estimates. While biased, the reduced variance of the estimate often results in a smaller mean squared error when compared to least-squares estimates.

```
regularizationParameter = 0.01;
```

Generate Random Least-Squares Matrices

For this example, use the helper function `realRandomLeastSquaresMatrices` to generate random matrices A and B for the least-squares problem $AX=B$. The matrices are generated such that the elements of A and B are between -1 and +1, and A has rank r.

```
rng('default')
r = 3; % Rank of A
[A,B] = fixed.example.realRandomLeastSquaresMatrices(m,n,p,r);
```

Select Fixed-Point Data Types

Use the helper function `realQRMatrixSolveFixedpointTypes` to select fixed-point data types for input matrices A and B, and output X such that there is a low probability of overflow during the computation.

```
max_abs_A = 1; % Upper bound on max(abs(A(:)))
max_abs_B = 1; % Upper bound on max(abs(B(:)))
```

```

precisionBits = 32; % Number of bits of precision
T = fixed.realQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,[],[],regularizationParameter);
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);

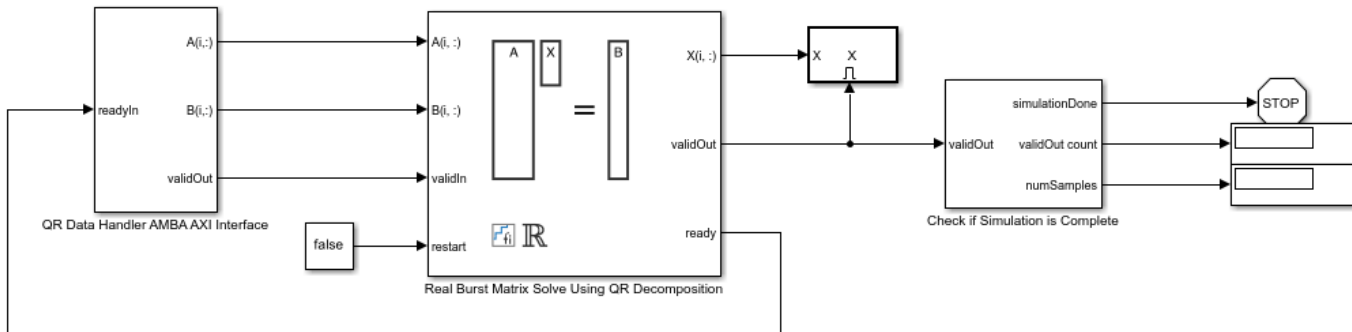
```

Open the Model

```

model = 'RealBurstQRMatrixSolveModel';
open_system(model);

```



Copyright 2021 The MathWorks, Inc.

The Data Handler subsystem in this model takes real matrices A and B as inputs. It sends rows of A and B to QR block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set a delay between the feeding in rows of A and B in the Data Handler to emulate the processing time of the upstream block. `validIn` remains high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Real Burst Matrix Solve Using QR Decomposition block.

```

numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A and B
fixed.example.setModelWorkspace(model,'A',A,'B',B,'m',m,'n',n,'p',p,...
    'regularizationParameter',regularizationParameter,...
    'numSamples',numSamples,'rowDelay',rowDelay,'OutputType',OutputType);

```

Simulate the Model

```

out = sim(model);

```

Construct the Solution from the Output Data

The Real Burst Matrix Solve Using QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of X are output in the order they are computed, last row first, so you must reconstruct the data to interpret the results. To reconstruct the matrix X from the output data, use the helper function `matrixSolveModelOutputToArray`.

```
X = fixed.example.matrixSolveModelOutputToArray(out.X,n,p,numSamples);
```

Verify the Accuracy of the Output

Verify that the relative error between the fixed-point output and builtin MATLAB in double-precision floating-point is small.

$$X_{\text{double}} = \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \setminus \begin{bmatrix} 0_{n,p} \\ B \end{bmatrix}$$

```
A_lambda = double([regularizationParameter*eye(n);A]);  
B_0 = [zeros(n,p);double(B)];  
X_double = A_lambda\B_0;  
relativeError = norm(X_double - double(X))/norm(X_double)
```

```
relativeError =
```

```
1.0202e-05
```

Suppress mlint warnings in this file.

```
#![ok<*NOPTS>  
#![ok<*NASGU>  
#![ok<*ASGLU>
```

See Also

Real Burst Matrix Solve Using QR Decomposition

Implement Hardware-Efficient Real Burst Matrix Solve Using Q-less QR Decomposition with Tikhonov Regularization

This example shows how to use the Real Burst Matrix Solve Using QR Decomposition block to solve the regularized least-squares matrix equation

$$\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}^T \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} X = (\lambda^2 I_n + A^T A) X = B$$

where A is an m -by- n matrix with $m \geq n$, B is n -by- p , X is n -by- p , $I_n = \text{eye}(n)$, and λ is a regularization parameter.

Define System Parameters

Define the matrix attributes and system parameters for this example.

m is the number of rows in matrix A . In a problem such as beamforming or direction finding, m corresponds to the number of samples that are integrated over.

`m = 100;`

n is the number of columns in matrix A and rows in matrices B and X . In a least-squares problem, m is greater than n , and usually m is much larger than n . In a problem such as beamforming or direction finding, n corresponds to the number of sensors.

`n = 10;`

p is the number of columns in matrices B and X . It corresponds to simultaneously solving a system with p right-hand sides.

`p = 1;`

In this example, set the rank of matrix A to be less than the number of columns. In a problem such as beamforming or direction finding, $\text{rank}(A)$ corresponds to the number of signals impinging on the sensor array.

`rankA = 3;`

`precisionBits` defines the number of bits of precision required for the matrix solve. Set this value according to system requirements.

`precisionBits = 32;`

Small, positive values of the regularization parameter can improve the conditioning of the problem and reduce the variance of the estimates. While biased, the reduced variance of the estimate often results in a smaller mean squared error when compared to least-squares estimates.

`regularizationParameter = 0.01;`

In this example, real-valued matrices A and B are constructed such that the magnitude of their elements is less than or equal to one. Your own system requirements will define what those values are. If you don't know what they are, and A and B are fixed-point inputs to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of A and B .

`max_abs_A` is an upper bound on the maximum magnitude element of A .

```
max_abs_A = 1;
```

`max_abs_B` is an upper bound on the maximum magnitude element of B .

```
max_abs_B = 1;
```

Thermal noise standard deviation is the square root of thermal noise power, which is a system parameter. A well-designed system has the quantization level lower than the thermal noise. Here, set `thermalNoiseStandardDeviation` to the equivalent of -50 dB noise power.

```
thermalNoiseStandardDeviation = sqrt(10^(-50/10))
```

```
thermalNoiseStandardDeviation =
```

```
    0.0032
```

The quantization noise standard deviation is a function of the required number of bits of precision. Use `fixed.realQuantizationNoiseStandardDeviation` to compute this. See that it is less than `thermalNoiseStandardDeviation`.

```
quantizationNoiseStandardDeviation = fixed.realQuantizationNoiseStandardDeviation(precisionBits)
```

```
quantizationNoiseStandardDeviation =
```

```
    6.7212e-11
```

Compute Fixed-Point Types

In this example, assume that the designed system matrix A does not have full rank (there are fewer signals of interest than number of columns of matrix A), and the measured system matrix A has additive thermal noise that is larger than the quantization noise. The additive noise makes the measured matrix A have full rank.

Set $\sigma_{\text{noise}} = \sigma_{\text{thermal noise}}$.

```
noiseStandardDeviation = thermalNoiseStandardDeviation;
```

Use the `fixed.realQlessQRMatrixSolveFixedpointTypes` function to compute fixed-point types.

```
T = fixed.realQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,noiseStandardDeviation,[],regularizationParameter)
```

```
T =
```

```
struct with fields:
```

```
    A: [0x0 embedded.fi]
```

```
    B: [0x0 embedded.fi]
```

```
    X: [0x0 embedded.fi]
```


T.A is the type computed for transforming $\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}$ to $R = Q^T \begin{bmatrix} \lambda I_n \\ A \end{bmatrix}$ in-place so that it does not overflow.

T.A

ans =

[]

```

        DataTypeMode: Fixed-point: binary point scaling
          Signedness: Signed
          WordLength: 39
        FractionLength: 32

```

T.B is the type computed for B so that it does not overflow.

T.B

ans =

[]

```

        DataTypeMode: Fixed-point: binary point scaling
          Signedness: Signed
          WordLength: 35
        FractionLength: 32

```

T.X is the type computed for the solution $X = \left(\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}^T \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \right) \setminus B$ so that there is a low probability that it overflows.

T.X

ans =

[]

```

        DataTypeMode: Fixed-point: binary point scaling
          Signedness: Signed
          WordLength: 50
        FractionLength: 32

```

Use the Specified Types to Solve the Matrix Equation

Create random matrices A and B such that rankA=rk(A). Add random measurement noise to A which will make it become full rank.

```

rng('default');
[A,B] = fixed.example.realRandomQlessQRMatrices(m,n,p,rkA);
A = A + fixed.example.realNormalRandomArray(0,noiseStandardDeviation,m,n);

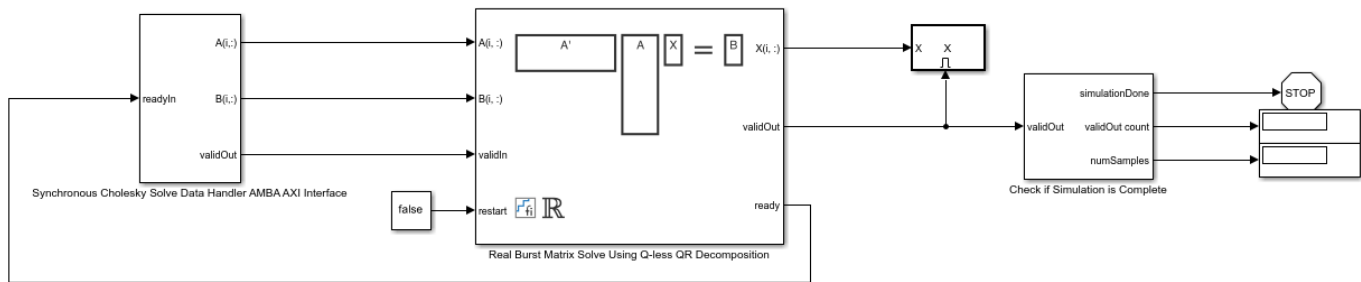
```

Cast the inputs to the types determined by fixed.realQlessQRMatrixSolveFixedpointTypes. Quantizing to fixed-point is equivalent to adding random noise.

```
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);
```

Open the Model

```
model = 'RealBurstQlessQRMatrixSolveModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Real Burst Matrix Solve Using QR Decomposition block.

```
numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A and B
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
    'regularizationParameter', regularizationParameter, ...
    'numSamples', numSamples, 'rowDelay', rowDelay, 'OutputType', OutputType);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Real Burst Matrix Solve Using Q-less QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of `X` are output in the order they are computed, last row first, so you must reconstruct the data to interpret the results. To reconstruct the matrix `X` from the output data, use the helper function `matrixSolveModelOutputToArray`.

```
X = fixed.example.matrixSolveModelOutputToArray(out.X, n, p);
```

Verify the Accuracy of the Output

Verify that the relative error between the fixed-point output and builtin MATLAB in double-precision floating-point is small.

$$X_{\text{double}} = \left(\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}^T \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \right) \setminus B$$

```
A_lambda = double([regularizationParameter*eye(n);A]);  
X_double = (A_lambda'*A_lambda)\double(B);  
relativeError = norm(X_double - double(X))/norm(X_double)
```

```
relativeError =
```

```
4.8612e-06
```

Suppress mlint warnings in this file.

```
 %#ok<*NASGU>  
 %#ok<*ASGLU>  
 %#ok<*NOPTS>
```

See Also

Real Burst Matrix Solve Using Q-less QR Decomposition

Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using QR Decomposition with Tikhonov Regularization

This example shows how to use the Real Partial-Systolic Matrix Solve Using QR Decomposition block to solve the regularized least-squares matrix equation

$$\begin{bmatrix} \lambda I_n \\ A \end{bmatrix} X = \begin{bmatrix} 0_{n,p} \\ B \end{bmatrix},$$

where A is an m-by-n matrix with $m \geq n$, B is m-by-p, X is n-by-p, $I_n = \text{eye}(n)$, $0_{n,p} = \text{zeros}(n,p)$, and λ is a regularization parameter.

The least-squares solution is

$$X_{ls} = (\lambda^2 I_n + A^T A)^{-1} A^T B$$

but is computed without squares or inverses.

Define Matrix Dimensions

Specify the number of rows in matrices A and B, the number of columns in matrix A, and the number of columns in matrix B.

```
m = 300; % Number of rows in matrices A and B
n = 10; % Number of columns in matrix A
p = 1; % Number of columns in matrix B
```

Define Tikhonov Regularization Parameter

Small, positive values of the regularization parameter can improve the conditioning of the problem and reduce the variance of the estimates. While biased, the reduced variance of the estimate often results in a smaller mean squared error when compared to least-squares estimates.

```
regularizationParameter = 0.01;
```

Generate Random Least-Squares Matrices

For this example, use the helper function `realRandomLeastSquaresMatrices` to generate random matrices A and B for the least-squares problem $AX=B$. The matrices are generated such that the elements of A and B are between -1 and +1, and A has rank r.

```
rng('default')
r = 3; % Rank of A
[A,B] = fixed.example.realRandomLeastSquaresMatrices(m,n,p,r);
```

Select Fixed-Point Data Types

Use the helper function `realQRMatrixSolveFixedpointTypes` to select fixed-point data types for input matrices A and B, and output X such that there is a low probability of overflow during the computation.

```
max_abs_A = 1; % Upper bound on max(abs(A(:)))
max_abs_B = 1; % Upper bound on max(abs(B(:)))
```

```

precisionBits = 32; % Number of bits of precision
T = fixed.realQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,[],[],regularizationParameter);
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);

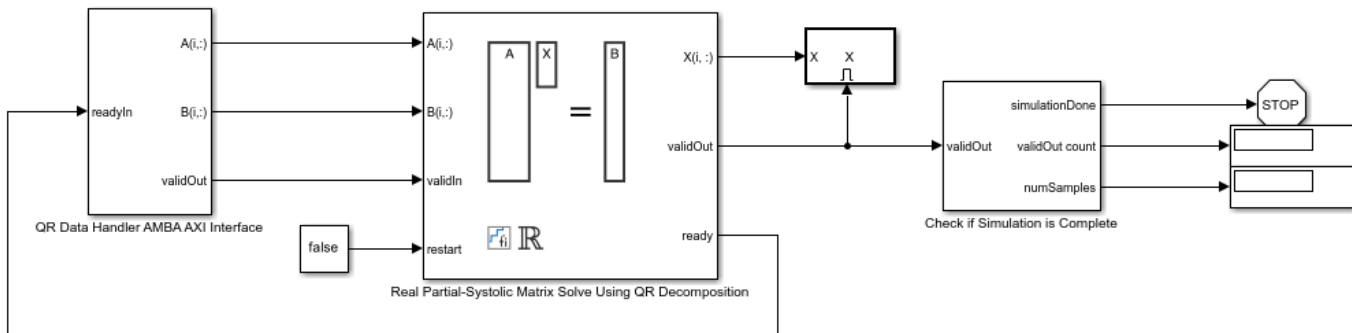
```

Open the Model

```

model = 'RealPartialSystolicQRMatrixSolveModel';
open_system(model);

```



Copyright 2021 The MathWorks, Inc.

The Data Handler subsystem in this model takes real matrices A and B as inputs. It sends rows of A and B to QR block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set a delay between the feeding in rows of A and B in the Data Handler to emulate the processing time of the upstream block. `validIn` remains high when `rowDelay` is set to 0 because this indicates the Data Handler always has data available.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Real Partial-Systolic Matrix Solve Using QR Decomposition block.

```

numSamples = 1; % Number of sample matrices
rowDelay = 1; % Delay of clock cycles between feeding in rows of A and B
fixed.example.setModelWorkspace(model,'A',A,'B',B,'m',m,'n',n,'p',p,...
    'regularizationParameter',regularizationParameter,...
    'numSamples',numSamples,'rowDelay',rowDelay,'OutputType',OutputType);

```

Simulate the Model

```

out = sim(model);

```

Construct the Solution from the Output Data

The Real Partial-Systolic Matrix Solve Using QR Decomposition block outputs data one row at a time. When a result row is output, the block sets `validOut` to true. The rows of X are output in the order they are computed, last row first, so you must reconstruct the data to interpret the results. To

reconstruct the matrix X from the output data, use the helper function `matrixSolveModelOutputToArray`.

```
X = fixed.example.matrixSolveModelOutputToArray(out.X,n,p,numSamples);
```

Verify the Accuracy of the Output

Verify that the relative error between the fixed-point output and builtin MATLAB in double-precision floating-point is small.

$$X_{\text{double}} = \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \setminus \begin{bmatrix} 0_{n,p} \\ B \end{bmatrix}$$

```
A_lambda = double([regularizationParameter*eye(n);A]);  
B_0 = [zeros(n,p);double(B)];  
X_double = A_lambda\B_0;  
relativeError = norm(X_double - double(X))/norm(X_double)
```

```
relativeError =
```

```
2.5211e-05
```

Suppress `mLint` warnings in this file.

```
##ok<*NOPTS>  
##ok<*NASGU>  
##ok<*ASGLU>
```

See Also

Real Partial-Systolic Matrix Solve Using QR Decomposition

Implement Hardware-Efficient Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Tikhonov Regularization

This example shows how to use the Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition block to solve the regularized least-squares matrix equation

$$\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}^T \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} X = (\lambda^2 I_n + A^T A) X = B$$

where A is an m -by- n matrix with $m \geq n$, B is n -by- p , X is n -by- p , $I_n = \text{eye}(n)$, and λ is a regularization parameter.

Define System Parameters

Define the matrix attributes and system parameters for this example.

m is the number of rows in matrix A . In a problem such as beamforming or direction finding, m corresponds to the number of samples that are integrated over.

```
m = 100;
```

n is the number of columns in matrix A and rows in matrices B and X . In a least-squares problem, m is greater than n , and usually m is much larger than n . In a problem such as beamforming or direction finding, n corresponds to the number of sensors.

```
n = 10;
```

p is the number of columns in matrices B and X . It corresponds to simultaneously solving a system with p right-hand sides.

```
p = 1;
```

In this example, set the rank of matrix A to be less than the number of columns. In a problem such as beamforming or direction finding, $\text{rank}(A)$ corresponds to the number of signals impinging on the sensor array.

```
rankA = 3;
```

`precisionBits` defines the number of bits of precision required for the matrix solve. Set this value according to system requirements.

```
precisionBits = 32;
```

Small, positive values of the regularization parameter can improve the conditioning of the problem and reduce the variance of the estimates. While biased, the reduced variance of the estimate often results in a smaller mean squared error when compared to least-squares estimates.

```
regularizationParameter = 0.01;
```

In this example, real-valued matrices A and B are constructed such that the magnitude of their elements is less than or equal to one. Your own system requirements will define what those values

are. If you don't know what they are, and A and B are fixed-point inputs to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of A and B.

`max_abs_A` is an upper bound on the maximum magnitude element of A.

```
max_abs_A = 1;
```

`max_abs_B` is an upper bound on the maximum magnitude element of B.

```
max_abs_B = 1;
```

Thermal noise standard deviation is the square root of thermal noise power, which is a system parameter. A well-designed system has the quantization level lower than the thermal noise. Here, set `thermalNoiseStandardDeviation` to the equivalent of -50 dB noise power.

```
thermalNoiseStandardDeviation = sqrt(10^(-50/10))
```

```
thermalNoiseStandardDeviation =
```

```
0.0032
```

The quantization noise standard deviation is a function of the required number of bits of precision. Use `fixed.realQuantizationNoiseStandardDeviation` to compute this. See that it is less than `thermalNoiseStandardDeviation`.

```
quantizationNoiseStandardDeviation = fixed.realQuantizationNoiseStandardDeviation(precisionBits)
```

```
quantizationNoiseStandardDeviation =
```

```
6.7212e-11
```

Compute Fixed-Point Types

In this example, assume that the designed system matrix A does not have full rank (there are fewer signals of interest than number of columns of matrix A), and the measured system matrix A has additive thermal noise that is larger than the quantization noise. The additive noise makes the measured matrix A have full rank.

Set $\sigma_{\text{noise}} = \sigma_{\text{thermal noise}}$.

```
noiseStandardDeviation = thermalNoiseStandardDeviation;
```

Use the `fixed.realQlessQRMatrixSolveFixedpointTypes` function to compute fixed-point types.

```
T = fixed.realQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,noiseStandardDeviation,[],regularizationParameter)
```

```
T =
```

```
struct with fields:
```

```
A: [0x0 embedded.fi]
```



```
B: [0x0 embedded.fi]
X: [0x0 embedded.fi]
```

T.A is the type computed for transforming $\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}$ to $R = Q^T \begin{bmatrix} \lambda I_n \\ A \end{bmatrix}$ in-place so that it does not overflow.

T.A

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 39
    FractionLength: 32
```

T.B is the type computed for B so that it does not overflow.

T.B

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 35
    FractionLength: 32
```

T.X is the type computed for the solution $X = \left(\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}^T \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \right) \setminus B$ so that there is a low probability that it overflows.

T.X

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 50
    FractionLength: 32
```

Use the Specified Types to Solve the Matrix Equation

Create random matrices A and B such that rankA=rank(A). Add random measurement noise to A which will make it become full rank.

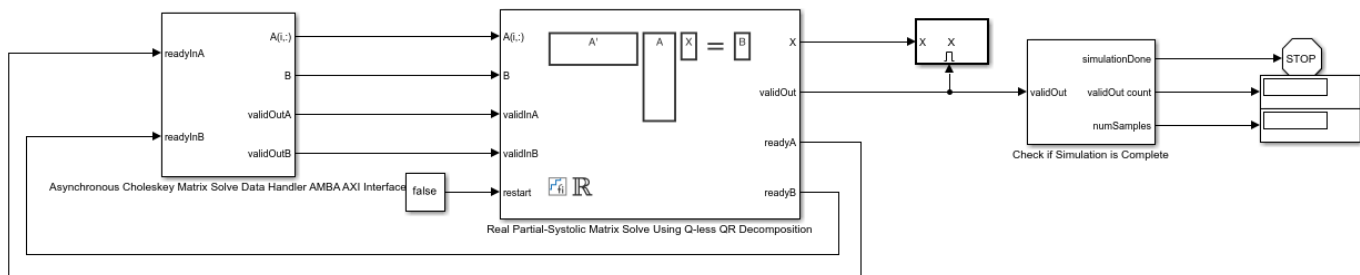
```
rng('default');
[A,B] = fixed.example.realRandomQlessQRMatrices(m,n,p,rankA);
A = A + fixed.example.realNormalRandomArray(0,noiseStandardDeviation,m,n);
```

Cast the inputs to the types determined by `fixed.realQlessQRMatrixSolveFixedpointTypes`. Quantizing to fixed-point is equivalent to adding random noise.

```
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);
```

Open the Model

```
model = 'RealPartialSystolicQlessQRMatrixSolveModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition block.

```
numOutputs = 1; % Number of recorded outputs
aDelay = 1; % Delay of clock cycles between feeding in rows of A
bDelay = 1; % Delay of clock cycles between feeding in B matrices
fixed.example.setModelWorkspace(model,'A',A,'B',B,'m',m,'n',n,'p',p,...
    'regularizationParameter',regularizationParameter,...
    'aDelay',aDelay,'bDelay',bDelay,...
    'numOutputs',numOutputs,'OutputType',OutputType);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition block outputs matrix `X` at each time step. When a valid result matrix is output, the block sets `validOut` to true.

```
X = out.X;
```

Verify the Accuracy of the Output

Verify that the relative error between the fixed-point output and builtin MATLAB in double-precision floating-point is small.

$$X_{\text{double}} = \left(\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}^T \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \right) \backslash B$$

```
A_lambda = double([regularizationParameter*eye(n);A]);
X_double = (A_lambda'*A_lambda)\double(B);
relativeError = norm(X_double - double(X))/norm(X_double)
```

```
relativeError =
```

```
4.8612e-06
```

Suppress mlint warnings in this file.

```
 %#ok<*NASGU>
 %#ok<*ASGLU>
 %#ok<*NOPTS>
```

See Also

Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition

Determine Fixed-Point Types for Complex Least-Squares Matrix Solve with Tikhonov Regularization

This example shows how to use the `fixed.complexQRMatrixSolveFixedpointTypes` function to analytically determine fixed-point types for the solution of the complex least-squares matrix equation

$$\begin{bmatrix} \lambda I_n \\ A \end{bmatrix} X = \begin{bmatrix} 0_{n,p} \\ B \end{bmatrix},$$

where A is an m -by- n matrix with $m \geq n$, B is m -by- p , X is n -by- p , $I_n = \text{eye}(n)$, $0_{n,p} = \text{zeros}(n, p)$, and λ is a regularization parameter.

The least-squares solution is

$$X_{LS} = (\lambda^2 I_n + A^T A)^{-1} A^T B$$

but is computed without squares or inverses.

Define System Parameters

Define the matrix attributes and system parameters for this example.

m is the number of rows in matrices A and B . In a problem such as beamforming or direction finding, m corresponds to the number of samples that are integrated over.

```
m = 300;
```

n is the number of columns in matrix A and rows in matrix X . In a least-squares problem, m is greater than n , and usually m is much larger than n . In a problem such as beamforming or direction finding, n corresponds to the number of sensors.

```
n = 10;
```

p is the number of columns in matrices B and X . It corresponds to simultaneously solving a system with p right-hand sides.

```
p = 1;
```

In this example, set the rank of matrix A to be less than the number of columns. In a problem such as beamforming or direction finding, $\text{rank}(A)$ corresponds to the number of signals impinging on the sensor array.

```
rankA = 3;
```

`precisionBits` defines the number of bits of precision required for the matrix solve. Set this value according to system requirements.

```
precisionBits = 32;
```

Small, positive values of the regularization parameter can improve the conditioning of the problem and reduce the variance of the estimates. While biased, the reduced variance of the estimate often results in a smaller mean squared error when compared to least-squares estimates.

```
regularizationParameter = 0.01;
```

In this example, complex-valued matrices A and B are constructed such that the magnitude of the real and imaginary parts of their elements is less than or equal to one, so the maximum possible absolute value of any element is $|1 + 1i| = \sqrt{2}$. Your own system requirements will define what those values are. If you don't know what they are, and A and B are fixed-point inputs to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of A and B.

`max_abs_A` is an upper bound on the maximum magnitude element of A.

```
max_abs_A = sqrt(2);
```

`max_abs_B` is an upper bound on the maximum magnitude element of B.

```
max_abs_B = sqrt(2);
```

Thermal noise standard deviation is the square root of thermal noise power, which is a system parameter. A well-designed system has the quantization level lower than the thermal noise. Here, set `thermalNoiseStandardDeviation` to the equivalent of -50dB noise power.

```
thermalNoiseStandardDeviation = sqrt(10^(-50/10))
```

```
thermalNoiseStandardDeviation = 0.0032
```

The quantization noise standard deviation is a function of the required number of bits of precision. Use `fixed.complexQuantizationNoiseStandardDeviation` to compute this. See that it is less than `thermalNoiseStandardDeviation`.

```
quantizationNoiseStandardDeviation = fixed.complexQuantizationNoiseStandardDeviation(precisionBits)
```

```
quantizationNoiseStandardDeviation = 9.5053e-11
```

Compute Fixed-Point Types

In this example, assume that the designed system matrix A does not have full rank (there are fewer signals of interest than number of columns of matrix A), and the measured system matrix A has additive thermal noise that is larger than the quantization noise. The additive noise makes the measured matrix A have full rank.

Set $\sigma_{\text{noise}} = \sigma_{\text{thermal noise}}$.

```
noiseStandardDeviation = thermalNoiseStandardDeviation;
```

Use `fixed.complexQRMatrixSolveFixedpointTypes` to compute fixed-point types.

```
T = fixed.complexQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,noiseStandardDeviation,[],regularizationParameter)
```

```
T = struct with fields:
    A: [0x0 embedded.fi]
    B: [0x0 embedded.fi]
    X: [0x0 embedded.fi]
```

`T.A` is the type computed for transforming $\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}$ to $R = Q^T \begin{bmatrix} \lambda I_n \\ A \end{bmatrix}$ in-place so that it does not overflow.

`T.A`

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 40
    FractionLength: 32
```

T.B is the type computed for transforming $\begin{bmatrix} 0_{n,p} \\ B \end{bmatrix}$ to $C = Q^T \begin{bmatrix} 0_{n,p} \\ B \end{bmatrix}$ in-place so that it does not overflow.

```
T.B
```

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 40
    FractionLength: 32
```

T.X is the type computed for the solution $X = \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \backslash \begin{bmatrix} 0_{n,p} \\ B \end{bmatrix}$, so that there is a low probability that it overflows.

```
T.X
```

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 44
    FractionLength: 32
```

Use the Specified Types to Solve the Matrix Equation

Create random matrices A and B such that B is in the range of A, and $\text{rankA} = \text{rank}(A)$. Add random measurement noise to A which will make it become full rank, but it will also affect the solution so that B is only close to the range of A.

```
rng('default');
[A,B] = fixed.example.complexRandomLeastSquaresMatrices(m,n,p,rankA);
A = A + fixed.example.complexNormalRandomArray(0,noiseStandardDeviation,m,n);
```

Cast the inputs to the types determined by `fixed.complexQRMatrixSolveFixedpointTypes`. Quantizing to fixed-point is equivalent to adding random noise [4,5].

```
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
```

Accelerate the `fixed.qrMatrixSolve` function by using `fiaccel` to generate a MATLAB executable (MEX) function.

```
fiaccel fixed.qrMatrixSolve -args {A,B,T.X,regularizationParameter} -o qrMatrixSolve_mex
```

Specify output type `T.X` and compute fixed-point $X = A \setminus B$ using the QR method.

```
X = qrMatrixSolve_mex(A,B,T.X,regularizationParameter);
```

Verify the Accuracy of the Output

Verify that the relative error between the fixed-point output and the output from MATLAB using the default double-precision floating-point values is small.

$$X_{\text{double}} = \begin{bmatrix} \lambda I_n & 0_{n,p} \\ A & B \end{bmatrix}$$

```
A_lambda = double([regularizationParameter*eye(n);A]);
B_0 = [zeros(n,p);double(B)];
X_double = A_lambda \ B_0;
relativeError = norm(X_double - double(X))/norm(X_double)
```

```
relativeError = 5.2634e-06
```

Suppress `mlint` warnings in this file.

```
 %#ok<*NASGU>
 %#ok<*ASGLU>
```

See Also

[fixed.complexQRMatrixSolveFixedpointTypes](#) | Complex Burst Matrix Solve Using QR Decomposition | Complex Partial-Systolic Matrix Solve Using QR Decomposition

Determine Fixed-Point Types for Complex Q-less QR Matrix Solve with Tikhonov Regularization

This example shows how to use the `fixed.complexQlessQRMatrixSolveFixedpointTypes` function to analytically determine fixed-point types for the solution of the complex least-squares matrix equation

$$\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}^H \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} X = (\lambda^2 I_n + A^H A) X = B$$

where A is an m -by- n matrix with $m \geq n$, B is n -by- p , X is n -by- p , $I_n = \text{eye}(n)$, and λ is a regularization parameter.

Define System Parameters

Define the matrix attributes and system parameters for this example.

m is the number of rows in matrix A . In a problem such as beamforming or direction finding, m corresponds to the number of samples that are integrated over.

```
m = 300;
```

n is the number of columns in matrix A and rows in matrices B and X . In a least-squares problem, m is greater than n , and usually m is much larger than n . In a problem such as beamforming or direction finding, n corresponds to the number of sensors.

```
n = 10;
```

p is the number of columns in matrices B and X . It corresponds to simultaneously solving a system with p right-hand sides.

```
p = 1;
```

In this example, set the rank of matrix A to be less than the number of columns. In a problem such as beamforming or direction finding, $\text{rank}(A)$ corresponds to the number of signals impinging on the sensor array.

```
rankA = 3;
```

`precisionBits` defines the number of bits of precision required for the matrix solve. Set this value according to system requirements.

```
precisionBits = 32;
```

Small, positive values of the regularization parameter can improve the conditioning of the problem and reduce the variance of the estimates. While biased, the reduced variance of the estimate often results in a smaller mean squared error when compared to least-squares estimates.

```
regularizationParameter = 0.01;
```

In this example, complex-valued matrices A and B are constructed such that the magnitude of the real and imaginary parts of their elements is less than or equal to one, so the maximum possible absolute value of any element is $|1 + 1i| = \sqrt{2}$. Your own system requirements will define what those values

are. If you don't know what they are, and A and B are fixed-point inputs to the system, then you can use the upperbound function to determine the upper bounds of the fixed-point types of A and B.

max_abs_A is an upper bound on the maximum magnitude element of A.

```
max_abs_A = sqrt(2);
```

max_abs_B is an upper bound on the maximum magnitude element of B.

```
max_abs_B = sqrt(2);
```

Thermal noise standard deviation is the square root of thermal noise power, which is a system parameter. A well-designed system has the quantization level lower than the thermal noise. Here, set thermalNoiseStandardDeviation to the equivalent of -50dB noise power.

```
thermalNoiseStandardDeviation = sqrt(10^(-50/10))
```

```
thermalNoiseStandardDeviation = 0.0032
```

The quantization noise standard deviation is a function of the required number of bits of precision. Use fixed.complexQuantizationNoiseStandardDeviation to compute this. See that it is less than thermalNoiseStandardDeviation.

```
quantizationNoiseStandardDeviation = fixed.complexQuantizationNoiseStandardDeviation(precisionBits)
```

```
quantizationNoiseStandardDeviation = 9.5053e-11
```

Compute Fixed-Point Types

In this example, assume that the designed system matrix A does not have full rank (there are fewer signals of interest than number of columns of matrix A), and the measured system matrix A has additive thermal noise that is larger than the quantization noise. The additive noise makes the measured matrix A have full rank.

Set $\sigma_{\text{noise}} = \sigma_{\text{thermal noise}}$.

```
noiseStandardDeviation = thermalNoiseStandardDeviation;
```

Use the fixed.complexQlessQRMatrixSolveFixedpointTypes function to compute fixed-point types.

```
T = fixed.complexQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,noiseStandardDeviation,[],regularizationParameter)
```

```
T = struct with fields:
    A: [0x0 embedded.fi]
    B: [0x0 embedded.fi]
    X: [0x0 embedded.fi]
```

T.A is the type computed for transforming $\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}$ to $R = Q^H \begin{bmatrix} \lambda I_n \\ A \end{bmatrix}$ in-place so that it does not overflow.

```
T.A
```

```
ans =
```

```
[]
```

```

        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 40
        FractionLength: 32

```

T.B is the type computed for B so that it does not overflow.

T.B

ans =

[]

```

        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 35
        FractionLength: 32

```

T.X is the type computed for the solution $X = \left(\begin{bmatrix} \lambda I_n & \\ & A \end{bmatrix}^H \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \right) \setminus B$ so that there is a low probability that it overflows.

T.X

ans =

[]

```

        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 48
        FractionLength: 32

```

Use the Specified Types to Solve the Matrix Equation

Create random matrices A and B such that $\text{rank}A = \text{rank}(A)$. Add random measurement noise to A which will make it become full rank.

```

rng('default');
[A,B] = fixed.example.complexRandomQlessQRMatrices(m,n,p,rankA);
A = A + fixed.example.complexNormalRandomArray(0,noiseStandardDeviation,m,n);

```

Cast the inputs to the types determined by `fixed.complexQlessQRMatrixSolveFixedpointTypes`. Quantizing to fixed-point is equivalent to adding random noise.

```

A = cast(A,'like',T.A);
B = cast(B,'like',T.B);

```

Accelerate the `fixed.qlessQRMatrixSolve` function by using `fiaccl` to generate a MATLAB executable (MEX) function.

```

fiaccl +fixed/qlessQRMatrixSolve -args {A,B,T.X,[],regularizationParameter} -o qlessQRMatrixSolve

```

Specify output type T.X and compute fixed-point $X = \left(\begin{bmatrix} \lambda I_n & \\ & A \end{bmatrix}^H \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \right) \setminus B$ using the QR method.

```
X = qlessQRMatrixSolve_mex(A,B,T,X,[],regularizationParameter);
```

Verify the Accuracy of the Output

Verify that the relative error between the fixed-point output and builtin MATLAB in double-precision floating-point is small.

$$X_{\text{double}} = \left(\begin{bmatrix} \lambda I_n & \\ & A \end{bmatrix}^H \begin{bmatrix} \lambda I_n \\ & A \end{bmatrix} \right) \setminus B$$

```
A_lambda = double([regularizationParameter*eye(n);A]);
X_double = (A_lambda'*A_lambda)\double(B);
relativeError = norm(X_double - double(X))/norm(X_double)

relativeError = 1.0591e-05
```

Suppress mlint warnings in this file.

```
##ok<*NASGU>
##ok<*ASGLU>
```

See Also

[fixed.complexQlessQRMatrixSolveFixedpointTypes](#) | Complex Burst Matrix Solve Using Q-less QR Decomposition | [Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition](#)

Determine Fixed-Point Types for Real Least-Squares Matrix Solve with Tikhonov Regularization

This example shows how to use the `fixed.realQRMatrixSolveFixedpointTypes` function to analytically determine fixed-point types for the solution of the real least-squares matrix equation

$$\begin{bmatrix} \lambda I_n \\ A \end{bmatrix} X = \begin{bmatrix} 0_{n,p} \\ B \end{bmatrix},$$

where A is an m -by- n matrix with $m \geq n$, B is m -by- p , X is n -by- p , $I_n = \text{eye}(n)$, $0_{n,p} = \text{zeros}(n, p)$, and λ is a regularization parameter.

The least-squares solution is

$$X_{LS} = (\lambda^2 I_n + A^T A)^{-1} A^T B$$

but is computed without squares or inverses.

Define System Parameters

Define the matrix attributes and system parameters for this example.

m is the number of rows in matrices A and B . In a problem such as beamforming or direction finding, m corresponds to the number of samples that are integrated over.

```
m = 300;
```

n is the number of columns in matrix A and rows in matrix X . In a least-squares problem, m is greater than n , and usually m is much larger than n . In a problem such as beamforming or direction finding, n corresponds to the number of sensors.

```
n = 10;
```

p is the number of columns in matrices B and X . It corresponds to simultaneously solving a system with p right-hand sides.

```
p = 1;
```

In this example, set the rank of matrix A to be less than the number of columns. In a problem such as beamforming or direction finding, $\text{rank}(A)$ corresponds to the number of signals impinging on the sensor array.

```
rankA = 3;
```

`precisionBits` defines the number of bits of precision required for the matrix solve. Set this value according to system requirements.

```
precisionBits = 32;
```

Small, positive values of the regularization parameter can improve the conditioning of the problem and reduce the variance of the estimates. While biased, the reduced variance of the estimate often results in a smaller mean squared error when compared to least-squares estimates.

```
regularizationParameter = 0.01;
```

In this example, real-valued matrices A and B are constructed such that the magnitude of their elements is less than or equal to one. Your own system requirements will define what those values are. If you don't know what they are, and A and B are fixed-point inputs to the system, then you can use the upperbound function to determine the upper bounds of the fixed-point types of A and B .

`max_abs_A` is an upper bound on the maximum magnitude element of A .

```
max_abs_A = 1;
```

`max_abs_B` is an upper bound on the maximum magnitude element of B .

```
max_abs_B = 1;
```

Thermal noise standard deviation is the square root of thermal noise power, which is a system parameter. A well-designed system has the quantization level lower than the thermal noise. Here, set `thermalNoiseStandardDeviation` to the equivalent of -50dB noise power.

```
thermalNoiseStandardDeviation = sqrt(10^(-50/10))
```

```
thermalNoiseStandardDeviation = 0.0032
```

The quantization noise standard deviation is a function of the required number of bits of precision. Use `fixed.realQuantizationNoiseStandardDeviation` to compute this. See that it is less than `thermalNoiseStandardDeviation`.

```
quantizationNoiseStandardDeviation = fixed.realQuantizationNoiseStandardDeviation(precisionBits)
```

```
quantizationNoiseStandardDeviation = 6.7212e-11
```

Compute Fixed-Point Types

In this example, assume that the designed system matrix A does not have full rank (there are fewer signals of interest than number of columns of matrix A), and the measured system matrix A has additive thermal noise that is larger than the quantization noise. The additive noise makes the measured matrix A have full rank.

Set $\sigma_{\text{noise}} = \sigma_{\text{thermal noise}}$.

```
noiseStandardDeviation = thermalNoiseStandardDeviation;
```

Use `fixed.realQRMatrixSolveFixedpointTypes` to compute fixed-point types.

```
T = fixed.realQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,noiseStandardDeviation,[],regularizationParameter)
```

```
T = struct with fields:
    A: [0x0 embedded.fi]
    B: [0x0 embedded.fi]
    X: [0x0 embedded.fi]
```

$T.A$ is the type computed for transforming $\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}$ to $R = Q^T \begin{bmatrix} \lambda I_n \\ A \end{bmatrix}$ in-place so that it does not overflow.

```
T.A
```

```
ans =
```

```

[]

```

```

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 39
    FractionLength: 32

```

T.B is the type computed for transforming $\begin{bmatrix} 0_{n,p} \\ B \end{bmatrix}$ to $C = Q^T \begin{bmatrix} 0_{n,p} \\ B \end{bmatrix}$ in-place so that it does not overflow.

```

T.B

```

```

ans =

```

```

[]

```

```

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 39
    FractionLength: 32

```

T.X is the type computed for the solution $X = \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \begin{bmatrix} 0_{n,p} \\ B \end{bmatrix}$, so that there is a low probability that it overflows.

```

T.X

```

```

ans =

```

```

[]

```

```

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 44
    FractionLength: 32

```

Use the Specified Types to Solve the Matrix Equation

Create random matrices A and B such that B is in the range of A, and $\text{rankA} = \text{rank}(A)$. Add random measurement noise to A which will make it become full rank, but it will also affect the solution so that B is only close to the range of A.

```

rng('default');
[A,B] = fixed.example.realRandomLeastSquaresMatrices(m,n,p,rankA);
A = A + fixed.example.realNormalRandomArray(0,noiseStandardDeviation,m,n);

```

Cast the inputs to the types determined by `fixed.realQRMatrixSolveFixedpointTypes`. Quantizing to fixed-point is equivalent to adding random noise [4,5].

```

A = cast(A,'like',T.A);
B = cast(B,'like',T.B);

```

Accelerate the `fixed.qrMatrixSolve` function by using `fiaccel` to generate a MATLAB executable (MEX) function.

```

fiaccel fixed.qrMatrixSolve -args {A,B,T.X,regularizationParameter} -o qrRealMatrixSolve_mex

```

Specify output type T.X and compute fixed-point $X = A \setminus B$ using the QR method.

```
X = qrRealMatrixSolve_mex(A,B,T.X,regularizationParameter);
```

Verify the Accuracy of the Output

Verify that the relative error between the fixed-point output and the output from MATLAB using the default double-precision floating-point values is small.

$$X_{\text{double}} = \begin{bmatrix} \lambda I_n & \begin{matrix} 0_{n,p} \\ B \end{matrix} \\ A \end{bmatrix}$$

```
A_lambda = double([regularizationParameter*eye(n);A]);
B_0 = [zeros(n,p);double(B)];
X_double = A_lambda\B_0;
relativeError = norm(X_double - double(X))/norm(X_double)
```

```
relativeError = 5.1152e-06
```

Suppress mlint warnings in this file.

```
 %#ok<NASGU>
 %#ok<ASGLU>
```

See Also

[fixed.realQRMatrixSolveFixedpointTypes](#) | Real Burst Matrix Solve Using QR Decomposition
| Real Partial-Systolic Matrix Solve Using QR Decomposition

Determine Fixed-Point Types for Real Q-less QR Matrix Solve with Tikhonov Regularization

This example shows how to use the `fixed.realQlessQRMatrixSolveFixedpointTypes` function to analytically determine fixed-point types for the solution of the real least-squares matrix equation

$$\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}^T \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} X = (\lambda^2 I_n + A^T A) X = B$$

where A is an m -by- n matrix with $m \geq n$, B is n -by- p , X is n -by- p , $I_n = \text{eye}(n)$, and λ is a regularization parameter.

Define System Parameters

Define the matrix attributes and system parameters for this example.

m is the number of rows in matrix A . In a problem such as beamforming or direction finding, m corresponds to the number of samples that are integrated over.

```
m = 300;
```

n is the number of columns in matrix A and rows in matrices B and X . In a least-squares problem, m is greater than n , and usually m is much larger than n . In a problem such as beamforming or direction finding, n corresponds to the number of sensors.

```
n = 10;
```

p is the number of columns in matrices B and X . It corresponds to simultaneously solving a system with p right-hand sides.

```
p = 1;
```

In this example, set the rank of matrix A to be less than the number of columns. In a problem such as beamforming or direction finding, $\text{rank}(A)$ corresponds to the number of signals impinging on the sensor array.

```
rankA = 3;
```

`precisionBits` defines the number of bits of precision required for the matrix solve. Set this value according to system requirements.

```
precisionBits = 32;
```

Small, positive values of the regularization parameter can improve the conditioning of the problem and reduce the variance of the estimates. While biased, the reduced variance of the estimate often results in a smaller mean squared error when compared to least-squares estimates.

```
regularizationParameter = 0.01;
```

In this example, real-valued matrices A and B are constructed such that the magnitude of their elements is less than or equal to one. Your own system requirements will define what those values are. If you don't know what they are, and A and B are fixed-point inputs to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of A and B .

`max_abs_A` is an upper bound on the maximum magnitude element of `A`.

```
max_abs_A = 1;
```

`max_abs_B` is an upper bound on the maximum magnitude element of `B`.

```
max_abs_B = 1;
```

Thermal noise standard deviation is the square root of thermal noise power, which is a system parameter. A well-designed system has the quantization level lower than the thermal noise. Here, set `thermalNoiseStandardDeviation` to the equivalent of -50dB noise power.

```
thermalNoiseStandardDeviation = sqrt(10^(-50/10))
```

```
thermalNoiseStandardDeviation = 0.0032
```

The quantization noise standard deviation is a function of the required number of bits of precision. Use `fixed.realQuantizationNoiseStandardDeviation` to compute this. See that it is less than `thermalNoiseStandardDeviation`.

```
quantizationNoiseStandardDeviation = fixed.realQuantizationNoiseStandardDeviation(precisionBits)
```

```
quantizationNoiseStandardDeviation = 6.7212e-11
```

Compute Fixed-Point Types

In this example, assume that the designed system matrix `A` does not have full rank (there are fewer signals of interest than number of columns of matrix `A`), and the measured system matrix `A` has additive thermal noise that is larger than the quantization noise. The additive noise makes the measured matrix `A` have full rank.

Set $\sigma_{\text{noise}} = \sigma_{\text{thermal noise}}$.

```
noiseStandardDeviation = thermalNoiseStandardDeviation;
```

Use the `fixed.realQlessQRMatrixSolveFixedpointTypes` function to compute fixed-point types.

```
T = fixed.realQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,...
    precisionBits,noiseStandardDeviation,[],regularizationParameter)
```

```
T = struct with fields:
    A: [0x0 embedded.fi]
    B: [0x0 embedded.fi]
    X: [0x0 embedded.fi]
```

`T.A` is the type computed for transforming $\begin{bmatrix} \lambda I_n \\ A \end{bmatrix}$ to $R = Q^T \begin{bmatrix} \lambda I_n \\ A \end{bmatrix}$ in-place so that it does not overflow.

```
T.A
```

```
ans =
```

```
[]
```

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
```

```

    WordLength: 39
    FractionLength: 32

```

T.B is the type computed for B so that it does not overflow.

```
T.B
```

```
ans =
```

```
[]
```

```

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 35
    FractionLength: 32

```

T.X is the type computed for the solution $X = \left(\begin{bmatrix} \lambda I_n & \\ & A \end{bmatrix}^T \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \right) \setminus B$ so that there is a low probability that it overflows.

```
T.X
```

```
ans =
```

```
[]
```

```

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 48
    FractionLength: 32

```

Use the Specified Types to Solve the Matrix Equation

Create random matrices A and B such that $\text{rank}A = \text{rank}(A)$. Add random measurement noise to A which will make it become full rank.

```

rng('default');
[A,B] = fixed.example.realRandomQlessQRMatrices(m,n,p,rankA);
A = A + fixed.example.realNormalRandomArray(0,noiseStandardDeviation,m,n);

```

Cast the inputs to the types determined by `fixed.realQlessQRMatrixSolveFixedpointTypes`. Quantizing to fixed-point is equivalent to adding random noise.

```

A = cast(A,'like',T.A);
B = cast(B,'like',T.B);

```

Accelerate the `fixed.qlessQRMatrixSolve` function by using `fiaccel` to generate a MATLAB executable (MEX) function.

```
fiaccel +fixed/qlessQRMatrixSolve -args {A,B,T.X,[],regularizationParameter} -o qlessQRMatrixSolve_mex
```

Specify output type T.X and compute fixed-point $X = \left(\begin{bmatrix} \lambda I_n & \\ & A \end{bmatrix}^T \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \right) \setminus B$ using the QR method.

```
X = qlessQRMatrixSolve_mex(A,B,T.X,[],regularizationParameter);
```

Verify the Accuracy of the Output

Verify that the relative error between the fixed-point output and builtin MATLAB in double-precision floating-point is small.

$$X_{\text{double}} = \left(\begin{bmatrix} \lambda I_n & \\ & \lambda I_n \end{bmatrix}^T \begin{bmatrix} \lambda I_n \\ A \end{bmatrix} \right) \setminus B$$

```
A_lambda = double([regularizationParameter*eye(n);A]);
X_double = (A_lambda'*A_lambda)\double(B);
relativeError = norm(X_double - double(X))/norm(X_double)
```

```
relativeError = 1.0133e-05
```

Suppress mlint warnings in this file.

```
 %#ok<*NASGU>
 %#ok<*ASGLU>
```

See Also

[fixed.realQlessQRMatrixSolveFixedpointTypes](#) | Real Burst Matrix Solve Using Q-less QR Decomposition | Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition

Implement Hardware-Efficient Real Burst Q-less QR with Forgetting Factor

This example shows how to use the hardware-efficient Real Burst Q-less QR Decomposition with Forgetting Factor Whole R Output block.

Q-less QR Decomposition with Forgetting Factor

The Real Burst Q-less QR Decomposition with Forgetting Factor Whole R Output block implements the following recursion to compute the upper-triangular factor R of continuously streaming n-by-1 row vectors $A(k,:)$ using forgetting factor α . It is as if matrix A is infinitely tall. The forgetting factor in the range $0 < \alpha < 1$ keeps it from integrating without bound.

$$\begin{aligned}
 R_0 &= \text{zeros}(n, n) \\
 [\sim, R_1] &= \text{qr} \left(\begin{bmatrix} R_0 \\ A(1,:) \end{bmatrix}, 0 \right) \\
 R_1 &= \alpha R_1 \\
 \\
 [\sim, R_2] &= \text{qr} \left(\begin{bmatrix} R_1 \\ A(2,:) \end{bmatrix}, 0 \right) \\
 R_2 &= \alpha R_2 \\
 \\
 &\vdots \\
 \\
 [\sim, R_k] &= \text{qr} \left(\begin{bmatrix} R_{k-1} \\ A(k,:) \end{bmatrix}, 0 \right) \\
 R_k &= \alpha R_k \\
 \\
 &\vdots
 \end{aligned}$$

AMBA AXI Handshaking Process

The Data Handler subsystem in this model takes real matrix A as input. It sends rows of A to the QR Decomposition block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set delay for the feeding in rows of A in the Data Handler to emulate the processing time of the upstream block. `validOut` signal of the Data Handler remain high when `delayLen` is set to 0 because this indicates the Data Handler always has data available.

Define System Parameters

n is the length of the row vectors $A(k,:)$ and the number of rows and columns in R.

n = 5;

m is the effective numbers of rows of A to integrate over.

```
m = 100;
```

Use the `fixed.forgettingFactor` function to compute the forgetting factor as a function of the number of rows that you are integrating over.

```
forgettingFactor = fixed.forgettingFactor(m)
```

```
forgettingFactor =
```

```
    0.9950
```

`precisionBits` defines the number of bits of precision required for the QR Decomposition. Set this value according to system requirements.

```
precisionBits = 24;
```

In this example, real-valued matrix `A` is constructed such that the magnitude of its elements are less than or equal to one. Your own system requirements will define what those values are. If you don't know what they are, and `A` is a fixed-point input to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of `A`.

`max_abs_A` is an upper bound on the maximum magnitude element of `A`.

```
max_abs_A = 1;
```

Select Fixed-Point Types

Use the `fixed.qlessqrFixedpointTypes` function to compute fixed-point types.

```
T = fixed.qlessqrFixedpointTypes(m,max_abs_A,precisionBits)
```

```
T =
```

```
  struct with fields:
```

```
    A: [0x0 embedded.fi]
```

`T.A` is the fixed-point type computed for transforming `A` to `R` in-place so that it does not overflow.

```
T.A
```

```
ans =
```

```
 []
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 31
    FractionLength: 24
```

Define Simulation Parameters

Create random matrix `A` to contain a specified number of inputs.

numInputs is the number of input rows A(k,:) for this example.

```
numInputs = 500;
rng('default')
A = fixed.example.realUniformRandomArray(-1,1,numInputs,n);
```

Cast the inputs to the types determined by fixed.qlessqrFixedpointTypes.

```
A = cast(A,'like',T.A);
```

Cast the forgetting factor to a fixed-point type with the same word length as A and best-precision scaling.

```
forgettingFactor = fi(forgettingFactor,1,T.A.WordLength);
```

Set delay for feeding in rows of A.

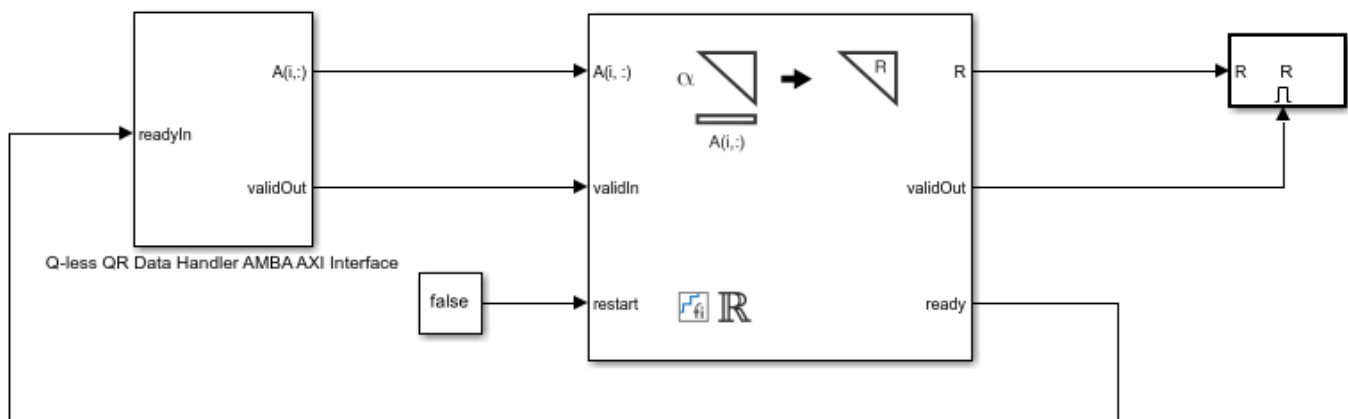
```
delayLen = 1;
```

Select a stop time for the simulation that is long enough to process all the inputs from A.

```
stopTime = 2*numInputs*T.A.WordLength;
```

Open the Model

```
model = 'RealBurstQlessQRForgettingFactorModel';
open_system(model);
```



Copyright 2022 The MathWorks, Inc.

Set Variables in the Model Workspace

Use the helper function setModelWorkspace to add the variables defined above to the model workspace.

```
fixed.example.setModelWorkspace(model,'A',A,'n',n,...
    'forgettingFactor',forgettingFactor,...
    'regularizationParameter',0,...
    'delayLen',delayLen,'stopTime',stopTime);
```

Simulate the Model

```
out = sim(model);
```

Verify the Accuracy of the Output

Define matrix A_k as follows

$$A_k = \begin{bmatrix} \alpha^k & & & & \\ & \alpha^{k-1} & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & \alpha \end{bmatrix} A(1:k, :).$$

Then using the formula for the computation of the k th output R_k , and the fact that $[Q, R] = \text{qr}(A, 0) \Rightarrow A'A = R'Q'QR = R'R$, you can show that

$$A_k'A_k = R_k'R_k.$$

So to verify the output, the difference between $A_k'A_k$ and $R_k'R_k$ should be small.

Choose the last output of the simulation.

```
R = double(out.R(:, :, end))
```

```
R =
```

```

5.2620    0.5072    0.1204    0.4957   -0.1677
         0    5.0102   -0.3150   -0.0484    0.6541
         0         0    5.2474   -0.2291    0.0964
         0         0         0    5.0981    0.0500
         0         0         0         0    5.0053
```

Verify that R is upper triangular.

```
isequal(R, triu(R))
```

```
ans =
```

```
logical
```

```
1
```

Verify that the diagonal is greater than or equal to zero.

```
diag(R)
```

```
ans =
```

```

5.2620
5.0102
```

```

5.2474
5.0981
5.0053

```

Synchronize the last output R with the input by finding the number of inputs that produced it.

```

A = double(A);
alpha = double(forgettingFactor);
relative_errors = nan(1,n);
for k = 1:numInputs
    A_k = alpha.^(k:-1:1)' .* A(1:k,:);
    relative_errors(k) = norm(A_k'*A_k - R'*R)/norm(A_k'*A_k);
end

```

k is the number of inputs A(k,:) that produced the last R.

```

k = find(relative_errors==min(relative_errors),1,'last')

```

```

k =
    165

```

Verify that

$$A_k' A_k = R_k' R_k$$

with a small relative error.

```

A_k = alpha.^(k:-1:1)' .* A(1:k,:);
relative_error = norm(A_k'*A_k - R'*R)/norm(A_k'*A_k)

```

```

relative_error =
    5.9156e-06

```

Suppress mlint warnings in this file.

```

%#ok< *NOPTS>

```

See Also

Blocks

Real Burst Q-less QR Decomposition with Forgetting Factor Whole R Output

Functions

fixed.forgettingFactor | upperbound | fixed.qlessqrFixedpointTypes

Implement Hardware-Efficient Complex Burst Q-less QR with Forgetting Factor

This example shows how to use the hardware-efficient Complex Burst Q-less QR Decomposition with Forgetting Factor Whole R Output block.

Q-less QR Decomposition with Forgetting Factor

The Complex Burst Q-less QR Decomposition with Forgetting Factor Whole R Output block implements the following recursion to compute the upper-triangular factor R of continuously streaming n-by-1 row vectors $A(k,:)$ using forgetting factor α . It is as if matrix A is infinitely tall. The forgetting factor in the range $0 < \alpha < 1$ keeps it from integrating without bound.

$$\begin{aligned}
 R_0 &= \text{zeros}(n, n) \\
 [\sim, R_1] &= \text{qr} \left(\begin{bmatrix} R_0 \\ A(1,:) \end{bmatrix}, 0 \right) \\
 R_1 &= \alpha R_1 \\
 \\
 [\sim, R_2] &= \text{qr} \left(\begin{bmatrix} R_1 \\ A(2,:) \end{bmatrix}, 0 \right) \\
 R_2 &= \alpha R_2 \\
 \\
 &\vdots \\
 \\
 [\sim, R_k] &= \text{qr} \left(\begin{bmatrix} R_{k-1} \\ A(k,:) \end{bmatrix}, 0 \right) \\
 R_k &= \alpha R_k \\
 \\
 &\vdots
 \end{aligned}$$

Define System Parameters

n is the length of the row vectors $A(k,:)$ and the number of rows and columns in R.

```
n = 5;
```

m is the effective numbers of rows of A to integrate over.

```
m = 100;
```

Use the `fixed.forgettingFactor` function to compute the forgetting factor as a function of the number of rows that you are integrating over.

```
forgettingFactor = fixed.forgettingFactor(m)
```

```
forgettingFactor =
```

```
0.9950
```

`precisionBits` defines the number of bits of precision required for the QR Decomposition. Set this value according to system requirements.

```
precisionBits = 24;
```

In this example, complex-valued matrix `A` is constructed such that the magnitude of the real and imaginary parts of its elements is less than or equal to one, so the maximum possible absolute value of any element is $|1 + 1i| = \sqrt{2}$. Your own system requirements will define what those values are. If you don't know what they are, and `A` is a fixed-point input to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of `A`.

`max_abs_A` is an upper bound on the maximum magnitude element of `A`.

```
max_abs_A = sqrt(2);
```

Select Fixed-Point Types

Use the `fixed.qlessqrFixedpointTypes` function to compute fixed-point types.

```
T = fixed.qlessqrFixedpointTypes(m,max_abs_A,precisionBits)
```

```
T =
```

```
struct with fields:
```

```
  A: [0x0 embedded.fi]
```

`T.A` is the fixed-point type computed for transforming `A` to `R` in-place so that it does not overflow.

```
T.A
```

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 31
    FractionLength: 24
```

AMBA AXI Handshaking Process

The Data Handler subsystem in this model takes complex matrix `A` as inputs. It sends rows of `A` to QR Decomposition block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set delay for the feeding in rows of `A` in the Data Handler to emulate the processing time of the upstream block. `validOut` signal of Data Handler remain high when `delayLen` is set to 0 because this indicates the Data Handler always has data available.

Define Simulation Parameters

Create random matrix A to contain a specified number of inputs.

numInputs is the number of input rows A(k,:) for this example.

```
numInputs = 500;
rng('default')
A = fixed.example.complexUniformRandomArray(-1,1,numInputs,n);
```

Cast the inputs to the types determined by fixed.qlessqrFixedpointTypes.

```
A = cast(A,'like',T.A);
```

Cast the forgetting factor to a fixed-point type with the same word length as A and best-precision scaling.

```
forgettingFactor = fi(forgettingFactor,1,T.A.WordLength);
```

Set delay for feeding in rows of A

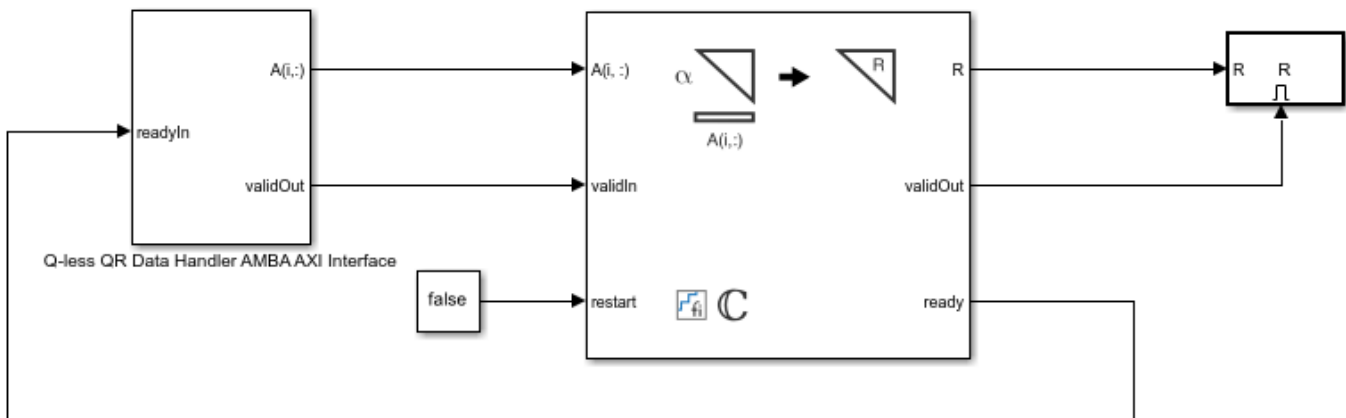
```
delayLen = 1;
```

Select a stop time for the simulation that is long enough to process all the inputs from A.

```
stopTime = 4*numInputs*T.A.WordLength;
```

Open the Model

```
model = 'ComplexBurstQlessQRForgettingFactorModel';
open_system(model);
```



Copyright 2022 The MathWorks, Inc.

Set Variables in the Model Workspace

Use the helper function setModelWorkspace to add the variables defined above to the model workspace.

```
fixed.example.setModelWorkspace(model,'A',A,'n',n,...
    'forgettingFactor',forgettingFactor,...
```

```
'regularizationParameter',0,...
'delayLen',delayLen,...
'stopTime',stopTime);
```

Simulate the Model

```
out = sim(model);
```

Verify the Accuracy of the Output

Define matrix A_k as follows

$$A_k = \begin{bmatrix} \alpha^k & & & \\ & \alpha^{k-1} & & \\ & & \ddots & \\ & & & \alpha \end{bmatrix} A(1:k, :).$$

Then using the formula for the computation of the k th output R_k , and the fact that $[Q, R] = \text{qr}(A, 0) \Rightarrow A'A = R'Q'QR = R'R$, you can show that

$$A_k'A_k = R_k'R_k.$$

So to verify the output, the difference between $A_k'A_k$ and $R_k'R_k$ should be small.

Choose the last output of the simulation.

```
R = double(out.R(:, :, end))
```

```
R =
```

```
Columns 1 through 4
```

```
7.4030 + 0.0000i    0.2517 - 0.3472i    0.4163 - 0.1448i    0.4088 + 0.5546i
0.0000 + 0.0000i    7.3291 + 0.0000i   -0.1239 - 0.3553i   -0.8237 + 0.2091i
0.0000 + 0.0000i    0.0000 + 0.0000i    7.3507 + 0.0000i    0.2622 - 0.6994i
0.0000 + 0.0000i    0.0000 + 0.0000i    0.0000 + 0.0000i    7.2422 + 0.0000i
0.0000 + 0.0000i    0.0000 + 0.0000i    0.0000 + 0.0000i    0.0000 + 0.0000i
```

```
Column 5
```

```
0.5450 + 0.7208i
0.1945 + 0.1716i
-0.5293 + 0.4192i
0.4574 - 0.1519i
7.0558 + 0.0000i
```

Verify that R is upper triangular.

```
isequal(R, triu(R))
```

```
ans =
```

```
logical
```

```
1
```

Verify that the diagonal is greater than or equal to zero.

```
diag(R)
```

```
ans =
```

```
7.4030
7.3291
7.3507
7.2422
7.0558
```

Synchronize the last output R with the input by finding the number of inputs that produced it.

```
A = double(A);
alpha = double(forgettingFactor);
relative_errors = nan(1,n);
for k = 1:numInputs
    A_k = alpha.^(k:-1:1)' .* A(1:k,:);
    relative_errors(k) = norm(A_k'*A_k - R'*R)/norm(A_k'*A_k);
end
```

k is the number of inputs A(k,:) that produced the last R.

```
k = find(relative_errors==min(relative_errors),1,'last')
```

```
k =
```

```
166
```

Verify that

$$A_k' A_k = R_k' R_k$$

with a small relative error.

```
A_k = alpha.^(k:-1:1)' .* A(1:k,:);
relative_error = norm(A_k'*A_k - R'*R)/norm(A_k'*A_k)
```

```
relative_error =
```

```
5.2882e-06
```

Suppress mlint warnings in this file.

`%#ok<*NOPTS>`

See Also

Blocks

Complex Burst Q-less QR Decomposition with Forgetting Factor Whole R Output

Functions

`fixed.forgettingFactor` | `upperbound` | `fixed.qlessqrFixedpointTypes`

Implement Hardware-Efficient Real Burst Matrix Solve Using Q-less QR Decomposition with Forgetting Factor

This example shows how to use the hardware-efficient Real Burst Matrix Solve Using Q-less QR Decomposition with Forgetting Factor block.

Q-less QR Decomposition with Forgetting Factor

The Real Burst Matrix Solve Using Q-less QR Decomposition with Forgetting Factor block implements the following recursion to compute the upper-triangular factor R of continuously streaming n-by-1 row vectors A(k,:) using forgetting factor α . It's as if matrix A is infinitely tall. The forgetting factor in the range $0 < \alpha < 1$ keeps it from integrating without bound.

$$\begin{aligned}
 R_0 &= \text{zeros}(n, n) \\
 [\sim, R_1] &= \text{qr} \left(\begin{bmatrix} R_0 \\ A(1, :) \end{bmatrix}, 0 \right) \\
 R_1 &= \alpha R_1 \\
 \\
 [\sim, R_2] &= \text{qr} \left(\begin{bmatrix} R_1 \\ A(2, :) \end{bmatrix}, 0 \right) \\
 R_2 &= \alpha R_2 \\
 \\
 &\vdots \\
 \\
 [\sim, R_k] &= \text{qr} \left(\begin{bmatrix} R_{k-1} \\ A(k, :) \end{bmatrix}, 0 \right) \\
 R_k &= \alpha R_k \\
 \\
 &\vdots
 \end{aligned}$$

Forward and Backward Substitution

When an upper triangular factor is ready, then forward and backward substitution are computed with the current input B to produce output X.

$$X = R_k \setminus (R'_k \setminus B)$$

AMBA AXI Handshaking Process

The Data Handler subsystem in this model takes real matrices A and B as inputs. It sends rows of A and full matrix of B to QR Decomposition block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set delays for the feeding in rows of A and the feeding in B matrices in the Data Handler to emulate the processing time of the upstream block. `validInA` and `validInB` remain high when `aDelay` and `bDelay` are set to 0 because this indicates the Data Handler always has data available.

Asynchronous Matrix Solver

This block operates asynchronously. First, Q-less QR decomposition is performed on the input A matrix and the resulting R matrix is put into a buffer. Then, the Forward Backward Substitute block uses the input B matrix and the buffered R matrix to compute $R'RX = B$. Because the R and B matrices are stored separately in buffers, the upstream Q-less QR decomposition block and the downstream Forward Backward Substitute block can run independently. The Forward Backward Substitute block starts processing when the first R and B matrices are available. Then it runs continuously using the latest buffered R and B matrices, regardless of the status of the Q-less QR Decomposition block. For example, if the upstream block stops providing A and B matrices, the Forward Backward Substitute block continues to generate the same output using the last pair of R and B matrices.

Define System Parameters

n is the length of the row vectors $A(k,:)$, the number of rows in B, and the number of rows and columns in R.

```
n = 5;
```

p is the number of columns in B

```
p = 1;
```

m is the effective numbers of rows of A to integrate over.

```
m = 100;
```

Use the `fixed.forgettingFactor` function to compute the forgetting factor as a function of the number of rows that you are integrating over.

```
forgettingFactor = fixed.forgettingFactor(m)
```

```
forgettingFactor =
```

```
0.9950
```

`precisionBits` defines the number of bits of precision required for the QR Decomposition. Set this value according to system requirements.

```
precisionBits = 24;
```

In this example, real-valued matrices A and B are constructed such that the magnitude of their real and imaginary parts of their elements is less than or equal to one, so the maximum possible absolute value of any element is $|1 + 1i| = \sqrt{2}$. Your own system requirements will define what those values are. If you don't know what they are, and A and B are fixed-point input to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of A and B.

`max_abs_A` is an upper bound on the maximum magnitude element of A.

```
max_abs_A = sqrt(2);
```

`max_abs_B` is an upper bound on the maximum magnitude element of B.

```
max_abs_B = sqrt(2);
```


Select Fixed-Point Types

Use the `fixed.realQlessQRMatrixSolveFixedpointTypes` function to compute fixed-point types.

```
T = fixed.realQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,precisionBits);
```

T.A is the fixed-point type computed for transforming A to R in-place so that it does not overflow.

T.A

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 31
    FractionLength: 24
```

T.B is the type computed for B so that it does not overflow.

T.B

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 27
    FractionLength: 24
```

T.X is the type computed for the output X so that there is a low probability of overflow.

T.X

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 77
    FractionLength: 24
```

Define Simulation Parameters

Create random matrix A to contain a specified number of inputs, and n-by-p random matrix B.

`numInputs` is the number of input rows $A(k,:)$ for this example.

```
numInputs = 500;
rng('default')
[A,B] = fixed.example.realRandomQlessQRMatrices(numInputs,n,p);
```

Cast the inputs to the types determined by `fixed.complexQlessQRMatrixSolveFixedpointTypes`.

```
A = cast(A, 'like', T.A);
B = cast(B, 'like', T.B);
```

Use the `fixed.extractNumericType` function to extract a `numericType` object to use as an input parameter to the block.

```
OutputType = fixed.extractNumericType(T.X)
```

```
OutputType =
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 77
    FractionLength: 24
```

Cast the forgetting factor to a fixed-point type with the same word length as A and best-precision scaling.

```
forgettingFactor = fi(forgettingFactor, 1, T.A.WordLength);
```

Set delay for feeding in rows of A

```
aDelay = 1;
```

Set delay for feeding in B matrices

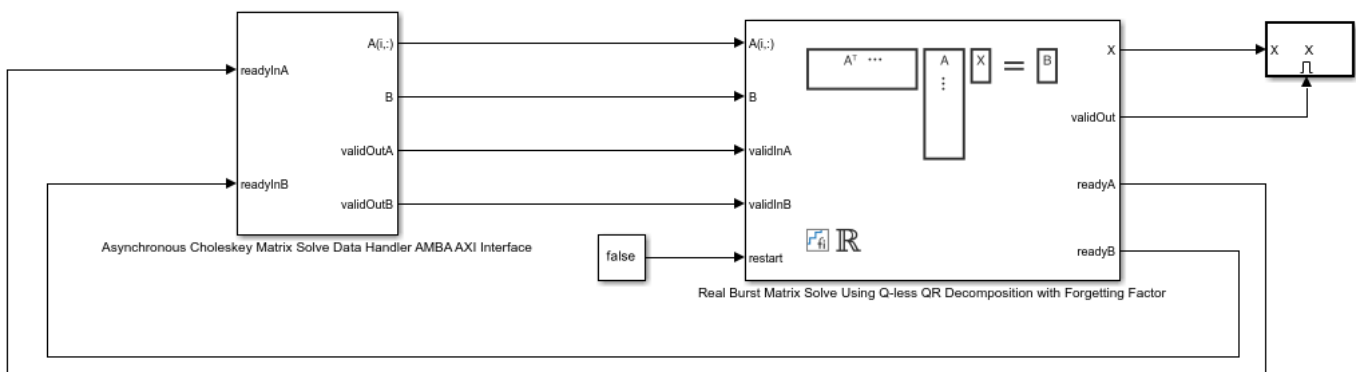
```
bDelay = 1;
```

Select a stop time for the simulation that is long enough to process all the inputs from A.

```
stopTime = 2*(2*numInputs + n)*T.A.WordLength;
```

Open the Model

```
model = 'RealBurstSolveQlessQRForgettingFactorModel';
open_system(model);
```



Copyright 2022 The MathWorks, Inc.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace.

```
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'n', n, 'p', p, ...
    'forgettingFactor', forgettingFactor, 'OutputType', OutputType, ...
    'regularizationParameter', 0, ...
    'aDelay', aDelay, 'bDelay', bDelay, ...
    'stopTime', stopTime);
```

Simulate the Model

```
out = sim(model);
```

Verify the Accuracy of the Output

Define matrix A_k as follows

$$A_k = \begin{bmatrix} \alpha^k & & & \\ & \alpha^{k-1} & & \\ & & \ddots & \\ & & & \alpha \end{bmatrix} A(1:k, :).$$

Then using the formula for the computation of the k th output R_k , and the fact that $[Q, R] = \text{qr}(A, 0) \Rightarrow A'A = R'Q'QR = R'R$, you can show that

$$A_k' A_k X = R_k' R_k X = B.$$

So to verify the output, the difference between $A_k' A_k X$ and B should be small.

Choose the last output of the simulation.

```
X = double(out.X(:, :, end));
```

Synchronize the last output X with the input by finding the number of inputs that produced it.

```
A = double(A);
B = double(B);
alpha = double(forgettingFactor);
relative_errors = nan(1, numInputs);
for k = 1:numInputs
    A_k = alpha.^(k:-1:1)' .* A(1:k, :);
    relative_errors(k) = norm(A_k'*A_k*X - B)/norm(B);
end
```

k is the number of inputs $A(k, :)$ that produced the last X .

```
k = find(relative_errors==min(relative_errors))
```

```
k =
```

```
329
```

Verify that

$$A_k' A_k X = B$$

with a small relative error.

```
A_k = alpha.^(k:-1:1)' .* A(1:k,:);  
relative_error = norm(A_k'*A_k*X - B)/norm(B)
```

```
relative_error =
```

```
7.6410e-04
```

Suppress mlint warnings in this file.

```
 %#ok<*NOPTS>
```

See Also

Blocks

Real Burst Matrix Solve Using Q-less QR Decomposition with Forgetting Factor

Functions

`fixed.forgettingFactor` | `upperbound` |
`fixed.realQlessQRMatrixSolveFixedpointTypes` | `fixed.extractNumericType`

Implement Hardware-Efficient Complex Burst Matrix Solve Using Q-less QR Decomposition with Forgetting Factor

This example shows how to use the hardware-efficient Complex Burst Matrix Solve Using Q-less QR Decomposition with Forgetting Factor block.

Q-less QR Decomposition with Forgetting Factor

The Complex Burst Matrix Solve Using Q-less QR Decomposition with Forgetting Factor block implements the following recursion to compute the upper-triangular factor R of continuously streaming n-by-1 row vectors A(k,:) using forgetting factor α . It's as if matrix A is infinitely tall. The forgetting factor in the range $0 < \alpha < 1$ keeps it from integrating without bound.

$$\begin{aligned}
 R_0 &= \text{zeros}(n, n) \\
 [\sim, R_1] &= \text{qr} \left(\begin{bmatrix} R_0 \\ A(1, :) \end{bmatrix}, 0 \right) \\
 R_1 &= \alpha R_1 \\
 \\
 [\sim, R_2] &= \text{qr} \left(\begin{bmatrix} R_1 \\ A(2, :) \end{bmatrix}, 0 \right) \\
 R_2 &= \alpha R_2 \\
 \\
 &\vdots \\
 \\
 [\sim, R_k] &= \text{qr} \left(\begin{bmatrix} R_{k-1} \\ A(k, :) \end{bmatrix}, 0 \right) \\
 R_k &= \alpha R_k \\
 \\
 &\vdots
 \end{aligned}$$

Forward and Backward Substitution

When an upper triangular factor is ready, then forward and backward substitution are computed with the current input B to produce output X.

$$X = R_k \setminus (R'_k \setminus B)$$

AMBA AXI Handshaking Process

The Data Handler subsystem in this model takes complex matrices A and B as inputs. It sends rows of A and full matrix of B to QR Decomposition block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set delays for the feeding in rows of A and the feeding in B matrices in the Data Handler to emulate the processing time of the upstream block. `validInA` and `validInB` remain high when `aDelay` and `bDelay` are set to 0 because this indicates the Data Handler always has data available.

Asynchronous Matrix Solver

This block operates asynchronously. First, Q-less QR decomposition is performed on the input A matrix and the resulting R matrix is put into a buffer. Then, the Forward Backward Substitute block uses the input B matrix and the buffered R matrix to compute $R'RX = B$. Because the R and B matrices are stored separately in buffers, the upstream Q-less QR decomposition block and the downstream Forward Backward Substitute block can run independently. The Forward Backward Substitute block starts processing when the first R and B matrices are available. Then it runs continuously using the latest buffered R and B matrices, regardless of the status of the Q-less QR Decomposition block. For example, if the upstream block stops providing A and B matrices, the Forward Backward Substitute block continues to generate the same output using the last pair of R and B matrices.

Define System Parameters

n is the length of the row vectors $A(k,:)$, the number of rows in B, and the number of rows and columns in R.

```
n = 5;
```

p is the number of columns in B

```
p = 1;
```

m is the effective numbers of rows of A to integrate over.

```
m = 100;
```

Use the `fixed.forgettingFactor` function to compute the forgetting factor as a function of the number of rows that you are integrating over.

```
forgettingFactor = fixed.forgettingFactor(m)
```

```
forgettingFactor =
```

```
0.9950
```

`precisionBits` defines the number of bits of precision required for the QR Decomposition. Set this value according to system requirements.

```
precisionBits = 24;
```

In this example, complex-valued matrices A and B are constructed such that the magnitude of their real and imaginary parts of their elements is less than or equal to one, so the maximum possible absolute value of any element is $|1 + 1i| = \sqrt{2}$. Your own system requirements will define what those values are. If you don't know what they are, and A and B are fixed-point input to the system, then you can use the `upperbound` function to determine the upper bounds of the fixed-point types of A and B.

`max_abs_A` is an upper bound on the maximum magnitude element of A.

```
max_abs_A = sqrt(2);
```

`max_abs_B` is an upper bound on the maximum magnitude element of B.

```
max_abs_B = sqrt(2);
```

Select Fixed-Point Types

Use the `fixed.complexQlessQRMatrixSolveFixedpointTypes` function to compute fixed-point types.

```
T = fixed.complexQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,precisionBits);
```

T.A is the fixed-point type computed for transforming A to R in-place so that it does not overflow.

T.A

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 31
    FractionLength: 24
```

T.B is the type computed for B so that it does not overflow.

T.B

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 27
    FractionLength: 24
```

T.X is the type computed for the output X so that there is a low probability of overflow.

T.X

```
ans =
```

```
[]
```

```
    DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 75
    FractionLength: 24
```

Define Simulation Parameters

Create random matrix A to contain a specified number of inputs, and n-by-p random matrix B.

`numInputs` is the number of input rows $A(k,:)$ for this example.

```
numInputs = 500;
rng('default')
[A,B] = fixed.example.complexRandomQlessQRMatrices(numInputs,n,p);
```

Cast the inputs to the types determined by `fixed.complexQlessQRMatrixSolveFixedpointTypes`.

```
A = cast(A, 'like', T.A);
B = cast(B, 'like', T.B);
```

Use the `fixed.extractNumericType` function to extract a `numericType` object to use as an input parameter to the block.

```
OutputType = fixed.extractNumericType(T.X)
```

```
OutputType =
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 75
    FractionLength: 24
```

Cast the forgetting factor to a fixed-point type with the same word length as A and best-precision scaling.

```
forgettingFactor = fi(forgettingFactor, 1, T.A.WordLength);
```

Set delay for feeding in rows of A

```
aDelay = 1;
```

Set delay for feeding in B matrices

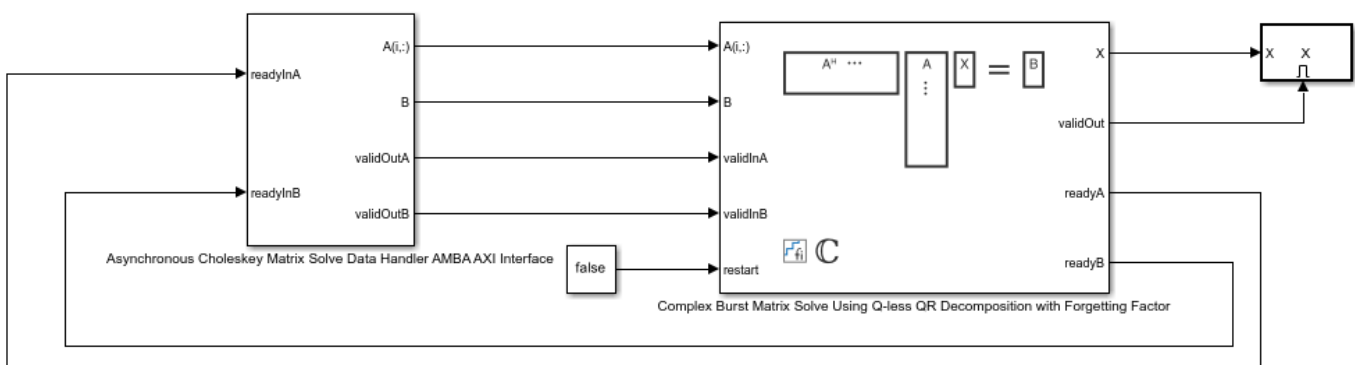
```
bDelay = 1;
```

Select a stop time for the simulation that is long enough to process all the inputs from A.

```
stopTime = 2*(2*numInputs + n)*T.A.WordLength;
```

Open the Model

```
model = 'ComplexBurstSolveQlessQRForgettingFactorModel';
open_system(model);
```



Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace.

```
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'n', n, 'p', p, ...
    'forgettingFactor', forgettingFactor, 'OutputType', OutputType, ...
    'regularizationParameter', 0, ...
    'aDelay', aDelay, 'bDelay', bDelay, ...
    'stopTime', stopTime);
```

Simulate the Model

```
out = sim(model);
```

Verify the Accuracy of the Output

Define matrix A_k as follows

$$A_k = \begin{bmatrix} \alpha^k & & & \\ & \alpha^{k-1} & & \\ & & \ddots & \\ & & & \alpha \end{bmatrix} A(1:k, :).$$

Then using the formula for the computation of the k th output R_k , and the fact that $[Q, R] = \text{qr}(A, 0) \Rightarrow A'A = R'Q'QR = R'R$, you can show that

$$A_k' A_k X = R_k' R_k X = B.$$

So to verify the output, the difference between $A_k' A_k X$ and B should be small.

Choose the last output of the simulation.

```
X = double(out.X(:, :, end));
```

Synchronize the last output X with the input by finding the number of inputs that produced it.

```
A = double(A);
B = double(B);
alpha = double(forgettingFactor);
relative_errors = nan(1, numInputs);
for k = 1:numInputs
    A_k = alpha.^(k:-1:1)' .* A(1:k, :);
    relative_errors(k) = norm(A_k'*A_k*X - B)/norm(B);
end
```

k is the number of inputs $A(k, :)$ that produced the last X .

```
k = find(relative_errors==min(relative_errors))
```

```
k =
```

```
165
```

Verify that

$$A_k' A_k X = B$$

with a small relative error.

```
A_k = alpha.^(k:-1:1)' .* A(1:k,:);  
relative_error = norm(A_k'*A_k*X - B)/norm(B)
```

```
relative_error =
```

```
3.2859e-05
```

Suppress mlint warnings in this file.

```
 %#ok<*NOPTS>
```

See Also

Blocks

Complex Burst Matrix Solve Using Q-less QR Decomposition with Forgetting Factor

Functions

`fixed.forgettingFactor` | `upperbound` |

`fixed.complexQlessQRMatrixSolveFixedpointTypes` | `fixed.extractNumericType`

Implement Hardware-Efficient Real Burst Asynchronous Matrix Solve Using Q-less QR Decomposition

This example shows how to implement a hardware-efficient solution to the real-valued matrix equation $A'AX=B$ using the Real Burst Asynchronous Matrix Solve Using Q-less QR Decomposition block.

Forward and Backward Substitution

When an upper triangular factor is ready, then forward and backward substitution are computed with the current input B to produce output X.

$$X = R_k \setminus (R'_k \setminus B)$$

Define Matrix Dimensions

Specify the number of rows in matrix A, the number of columns in matrix A and rows in B, and the number of columns in matrix B.

```
m = 30; % Number of rows in A
n = 10; % Number of columns in A and rows in B
p = 1; % Number of columns in B
numInputs = 3; % Number of A and B matrices
```

Generate Matrices

For this example, use the helper function `realRandomQlessQRMatrices` to generate random matrices A and B for the problem $A'AX=B$. The matrices are generated such that the elements of A and B are between -1 and +1, and A is full rank.

```
rng('default')
[A,B] = fixed.example.realRandomQlessQRMatrices(m,n,p);
if numInputs > 1
    for i = 2:numInputs
        [Atemp,Btemp] = fixed.example.realRandomQlessQRMatrices(m,n,p);
        A = cat(3,A,Atemp);
        B = cat(3,B,Btemp);
    end
end
```

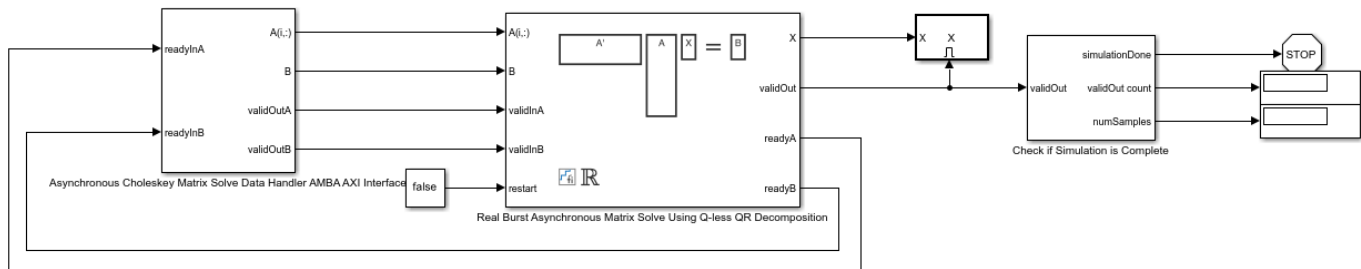
Select Fixed-Point Data Types

Use the helper function `realQlessQRMatrixSolveFixedpointTypes` to select fixed-point data types for input matrices A and B, and output X such that there is a low probability of overflow during the computation.

```
max_abs_A = 1; % Upper bound on max(abs(A(:)))
max_abs_B = 1; % Upper bound on max(abs(B(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.realQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,precisionBits);
A = cast(A,'like',T.A);
B = cast(B,'like',T.B);
OutputType = fixed.extractNumericType(T.X);
```

Open the Model

```
model = 'RealBurstAsyncQlessQRMatrixSolveModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

AMBA AXI Handshaking Process

The Data Handler subsystem in this model takes real matrices A and B as inputs. It sends rows of A and full matrix of B to the QR Decomposition block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set delays for the feeding in rows of A and the feeding in of B matrices in the Data Handler to emulate the processing time of the upstream block. `validInA` and `validInB` remain high when `aDelay` and `bDelay` are set to 0 because this indicates the Data Handler always has data available. When all matrices A and B are sent, the Data Handler loops back to the first A and B matrices.

Asynchronous Matrix Solver

This block operates asynchronously. First, Q-less QR decomposition is performed on the input A matrix and the resulting R matrix is put into a buffer. Then, the Forward Backward Substitute block uses the input B matrix and the buffered R matrix to compute $R'RX = B$. Because the R and B matrices are stored separately in buffers, the upstream Q-less QR decomposition block and the downstream Forward Backward Substitute block can run independently. The Forward Backward Substitute block starts processing when the first R and B matrices are available. Then it runs continuously using the latest buffered R and B matrices, regardless of the status of the Q-less QR Decomposition block. For example, if the upstream block stops providing A and B matrices, the Forward Backward Substitute block continues to generate the same output using the last pair of R and B matrices.

The Data Handler sends A and B matrices to the QR decomposition block iteratively. After sending out the last A matrix, the Data Handler resets its internal counter and sends out first A matrix. The B matrix is handled in a similar fashion.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Real Burst Asynchronous Matrix Solve Using Q-less QR Decomposition block.

```
numOutputs = 10; % Number of recorded outputs
aDelay = 1; % Delay of clock cycles between feeding in rows of A
bDelay = 1; % Delay of clock cycles between feeding in B matrices
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
```

```
'regularizationParameter',0,...
'aDelay',aDelay,'bDelay',bDelay,...
'numOutputs',numOutputs,'OutputType',OutputType);
```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Real Burst Asynchronous Matrix Solve Using Q-less QR Decomposition block outputs matrix X at each time step. When a valid result matrix is output, the block sets `validOut` to true.

```
X = out.X;
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Real Burst Asynchronous Matrix Solve Using Q-less QR Decomposition block, compute the relative error. Choose the last output of the simulation.

```
X = double(X(:,:,end));
```

Synchronize the last output X with the input by finding the inputs A and B that produced it.

```
A = double(A);
B = double(B);
relative_errors = zeros(size(A,3),size(B,3));
for k = 1:size(A,3)
    for g = 1:size(B,3)
        relative_errors(k,g) = norm(A(:,:,k)'*A(:,:,k)*X - B(:,:,g))/norm(B(:,:,g));
    end
end
```

```
[AUsed,Bused] = find(relative_errors==min(relative_errors,[],'all')) %#ok<NOPTS>
```

```
relative_error = norm(double(A(:,:,AUsed))*A(:,:,AUsed)*X - B(:,:,Bused))/norm(double(B(:,:,Bused))
```

```
AUsed =
```

```
2
```

```
Bused =
```

```
3
```

```
relative_error =
```

```
0.0012
```

See Also

Real Burst Asynchronous Matrix Solve Using Q-less QR Decomposition

Implement Hardware-Efficient Complex Burst Asynchronous Matrix Solve Using Q-less QR Decomposition

This example shows how to implement a hardware-efficient solution to the complex-valued matrix equation $A'AX=B$ using the Complex Burst Asynchronous Matrix Solve Using Q-less QR Decomposition block.

Forward and Backward Substitution

When an upper triangular factor is ready, then forward and backward substitution are computed with the current input B to produce output X.

$$X = R_k \setminus (R'_k \setminus B)$$

Define Matrix Dimensions

Specify the number of rows in matrix A, the number of columns in matrix A and rows in B, and the number of columns in matrix B.

```
m = 30; % Number of rows in A
n = 10; % Number of columns in A and rows in B
p = 1; % Number of columns in B
numInputs = 3; % Number of A and B matrices
```

Generate Matrices

For this example, use the helper function `complexRandomQlessQRMatrices` to generate random matrices A and B for the problem $A'AX=B$. The matrices are generated such that the real and imaginary parts of the elements of A and B are between -1 and +1, and A is full rank.

```
rng('default')
[A,B] = fixed.example.complexRandomQlessQRMatrices(m,n,p);
if numInputs > 1
    for i = 2:numInputs
        [Atemp,Btemp] = fixed.example.complexRandomQlessQRMatrices(m,n,p);
        A = cat(3,A,Atemp);
        B = cat(3,B,Btemp);
    end
end
```

Select Fixed-Point Data Types

Use the helper function `complexQlessQRMatrixSolveFixedpointTypes` to select fixed-point data types for input matrices A and B, and output X such that there is a low probability of overflow during the computation.

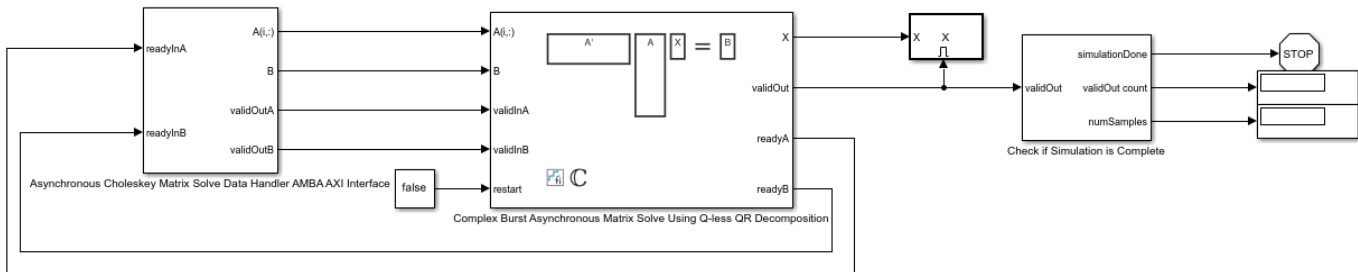
The real and imaginary parts of the elements of A and B are between -1 and 1, so the maximum possible absolute value of any element is $\sqrt{2}$.

```
max_abs_A = sqrt(2); % Upper bound on max(abs(A(:)))
max_abs_B = sqrt(2); % Upper bound on max(abs(B(:)))
precisionBits = 24; % Number of bits of precision
T = fixed.complexQlessQRMatrixSolveFixedpointTypes(m,n,max_abs_A,max_abs_B,precisionBits);
A = cast(A,'like',T.A);
```

```
B = cast(B, 'like', T.B);
OutputType = fixed.extractNumericType(T.X);
```

Open the Model

```
model = 'ComplexBurstAsyncQlessQRMatrixSolveModel';
open_system(model);
```



Copyright 2022 The MathWorks, Inc.

AMBA AXI Handshaking Process

The Data Handler subsystem in this model takes complex matrices A and B as inputs. It sends rows of A and full matrix of B to the QR Decomposition block using the AMBA AXI handshake protocol. The `validIn` signal indicates when data is available. The `ready` signal indicates that the block can accept the data. Transfer of data occurs only when both the `validIn` and `ready` signals are high. You can set delays for the feeding in rows of A and the feeding in of B matrices in the Data Handler to emulate the processing time of the upstream block. `validInA` and `validInB` remain high when `aDelay` and `bDelay` are set to 0 because this indicates the Data Handler always has data available. When all matrices A and B are sent, the Data Handler loops back to the first A and B matrices.

Asynchronous Matrix Solver

This block operates asynchronously. First, Q-less QR decomposition is performed on the input A matrix and the resulting R matrix is put into a buffer. Then, the Forward Backward Substitute block uses the input B matrix and the buffered R matrix to compute $R'RX = B$. Because the R and B matrices are stored separately in buffers, the upstream Q-less QR decomposition block and the downstream Forward Backward Substitute block can run independently. The Forward Backward Substitute block starts processing when the first R and B matrices are available. Then it runs continuously using the latest buffered R and B matrices, regardless of the status of the Q-less QR Decomposition block. For example, if the upstream block stops providing A and B matrices, the Forward Backward Substitute block continues to generate the same output using the last pair of R and B matrices.

The Data Handler sends A and B matrices to the QR decomposition block iteratively. After sending out the last A matrix, the Data Handler resets its internal counter and sends out first A matrix. The B matrix is handled in a similar fashion.

Set Variables in the Model Workspace

Use the helper function `setModelWorkspace` to add the variables defined above to the model workspace. These variables correspond to the block parameters for the Complex Burst Asynchronous Matrix Solve Using Q-less QR Decomposition block.

```
numOutputs = 10; % Number of recorded outputs
aDelay = 1; % Delay of clock cycles between feeding in rows of A
```

```

bDelay = 1; % Delay of clock cycles between feeding in B matrices
fixed.example.setModelWorkspace(model, 'A', A, 'B', B, 'm', m, 'n', n, 'p', p, ...
    'regularizationParameter', 0, ...
    'aDelay', aDelay, 'bDelay', bDelay, ...
    'numOutputs', numOutputs, 'OutputType', OutputType);

```

Simulate the Model

```
out = sim(model);
```

Construct the Solution from the Output Data

The Complex Burst Asynchronous Matrix Solve Using Q-less QR Decomposition block outputs matrix X at each time step. When a valid result matrix is output, the block sets `validOut` to true.

```
X = out.X;
```

Verify the Accuracy of the Output

To evaluate the accuracy of the Complex Burst Asynchronous Matrix Solve Using Q-less QR Decomposition block, compute the relative error. Choose the last output of the simulation.

```
X = double(X(:,:,end));
```

Synchronize the last output X with the input by finding the inputs A and B that produced it.

```

A = double(A);
B = double(B);
relative_errors = zeros(size(A,3),size(B,3));
for k = 1:size(A,3)
    for g = 1:size(B,3)
        relative_errors(k,g) = norm(A(:,:,k)'*A(:,:,k)*X - B(:,:,g))/norm(B(:,:,g));
    end
end
[AUsed,Bused] = find(relative_errors==min(relative_errors,[],'all')) %#ok<NOPTS>

```

```
relative_error = norm(double(A(:,:,AUsed)'*A(:,:,AUsed)*X - B(:,:,Bused)))/norm(double(B(:,:,Bused)));
```

```
AUsed =
```

```
1
```

```
Bused =
```

```
3
```

```
relative_error =
```

```
2.2034e-04
```

See Also

Complex Burst Asynchronous Matrix Solve Using Q-less QR Decomposition

How to Use Square Jacobi SVD HDL Optimized Block

This example shows how to use the Square Jacobi SVD HDL Optimized block to compute the singular value decomposition (SVD) of square matrices.

Two-Sided Jacobi SVD

The Square Jacobi HDL Optimized block uses the two-sided Jacobi algorithm to perform singular value decomposition. Given an input square matrix A, the block first computes the two-by-two SVD for off-diagonal elements, then applies the rotation to the A, U, and V matrices. Because the Jacobi algorithm can perform such computations in parallel, it is suitable for FPGA and ASIC applications. For more information, see Square Jacobi SVD HDL Optimized.

Define Simulation Parameters

Specify the dimension of the sample matrices, the number of input sample matrices, and the number of iterations of the Jacobi algorithm.

```
n = 8;
numSamples = 3;
nIterations = 10;
```

Generate Input A Matrices

Use the specified simulation parameters to generate the input matrix A.

```
rng('default');
A = randn(n,n,numSamples);
```

The Square Jacobi SVD HDL Optimized block supports both real and complex inputs. Set the complexity of the input in the block mask accordingly.

```
% A = complex(randn(n,n,numSamples),randn(n,n,numSamples));
```

Select Fixed-Point Data Types

Define the desired word length.

```
wordLength = 32;
```

Use the upper bound on the singular values to define fixed-point types that will never overflow. First, use the `fixed.singularValueUpperBound` function to determine the upper bound on the singular values.

```
svdUpperBound = fixed.singularValueUpperBound(n,n,max(abs(A(:))));
```

Define the integer length based on the value of the upper bound, with one additional bit for the sign, another additional bit for intermediate CORDIC growth, and one more bit for intermediate growth to compute the Jacobi rotations.

```
additionalBitGrowth = 3;
integerLength = ceil(log2(svdUpperBound)) + additionalBitGrowth;
```

Compute the fraction length based on the integer length and the desired word length.

```
fractionLength = wordLength - integerLength;
```

Define the signed fixed-point data type to be 'Fixed' or 'ScaledDouble'. You can also define the type to be 'double' or 'single'.

```
dataType = 'Fixed';
T.A = fi([],1,wordLength,fractionLength,'DataType',dataType);
disp(T.A)
```

```
[]
```

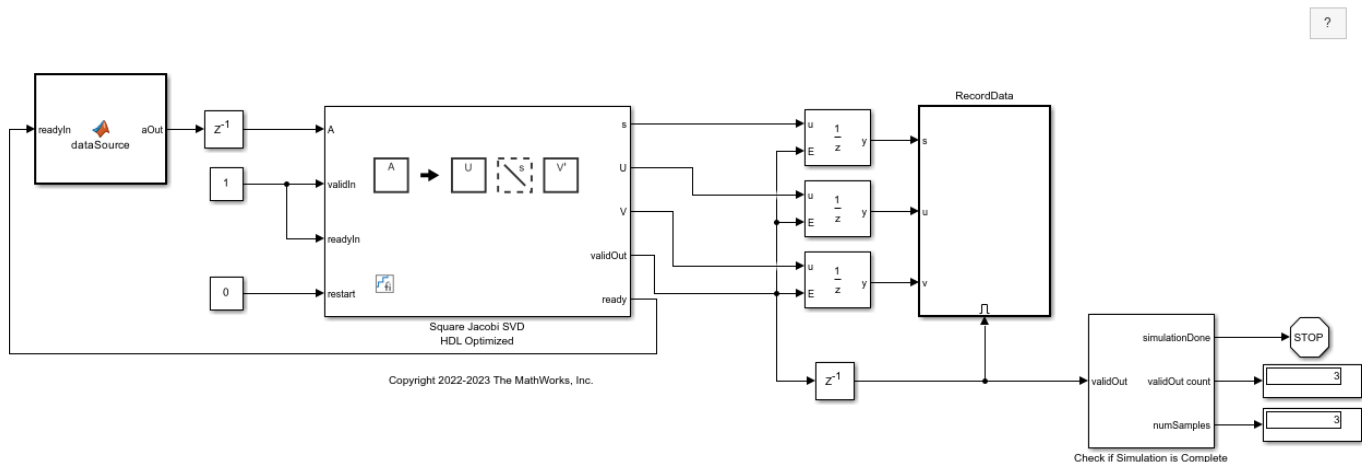
```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 32
FractionLength: 24
```

Cast the matrix A to the signed fixed-point type.

```
A = cast(A,'like',T.A);
```

Configure Model Workspace and Run Simulation

```
model = 'SquareJacobiSVDModel';
open_system(model);
setModelWorkspace(model,'A',A,'n',n,...
    'nIterations',nIterations,'numSamples',numSamples);
out = sim(model);
```



Verify Output Solutions

Verify the output solutions. In these steps, "identical" means within roundoff error.

- 1 Verify that $U \cdot \text{diag}(s) \cdot V'$ is identical to A. `relativeErrorUSV` represents the relative error between $U \cdot \text{diag}(s) \cdot V'$ and A.
- 2 Verify that the singular values `s` are identical to the floating-point SVD solution. `relativeErrorS` represents the relative error between `s` and the singular values calculated by the MATLAB® `svd` function.
- 3 Verify that `U` and `V` are unitary matrices. `relativeErrorUU` represents the relative error between $U' \cdot U$ and the identity matrix. `relativeErrorVV` represents the relative error between $V' \cdot V$ and the identity matrix.

```

for i = 1:numSamples
    disp(['Sample #',num2str(i),':']);
    a = A(:,:,i);
    U = out.U(:,:,i);
    V = out.V(:,:,i);
    s = out.s(:,:,i);

    % Verify U*diag(s)*V'
    if norm(double(a)) > 1
        relativeErrorUSV = norm(double(U*diag(s)*V')-double(a))/norm(double(a));
    else
        relativeErrorUSV = norm(double(U*diag(s)*V')-double(a));
    end
    relativeErrorUSV %#ok

    % Verify s
    s_expected = svd(double(a));
    normS = norm(s_expected);
    relativeErrorS = norm(double(s) - s_expected);
    if normS > 1
        relativeErrorS = relativeErrorS/normS;
    end
    relativeErrorS %#ok

    % Verify U'*U
    U = double(U);
    UU = U'*U;
    relativeErrorUU = norm(UU - eye(size(UU))) %#ok

    % Verify V'*V
    V = double(V);
    VV = V'*V;
    relativeErrorVV = norm(VV - eye(size(VV))) %#ok

    disp('-----');
end

```

Sample #1:

relativeErrorUSV =

4.9236e-06

relativeErrorS =

2.4379e-06

relativeErrorUU =

5.9432e-07

relativeErrorVV =

6.9467e-07

```
-----  
Sample #2:  
  
relativeErrorUSV =  
    6.0158e-06  
  
relativeErrorS =  
    2.4712e-06  
  
relativeErrorUU =  
    6.0220e-07  
  
relativeErrorVV =  
    5.2963e-07  
  
-----  
Sample #3:  
  
relativeErrorUSV =  
    5.7222e-06  
  
relativeErrorS =  
    2.9780e-06  
  
relativeErrorUU =  
    5.3064e-07  
  
relativeErrorVV =  
    5.2115e-07  
  
-----
```

See Also

Blocks

Square Jacobi SVD HDL Optimized

Functions

`fixed.singularValueUpperBound`

Implement HDL Optimized SVD in Feedforward Fashion Without Backpressure

This example shows how to implement a hardware-efficient singular value decomposition (SVD) using the Square Jacobi SVD HDL Optimized block in a feedforward fashion without backpressure.

The Square Jacobi SVD HDL Optimized block uses the AMBA AXI handshake protocol for both input and output. To use the block without backpressure control, feed a constant Boolean 'true' to the **readyIn** port, then configure the upstream input rate according to the block latency specified in Square Jacobi SVD HDL Optimized.

This model uses a Rate Transition block between the data source and the Square Jacobi SVD HDL Optimized block to emulate an upstream block with a lower sample rate. The sample time of the Square Jacobi SVD HDL Optimized block is normalized to 1, and the sample time of data source is the block latency + 1.

Define Simulation Parameters

Specify the dimension of the sample matrices, the number of input sample matrices, and the number of iterations of the Jacobi algorithm.

```
n = 8;
numSamples = 3;
nIterations = 6;
```

Generate Input A Matrices

Use the specified simulation parameters to generate the input A matrices.

```
rng('default');
A = randn(n,n,numSamples);
```

The Square Jacobi SVD HDL Optimized block supports both real and complex inputs. Set the complexity of the input in the block mask accordingly.

```
% A = complex(randn(n,n,numSamples),randn(n,n,numSamples));
```

Select Fixed-Point Data Types

Define the desired word length.

```
wordLength = 32;
```

Use the upper bound on the singular values to define fixed-point types that will never overflow. First, use the `fixed.singularValueUpperBound` function to determine the upper bound on the singular values.

```
svdUpperBound = fixed.singularValueUpperBound(n,n,max(abs(A(:))));
```

Define the integer length based on the value of the upper bound, with one additional bit for the sign, another additional bit for intermediate CORDIC growth, and one more bit for intermediate growth to compute the Jacobi rotations.

```
additionalBitGrowth = 3;
integerLength = ceil(log2(svdUpperBound)) + additionalBitGrowth;
```

Compute the fraction length based on the integer length and the desired word length.

```
fractionLength = wordLength - integerLength;
```

Define the signed fixed-point data type to be 'Fixed' or 'ScaledDouble'. You can also define the type to be 'double' or 'single'.

```
dataType = 'Fixed';
T.A = fi([],1,wordLength,fractionLength,'DataType',dataType);
disp(T.A)
```

```
[]
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 32
      FractionLength: 24
```

Cast the matrix A to the signed fixed-point type.

```
A = cast(A,'like',T.A);
```

Set Input Rate

The latency of the Square Jacobi SVD HDL Optimized block depends on the size n , complexity, and word length wl of the input matrix A, and the number of iterations $nIterations$ of the two-sided Jacobi algorithm.

- If signal complexity is real, the block delay is $(wl*2+31)*(n-1+rem(n,2))*nIterations + 2 + nextpow2(n)*(nextpow2(n)+1)/2+3$.
- If signal complexity is complex, the block delay is $(wl*6+48)*(n-1+rem(n,2))*nIterations + 2 + nextpow2(n)*(nextpow2(n)+1)/2+3$.
- If A is double precision, then wl is 53.
- if A is single precision, then wl is 24.

The Square Jacobi SVD HDL Optimized block will be ready in the next clock after a successful solution output. Therefore, the input rate transition ratio should be at least $1 + \text{block latency}$.

For real input, compute the rate transition ratio.

```
rateTransitionRatio = 1+(wordLength*2+31)*(n-1+rem(n,2))*nIterations+2+nextpow2(n)*(nextpow2(n)+1)/2+3;
```

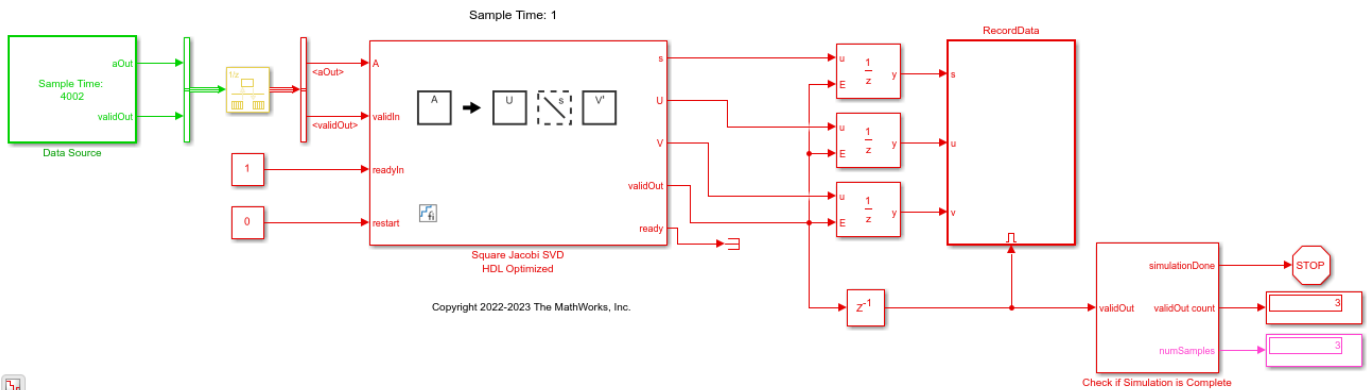
For complex input, compute the rate transition ratio.

```
% rateTransitionRatio = 1+(wordLength*6+48)*(n-1+rem(n,2))*nIterations+2+nextpow2(n)*(nextpow2(n)+1)/2+3;
```

Configure Model Workspace and Run Simulation

```
model = 'SquareJacobiSVDFeedforwardModel';
open_system(model);
setModelWorkspace(model,'A',A,'n',n,...
    'nIterations',nIterations,'numSamples',numSamples,...
    'rateTransitionRatio',rateTransitionRatio);
out = sim(model);
```

?



Verify Output Solutions

Verify the output solutions. In these steps, "identical" means within roundoff error.

- 1 Verify that $U \cdot \text{diag}(s) \cdot V'$ is identical to A . `relErrUSV` represents the relative error between $U \cdot \text{diag}(s) \cdot V'$ and A .
- 2 Verify that the singular values s are identical to the floating point SVD solution. `relativeErrors` represents the relative error between s and the singular values calculated by MATLAB® `svd` function.
- 3 Verify that U and V are unitary matrices. `relativeErrorUU` represents the relative error between $U' \cdot U$ and the identity matrix. `relativeErrorVV` represents the relative error between $V' \cdot V$ and the identity matrix.

```

for i = 1:numSamples
    disp(['Sample #', num2str(i), ':']);
    a = A(:,:,i);
    U = out.U(:,:,i);
    V = out.V(:,:,i);
    s = out.s(:,:,i);

    % Verify U*diag(s)*V'
    if norm(double(a)) > 1
        relativeErrorUSV = norm(double(U*diag(s)*V')-double(a))/norm(double(a));
    else
        relativeErrorUSV = norm(double(U*diag(s)*V')-double(a));
    end
    relativeErrorUSV %#ok

    % Verify s
    s_expected = svd(double(a));
    normS = norm(s_expected);
    relativeErrorS = norm(double(s) - s_expected);
    if normS > 1
        relativeErrorS = relativeErrorS/normS;
    end
    relativeErrorS %#ok

    % Verify U'*U
    U = double(U);

```

```
UU = U'*U;
relativeErrorUU = norm(UU - eye(size(UU))) %#ok

% Verify V'*V
V = double(V);
VV = V'*V;
relativeErrorVV = norm(VV - eye(size(VV))) %#ok

disp('-----');
end
```

Sample #1:

```
relativeErrorUSV =
    2.9950e-06
```

```
relativeErrorS =
    1.3585e-06
```

```
relativeErrorUU =
    3.3832e-07
```

```
relativeErrorVV =
    4.0781e-07
```

```
-----
Sample #2:
```

```
relativeErrorUSV =
    3.9727e-06
```

```
relativeErrorS =
    1.4026e-06
```

```
relativeErrorUU =
    3.4657e-07
```

```
relativeErrorVV =
    3.0048e-07
```

```
-----
Sample #3:
```

```
relativeErrorUSV =
```


3.7357e-06

relativeErrors =

1.7264e-06

relativeErrorUU =

2.7650e-07

relativeErrorVV =

3.0022e-07

See Also

Blocks

Square Jacobi SVD HDL Optimized

Functions

fixed.singularValueUpperBound

Implement HDL Optimized SVD with Backpressure Signal and HDL FIFO Block

This example shows how to implement hardware-efficient singular value decomposition (SVD) using the Square Jacobi SVD HDL Optimized block with backpressure control and an HDL FIFO block.

The Square Jacobi SVD HDL Optimized block uses the AMBA AXI handshake protocol for both input and output. The valid/ready handshake process is used to transfer data and control information. For more details about the handshake process, see Square Jacobi SVD HDL Optimized.

In this example, the data handler has an arbitrary delay to emulate upstream delay. The dummy receiver emulates downstream delay. The model has two synchronous FIFO blocks inserted between the upstream data handler block and Square Jacobi SVD HDL Optimized block, as well as between the Square Jacobi SVD HDL Optimized block and the downstream receiver. For all the blocks, the downstream ready signal connects to the upstream **readyIn** port, and the upstream **validOut** signal connects to the downstream **validIn** port.

The Square Jacobi SVD HDL Optimized block supports backpressure control through the AMBA AXI protocol. The Synchronous FIFO block uses the HDL FIFO block with glue logic to support the AMBA AXI protocol. This example uses the FIFO blocks to demonstrate how to interface the Square Jacobi SVD HDL Optimized block and the FIFO block with backpressure control.

Define Simulation Parameters

Specify the dimension of the sample matrices, the number of input sample matrices, and the number of iterations of the Jacobi algorithm.

```
n = 2;
numSamples = 20;
nIterations = 6;
```

Specify the upstream and downstream delays.

```
upstreamDelay = 200;
downstreamDelay = 1000;
```

Generate Input A Matrices

Use the specified simulation parameters to generate the input A matrices.

```
rng('default');
A = randn(n,n,numSamples);
```

The Square Jacobi SVD HDL Optimized block supports both real and complex inputs. Set the complexity of the input in the block mask accordingly.

```
% A = complex(randn(n,n,numSamples), randn(n,n,numSamples));
```

Select Fixed-Point Data Types

Define the desired word length.

```
wordLength = 32;
```

Use the upper bound on the singular values to define fixed-point types that will never overflow. First, use the `fixed.singularValueUpperBound` function to determine the upper bound on the singular values.

```
svdUpperBound = fixed.singularValueUpperBound(n,n,max(abs(A(:))));
```

Define the integer length based on the value of the upper bound, with one additional bit for the sign, another additional bit for intermediate CORDIC growth, and one more bit for intermediate growth to compute the Jacobi rotations.

```
additionalBitGrowth = 3;
integerLength = ceil(log2(svdUpperBound)) + additionalBitGrowth;
```

Compute the fraction length based on the integer length and the desired word length.

```
fractionLength = wordLength - integerLength;
```

Define the signed fixed-point data type to be 'Fixed' or 'ScaledDouble'. You can also define the type to be 'double' or 'single'.

```
dataType = 'Fixed';
T.A = fi([],1,wordLength,fractionLength,'DataType',dataType);
disp(T.A)
```

```
[]
```

```

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 32
    FractionLength: 26
    
```

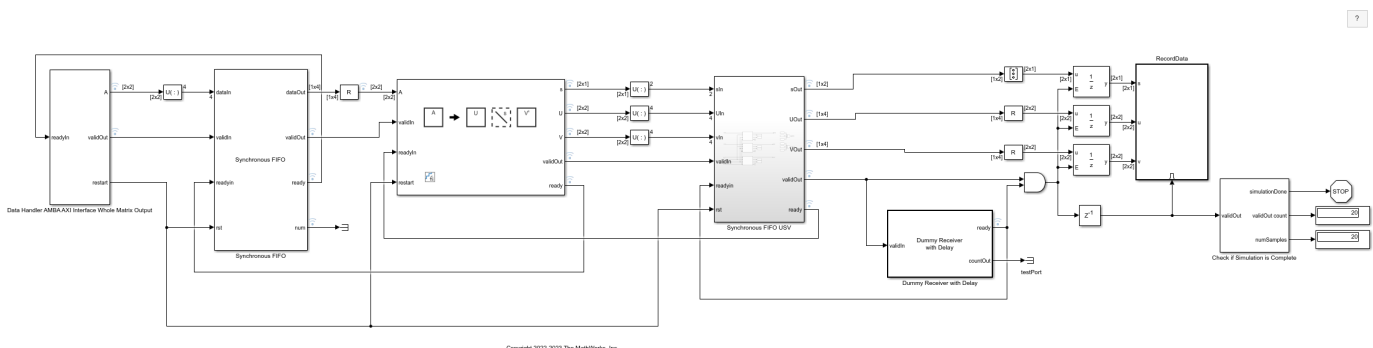
Cast the matrix A to the signed fixed-point type.

```
A = cast(A,'like',T.A);
```

Configure Model Workspace and Run Simulation

```

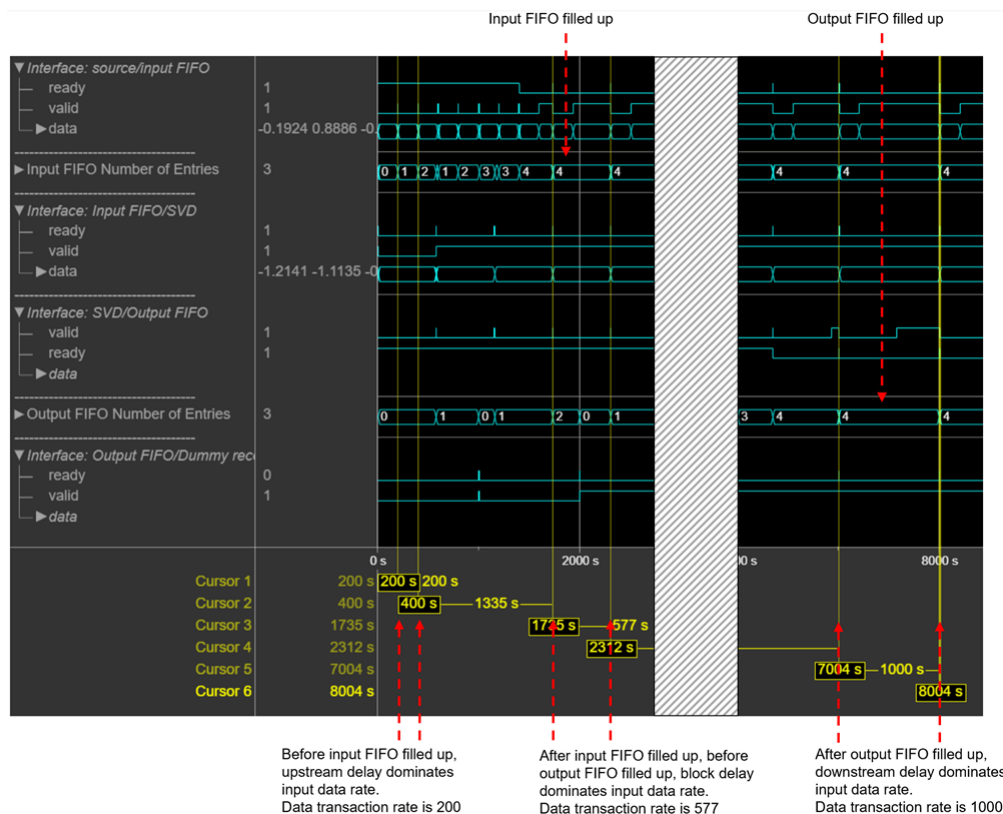
model = 'SquareJacobiSVDFFIFOModel';
open_system(model);
setModelWorkspace(model,'A',A,'n',n,...
    'nIterations',nIterations,'numSamples',numSamples,...
    'upstreamDelay',upstreamDelay,'downstreamDelay',downstreamDelay);
out = sim(model);
    
```



Handshake Process and Backpressure Signal

In this model, the data input rate is 200, the block latency is 577, and the dummy receiver delay is 1000. The FIFO depth is 4. Because the downstream is slower than upstream, the expected behavior is:

- 1 Before the input FIFO is full, the data source rate determines the data transaction rate. The input FIFO accepts data every 200 clocks.
- 2 After the input FIFO is full, it can only accept data when the Square Jacobi SVD HDL Optimized block is ready. The data transaction rate reduces to the block delay of 577.
- 3 The Square Jacobi SVD HDL Optimized block outputs data into the output FIFO, and the dummy receiver consumes the solution every 1000 clocks. Before the output FIFO is full, the data transaction rate is the same as the Square Jacobi SVD HDL Optimized block delay of 577.
- 4 After the output FIFO is full, it can only accept data when the dummy receiver is ready. The Square Jacobi SVD HDL Optimized block waits for the dummy receiver, and the data source waits for the Square Jacobi SVD HDL Optimized block. The data transaction rate of the signal chain reduces to every 1000 clocks.



By using backpressure signals and the handshake protocol, the upstream block waits for the downstream block without a hardcoded delay. The slowest block determines the throughput of the whole system to ensure data integrity.

Verify Output Solutions

Verify the output solutions. In these steps, "identical" means within roundoff error.

- 1 Verify that $U \cdot \text{diag}(s) \cdot V'$ is identical to A . `relErrUSV` represents the relative error between $U \cdot \text{diag}(s) \cdot V'$ and A .
- 2 Verify that the singular values s are identical to the floating-point SVD solution. `relativeErrorS` represents the relative error between s and singular values calculated by the MATLAB® `svd` function.
- 3 Verify that U and V are unitary matrices. `relativeErrorUU` represents the relative error between $U' \cdot U$ and the identity matrix. `relativeErrorVV` represents the relative error between $V' \cdot V$ and the identity matrix.

```

relativeErrorUSV = zeros(numSamples,1);
relativeErrorUU = zeros(numSamples,1);
relativeErrorVV = zeros(numSamples,1);
relativeErrorS = zeros(numSamples,1);

for i = 1:numSamples
    a = A(:, :, i);
    U = out.U(:, :, i);
    V = out.V(:, :, i);
    s = out.s(:, :, i);

    % Verify U*diag(s)*V'
    if norm(double(a)) > 1
        relErrUSV = norm(double(U*diag(s)*V')-double(a))/norm(double(a));
    else
        relErrUSV = norm(double(U*diag(s)*V')-double(a));
    end
    relativeErrorUSV(i) = relErrUSV;

    % Verify s
    s_expected = svd(double(a));
    normS = norm(s_expected);
    relErrS = norm(double(s) - s_expected);
    if normS > 1
        relErrS = relErrS/normS;
    end
    relativeErrorS(i) = relErrS;

    % Verify U'*U
    U = double(U);
    UU = U'*U;
    relativeErrorUU(i) = norm(UU - eye(size(UU)));

    % Verify V'*V
    V = double(V);
    VV = V'*V;
    relativeErrorVV(i) = norm(VV - eye(size(VV)));
end
disp(['Maximum s error: ', num2str(max(relativeErrorS))]);
disp(['Maximum UsV error: ', num2str(max(relativeErrorUSV))]);
disp(['Maximum UU error: ', num2str(max(relativeErrorUU))]);
disp(['Maximum VV error: ', num2str(max(relativeErrorVV))]);

Maximum s error: 2.0123e-07
Maximum UsV error: 2.1179e-07

```

Maximum UU error: 4.2175e-08
Maximum VV error: 4.4717e-08

See Also

Blocks

Square Jacobi SVD HDL Optimized

Functions

fixed.singularValueUpperBound

Troubleshooting

- “Fixed-Point Versus Built-in Integer Types” on page 49-2
- “Negative Fraction Length” on page 49-3
- “Fraction Length Greater Than Word Length” on page 49-5
- “fi Constructor Does Not Follow globalfimath Rules” on page 49-7
- “Decide Which Workflow is Right for Your Application” on page 49-8
- “Tips for Making Generated Code More Efficient” on page 49-9
- “Know When a Function is Supported for Instrumentation and Acceleration” on page 49-11
- “Resolve Error: Function is not Supported for Fixed-Point Conversion” on page 49-12
- “Resolve Error: fi*non-fi” on page 49-14
- “Resolve Error: Data Type Mismatch” on page 49-15
- “Resolve Error: Mismatched fimath” on page 49-16
- “Why Does the Fixed-Point Converter App Not Propose Data Types for System Objects?” on page 49-17
- “Slow Operations in the Fixed-Point Converter App” on page 49-18
- “Blocks That Do Not Support Fixed-Point Data Types” on page 49-19
- “Prevent the Fixed-Point Tool from Overriding Integer Data Types” on page 49-21
- “The Fixed-Point Tool did not Propose Data Types” on page 49-22
- “Fraction Lengths and Fixed-Point Numbers” on page 49-23
- “Why am I missing data type proposals for MATLAB Function block variables?” on page 49-24
- “Data Type Propagation Errors After Applying Proposed Data Types” on page 49-25
- “Resolve Range Analysis Issues” on page 49-27
- “Data Type Mismatch and Structure Initial Conditions” on page 49-28
- “Reconversion Using the Fixed-Point Tool” on page 49-30
- “Data Type Optimization Not Successful” on page 49-31
- “Compile-Time Recursion Limit Reached” on page 49-33
- “Output Variable Must Be Assigned Before Run-Time Recursive Call” on page 49-36
- “Unable to Determine That Every Element of Cell Array Is Assigned” on page 49-39
- “Nonconstant Index into varargin or varargout in a for-Loop” on page 49-43

Fixed-Point Versus Built-in Integer Types

There are several distinct differences between fixed-point data types and the built-in integer types in MATLAB. The most notable difference is that the built-in integer data types can only represent whole numbers, while the fixed-point data types also contain information on the position of the binary point, or the scaling of the number. This scaling allows the fixed-point data types to represent both integers and non-integers.

There are also slight differences in how math is performed with these types. Fixed-point types allow you to specify rules for math using the `fimath` object, including overflow and rounding modes. However, the built-in types have their own internal rules for arithmetic operations. See “Integers” for more information on how math is performed using built-in types.

See Also

`fi` | `fimath`

Negative Fraction Length

A negative fraction length occurs when the input value of a `fi` object contains trailing zeros before the decimal point. For example,

```
x = fi(16000,1,8)
```

produces a signed fixed-point number with a word length of 8 bits and best precision fraction length.

```
x =
```

```
16000
```

```

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 8
    FractionLength: -7

```

View the binary representation of `x`.

```
disp(bin(x))
```

```
01111101
```

There are seven implicit zeros at the end of this number before the binary point because the fraction length of `x` is `-7`.

Convert from binary to decimal the binary representation of `x` with seven zero bits appended to the end.

```
bin2dec('011111010000000')
```

```
ans =
```

```
16000
```

The result is the real world value of `x`.

You can also find the real world value using the equation

$$\text{Real World Value} = \text{Stored Integer Value} \times 2^{-\text{Fraction Length}}$$

Start by finding the stored integer of `x`.

```
Q = storedInteger(x)
```

```
Q =
```

```
125
```

Use the stored integer to find the real world value of `x`.

```
real_world_value = double(Q) * 2^-x.FractionLength
```

```
real_world_value =
```

```
16000
```

See Also

fi

Fraction Length Greater Than Word Length

A fraction length greater than the word length of a fixed-point number occurs when the number has an absolute value less than one and contains leading zeros.

```
x = fi(.0234,1,8)
```

```
x =
```

```
0.0234
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 8
    FractionLength: 12
```

View the binary representation of x.

```
disp(bin(x))
```

```
01100000
```

There are four implicit leading zeros after the binary point and before the binary representation of the stored integer because the fraction length of x is four greater than the word length.

Convert from binary to decimal the binary representation of x with four leading zeros, and scale this value according to the fraction length.

```
bin2dec('000001100000')*2^(-12)
```

```
ans =
```

```
0.0234
```

The result is the real world value of x.

You can also find the real world value using the equation
 Real World Value = Stored Integer Value $\times 2^{-\text{Fraction Length}}$.

Start by finding the stored integer of x.

```
Q = storedInteger(x)
```

```
Q =
```

```
96
```

Use the stored integer to find the real world value of x.

```
real_world_value = double(Q) * 2^-x.FractionLength
```

```
real_world_value =
```

```
0.0234
```

See Also

fi

fi Constructor Does Not Follow globalfimath Rules

Issue

If no fimath properties are used in the argument of the `fi` constructor, then it always uses nearest rounding and saturates on overflow for the creation of the `fi` object, regardless of any `globalfimath` settings.

Possible Solutions

If this behavior is undesirable for your application, you can do one of the following:

Use the cast Function to Create a fi Object Using the globalfimathrules

```
G = globalfimath('RoundingMethod', 'Floor', 'OverflowAction', 'Wrap');  
cast(x, 'like', fi([],1,16,10))
```

When you create a `fi` object using the `cast` function, the resulting `fi` object does not have a local `fimath`.

Specify fimath Properties in the fi Constructor

```
fi(x,1,16,10,'RoundingMethod','Floor','OverflowAction','Wrap');
```

When you create a `fi` object with `fimath` properties in the constructor, the `fi` object does have a local `fimath`.

See Also

`fi` | `fimath` | `globalfimath`

Decide Which Workflow is Right for Your Application

There are two primary workflows available for converting MATLAB code to fixed-point code.

- **Manual Workflow**

The manual workflow provides the most control to optimize the fixed-point types, but requires a greater understanding of fixed-point concepts.

For more information, see “Manual Fixed-Point Conversion Best Practices” on page 11-3.

- **Automated Workflow**

The Fixed-Point Converter app enables you to convert your MATLAB code to fixed-point code without requiring extensive preexisting knowledge of fixed-point concepts. However, this workflow provides less control over your data types.

For more information, see “Automated Fixed-Point Conversion Best Practices” on page 7-44.

| | Manual Workflow | Automated Workflow |
|---|-----------------|--------------------|
| Fully automated conversion | | ✓ |
| Less fixed-point expertise required | | ✓ |
| Quick turnaround time | | ✓ |
| Simulation range analysis | ✓ | ✓ |
| Static range analysis | | ✓ |
| Iterative workflow | ✓ | |
| Portable design | ✓ | ✓ |
| Greatest control and optimization of data types | ✓ | |
| Data type proposal | ✓ | ✓ |
| Histogram logging | ✓ | ✓ |
| Code coverage | | ✓ |
| Automatic plotting of output variables for comparison | | ✓ |

See Also

More About

- “Manual Fixed-Point Conversion in MATLAB”
- “Automated Fixed-Point Conversion in MATLAB”

Tips for Making Generated Code More Efficient

| In this section... |
|--|
| “fimath Settings for Efficient Code” on page 49-9 |
| “Replace Functions With More Efficient Fixed-Point Implementations” on page 49-9 |

fimath Settings for Efficient Code

The default settings of the `fimath` object offer the smallest rounding error and prevent overflows. However, they can result in extra logic in generated code. These default settings are:

- `RoundingMethod`: `Nearest`
- `OverflowAction`: `Saturate`
- `ProductMode`: `FullPrecision`
- `SumMode`: `FullPrecision`

For leaner code, it is recommended that you match the `fimath` settings to the settings of your processor.

- The `KeepLSB` setting for `ProductMode` and `SumMode` models the behavior of integer operations in the C language. `KeepMSB` for `ProductMode` models the behavior of many DSP devices.
- Different rounding methods require different amounts of overhead code. Setting the `RoundingMethod` property to `Floor`, which is equivalent to two's complement truncation, provides the most efficient rounding implementation for most operations. For the divide function, the most efficient `RoundingMethod` is `Zero`.
- The standard method for handling overflows is to wrap using modulo arithmetic. Other overflow handling methods create costly logic. It is recommended that you set the `OverflowAction` property to `Wrap` when possible.

Replace Functions With More Efficient Fixed-Point Implementations

CORDIC

The CORDIC-based algorithms are among the most hardware friendly because they require only iterative shift-add operations. Replacing functions with one of the CORDIC implementations can make your generated code more efficient. For a list of the CORDIC functions, and examples of them being implemented, see “CORDIC Algorithms”.

Lookup tables

You can implement some functions more efficiently by using a lookup table approach. For an example, see “Implement Fixed-Point Log2 Using Lookup Table” on page 54-100.

Division

Division is often not supported by hardware. When possible, it is best to avoid division operations.

When the denominator is a power of two, you can rewrite the division as a bit shift operation.

$x/8$

can be rewritten as

`bitsra(x,3)`

Other times it is more efficient to implement division as a multiplication by a reciprocal.

`x/5`

can be rewritten as

`x*0.2`

Know When a Function is Supported for Instrumentation and Acceleration

There are several steps you can take to identify the features which could result in errors during conversion.

- `%#codegen` and `coder.screener`

Add the `%#codegen` pragma to the top of the MATLAB file that is being converted to fixed point. Adding this directive instructs the MATLAB Code Analyzer to help you diagnose and fix violations that would result in errors during when you try to accelerate or instrument your code.

The `coder.screener` function takes your function as its input argument and warns you of anything in your code that is not supported for codegen. Codegen support is essential for minimum and maximum logging and data type proposals.

- Consult the table of supported functions

See “Language Support” for a table of features supported for code generation and fixed-point conversion.

See Also

More About

- “Resolve Error: Function is not Supported for Fixed-Point Conversion” on page 49-12
- “Compilation Directive `%#codegen`” on page 14-5

Resolve Error: Function is not Supported for Fixed-Point Conversion

Issue

Some functions are not supported for fixed-point conversion and could result in errors during conversion.

Possible Solutions

Isolate the Unsupported Functions

When you encounter a function that is not supported for conversion, you can temporarily leave that part of the algorithm in floating point.

The following code returns an error because the `log` function is not supported for fixed-point inputs.

```
x = fi(rand(3), 1, 16, 15);
y = log(x)
```

Cast the input, `x`, to a double, and then cast the output back to a fixed-point data type.

```
y = fi(log(double(x)), 1, 16)

y =

    -0.2050    -0.0906    -1.2783
    -0.0990    -0.4583    -0.6035
    -2.0637    -2.3275    -0.0435
```

```
        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 16
        FractionLength: 13
```

This casting allows you to continue with your conversion until you can find a replacement.

Create a Replacement Function

You can replace the unsupported function with an alternative that is supported for fixed-point conversion.

- **Lookup Table Approximation** — You can replace many functions that are not supported for fixed-point conversion with a lookup table. For an example, see “Implement Fixed-Point Log2 Using Lookup Table” on page 54-100.
- **Polynomial Approximation** — You can approximate the results of a function that is not supported for fixed-point with a polynomial approximation.
- **User-Authored Function** — You can write your own function that supports fixed-point inputs. For example, using the `mod` function, which does support fixed-point inputs, you can write your own version of the `rem` function, which does not support fixed-point inputs.

See Also

More About

- “Know When a Function is Supported for Instrumentation and Acceleration” on page 49-11

Resolve Error: fi*non-fi

Issue

When multiplying a fixed-point variable by a non-fixed-point variable, the variable that does not have a fixed-point type can only be a constant

Possible Solutions

Before instrumenting your code, cast the non-fi variable to an acceptable fixed-point type.

| Original Algorithm | New Algorithm |
|--|---|
| <pre>function y = myProduct(x) y = 1; for n = 1:length(x) y(:) = y*x(n); end end</pre> | <pre>function y = myProduct(x) y = ones(1,1, 'like', x(1)*x(1)); for n = 1:length(x) y(:) = y*x(n); end end</pre> |

See Also

Functions

cast | ones | zeros

Resolve Error: Data Type Mismatch

Issue

In this example, `y` uses the default `fimath` setting `FullPrecision` for the `SumMode` property. At each iteration of the for-loop in the function `mysum`, the word length of `y` grows by one bit.

During simulation in MATLAB, there is no issue because data types can easily change in MATLAB. However, a data type mismatch error occurs at build time because data types must remain static in C.

Possible Solutions

Rewrite the function to use subscripted assignment within the for-loop.

In this example, rewrite `y = y + x(n)` as `y(:) = y + x(n)`, so that the value on the right is assigned in to the data type of `y`. This assignment preserves the `numericType` of `y` and avoids the type mismatch error.

| Original Algorithm | New Algorithm |
|---|---|
| <p><i>Function:</i></p> <pre>function y = mysum(x,T) %#codegen y = zeros(size(x), 'like', T.y); for n = 1:length(x) y = y + x(n); end end</pre> | <p><i>Function:</i></p> <pre>function y = mysum(x,T) %#codegen y = zeros(size(x), 'like', T.y); for n = 1:length(x) y(:) = y + x(n); end end</pre> |
| <p><i>Types Table:</i></p> <pre>function T = mytypes(dt) switch(dt) case 'fixed' F = fimath('RoundingMethod', 'Floor') T.x = fi([],1,16,11, F); T.y = fi([],1,16,6, F); end end end</pre> | <p><i>Types Table:</i></p> <pre>function T = mytypes(dt) switch(dt) case 'fixed' F = fimath('RoundingMethod', 'Floor') T.x = fi([],1,16,11, F); T.y = fi([],1,16,6, F); end end end</pre> |

See Also

`subsasgn`

More About

- “Compilation Directive `%#codegen`” on page 14-5

Resolve Error: Mismatched fimath

Issue

If two `fi` object operands have an attached `fimath`, the `fimaths` must be equal.

Possible Solutions

Use the `removefimath` function to remove the `fimath` of one of the variables in just one instance. By removing the `fimath`, you avoid the “mismatched `fimath`” error without permanently changing the `fimath` of the variable.

| Original Algorithm | New Algorithm |
|--|--|
| <p><i>Function:</i></p> <pre>function y = mysum(x,T) %#codegen y = zeros(size(x), 'like', T.y); for n = 1:length(x) y(:) = y + x(n); end end</pre> | <p><i>Function:</i></p> <pre>function y = mysum(x,T) %#codegen y = zeros(size(x), 'like', T.y); for n = 1:length(x) y(:) = removefimath(y) + x(n); end end</pre> |
| <p><i>Types Table:</i></p> <pre>function T = mytypes(dt) switch(dt) case 'fixed' T.x = fi([],1,16,0, 'RoundingMethod', 'Floor', 'OverflowAction', 'Wrap'); T.y = fi([],1,16,0, 'RoundingMethod', 'Nearest'); end end end</pre> | <p><i>Types Table:</i></p> <pre>function T = mytypes(dt) switch(dt) case 'fixed' T.x = fi([],1,16,0, 'RoundingMethod', 'Floor', 'OverflowAction', 'Wrap'); T.y = fi([],1,16,0, 'RoundingMethod', 'Nearest'); end end end</pre> |

See Also

`removefimath` | `fimath` | `fi`

Why Does the Fixed-Point Converter App Not Propose Data Types for System Objects?

The Fixed-Point Converter app might not display simulation range data or data type proposals for a System object because:

- The app displays range information for a subset of DSP System Toolbox System objects only. For a list of supported System objects, see [Converting System Objects to Fixed-Point Using the Fixed-Point Converter App](#) on page 7-73.
- The System object is not configured to use custom fixed-point settings.

If the system object is not configured correctly, the proposed data type column appears dimmed and displays `Full precision` or `Same as...` to show the current property setting.

See Also

Related Examples

- [“Convert dsp.FIRFilter Object to Fixed-Point Using the Fixed-Point Converter App”](#) on page 7-74

Slow Operations in the Fixed-Point Converter App

By default, the Fixed-Point Converter app screens your entry-point functions for code generation readiness. For some large entry-point functions, or functions with many calls, screening can take a long time. If the screening takes a long time, certain app or MATLAB operations can be slower than expected or appear to be unresponsive.

To determine if slow operations are due to the code generation readiness screening, disable the screening.

Blocks That Do Not Support Fixed-Point Data Types

Issue

Some blocks do not support fixed-point data types and can result in an error during fixed-point conversion.

The Simulink Block Data Type Support table summarizes characteristics of blocks in the Simulink block library, including whether or not they support fixed-point data types. To view the table, at the MATLAB command line, enter:

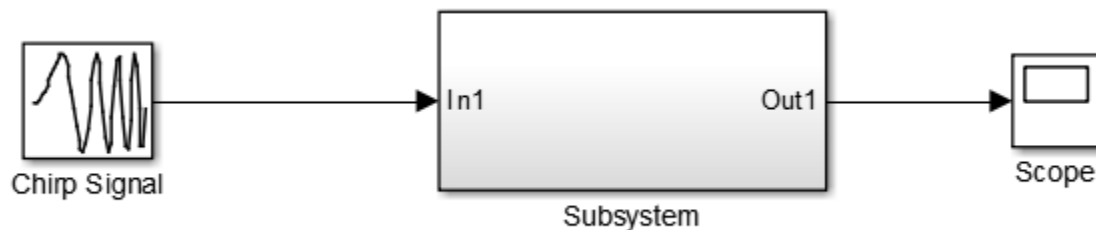
```
showblockdatatypetable
```

Possible Solutions

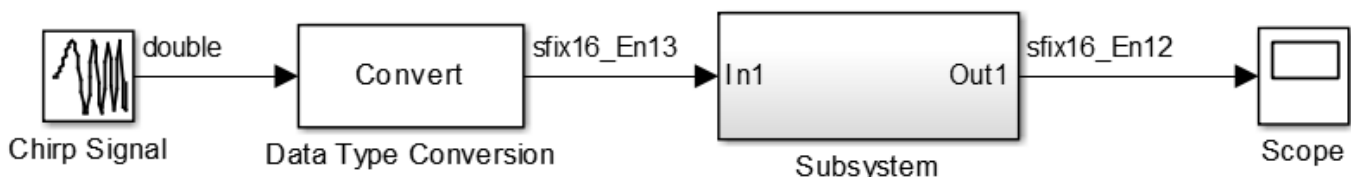
Isolate the Unsupported Block

If you encounter a block that is not supported for fixed-point conversion, you can isolate the unsupported block by decoupling it with a Data Type Conversion block. This workaround is useful when you do not intend to use the unsupported block on an embedded processor.

One example of this is using the Chirp Signal block, which does not support fixed-point outputs, to generate a signal for simulation data.



The subsystem shown is designed for use on an embedded processor and must be converted to fixed point. The Chirp Signal block creates simulation data. The Chirp Signal block supports only floating-point double outputs. However, if you decouple the Chirp Signal from the rest of the model by inserting a data type conversion block after the Chirp Signal block, you can use the Fixed-Point Tool to continue converting the subsystem to fixed point.



Isolate Unsupported Blocks with the Fixed-Point Tool

During the preparation stage of conversion, the Fixed-Point Tool identifies any blocks or constructs in your system under design that do not support fixed-point types. When the system under design

contains unsupported constructs, the Fixed-Point Tool encapsulates any unsupported elements in a subsystem containing the unsupported block surrounded by Data Type Conversion blocks. For more information, see “Use the Fixed-Point Tool to Prepare a System for Conversion” on page 38-2.

Isolate Unsupported Blocks with fxpopt

By default, fxpopt isolates any unsupported blocks by encapsulating the unsupported block in a subsystem surrounded by Data Type Conversion blocks. Isolated blocks are ignored by the optimizer.

Note When fxpopt isolates unsupported blocks, iterations of the optimization cannot be run in parallel. To run iterations of the optimization in parallel, use the Fixed-Point Tool to isolate unsupported blocks in the prepare stage, save the model, then run fxpopt with 'UseParallel' enabled.

Lookup Table Block Implementation

Many blocks that are not supported by the Fixed-Point Tool can be approximated with a lookup table block. Design an efficient fixed-point implementation of an unsupported block by using the **Lookup Table Optimizer**. For an example, see “Convert Floating-Point Model to Fixed Point” on page 40-2.

User-Authored Blocks

You can create your own block which is supported by the Fixed-Point Tool from one of the blocks in the User-Defined Functions Library.

See Also

Data Type Conversion

More About

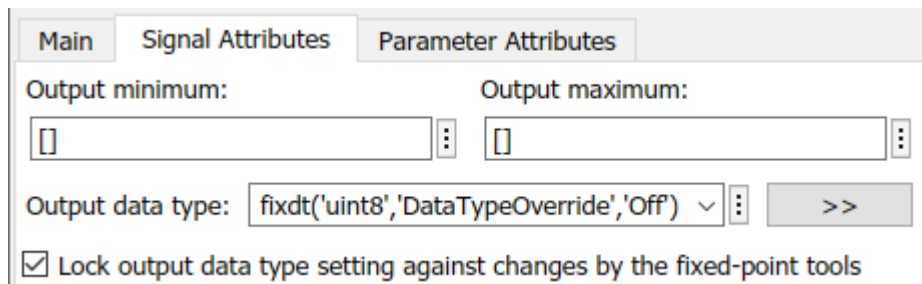
- “Approximate Functions with Lookup Tables”
- “User-Defined MATLAB Functions” (HDL Coder)

Prevent the Fixed-Point Tool from Overriding Integer Data Types

When performing data type override (DTO) on a selected system, the Fixed-Point Tool overrides the output data types of each block in the system. The only blocks that are never affected by DTO are blocks with `boolean` or enumerated output data types, and blocks that are untouched by DTO by design (for example, lookup table blocks). Depending on your application, you might want to preserve the data type of certain signals, for example, blocks that represent indices.

To prevent the Fixed-Point Tool from overriding the data type of a specific block, set the `DataTypeOverride` setting of the numeric type of the block to `Off`.

- 1 Open the Block Parameters dialog box by double-clicking the block.
- 2 Under the **Signal Attributes** tab, in the **Output data type** field, specify the desired data type and set the `et DataTypeOverride` property to `Off`.



You can set this override to off at the command line by changing the Data Type Override setting of a signal's `numericType`. In this example, the output data type of this block remains a built-in `uint8` even after performing data type override.

Alternatively, you can prevent the Fixed-Point Tool from replacing the current data type by using the **Lock output data type setting against changes by the fixed-point tools** parameter that is available on many blocks.

See Also

More About

- “Best Practices for Fixed-Point Conversion Workflow” on page 42-5

The Fixed-Point Tool did not Propose Data Types

Issue

In some cases, the Fixed-Point Tool does not propose data types for the system under design, or for one or more blocks within the system.

Possible Solutions

Inadequate Range Information

The Fixed-Point Tool bases its data type proposition on range information collected through simulation, derivation, and design ranges that you provide. Before proposing data types, you must collect range information which the Fixed-Point Tool uses to propose data types.

To collect range information, in the Fixed-Point Tool, select the desired **Range Collection Mode** in the **Setup** pane, and then click **Collect Ranges**.

Inherited Output Data Types

Blocks with inherited output data types use internal block rules to determine the output data type of the block. To enable proposals for results that specify an inherited output data type, in the Fixed-Point Tool, under **Settings**, set the **Convert inherited types** setting to Yes.

If this setting is set to No, the Fixed-Point Tool marks the proposal for these blocks as N/A.

See Also

More About

- “Collect Ranges” on page 42-20

Fraction Lengths and Fixed-Point Numbers

In this section...

“Fraction Length Greater Than Word Length” on page 49-23

“Negative Fraction Length” on page 49-23

Fraction Length Greater Than Word Length

A fraction length greater than the word length of a fixed-point number occurs when the number has an absolute value less than one and contains leading zeros.

In the following example, the fixed-point tool proposed a data type with a fraction length that is four greater than the word length. A binary representation of this number consists of the binary point, four implied leading zeros, followed by the binary representation of the stored integer: . X X X X 0 1 1 0 0 0 0 0, where the x’s represent the implied zeros.

| Name | Run | CompiledDT | Accept | ProposedDT | SimMin | SimMax |
|------|-------|------------|-------------------------------------|---------------|--------|--------|
| Gain | Run 1 | double | <input checked="" type="checkbox"/> | fixdt(1,8,12) | 0.0234 | 0.0234 |

Negative Fraction Length

A negative fraction length occurs when the number contains trailing zeros before the decimal point.

In the following example, the fixed-point tool proposed a data type with a negative fraction length. A binary representation of this number consists of the binary representation of the stored integer, followed by seven implied zeros, and then the binary point: 0 1 1 1 1 1 0 1 X X X X X X X ., where the x’s represent the implied zeros.

| Name | Run | CompiledDT | Accept | ProposedDT | SimMin | SimMax |
|------|-------|------------|-------------------------------------|---------------|--------|--------|
| Gain | Run 1 | double | <input checked="" type="checkbox"/> | fixdt(1,8,-7) | 16000 | 16000 |
| Out1 | Run 1 | | <input type="checkbox"/> | n/a | | |

See Also

More About

- “Fraction Length Greater Than Word Length” on page 49-5
- “Negative Fraction Length” on page 49-3

Why am I missing data type proposals for MATLAB Function block variables?

Fixed-Point Tool will not propose data types for variables in code inside a MATLAB Function block that is not executed during simulation. If your MATLAB Function block contains dead code, the variables will not appear in the Fixed-Point Tool.

- Update your input source so that all sections of your code are executed during simulation.
- This section of code may not be necessary. Delete the portion of code that is not exercised during simulation.

Data Type Propagation Errors After Applying Proposed Data Types


Under certain conditions, the Fixed-Point Tool may propose a data type that is not compatible with the model. The following topic outlines model configurations that may cause this issue, and how you can resolve the issue.

Tip Before attempting to autoscale a model, always ensure that you can update diagram successfully without data type override turned on.

Shared Data Type Groups

View Shared Data Type Groups

Organizing Fixed-Point Tool results into groups that must share the same data type can aid in the debugging process.

To view the data type group that a result belongs to, add the **DTGroup** column to the spreadsheet. Click the add column button . Select **DTGroup** in the menu.

Click the **DTGroup** column header to sort the results by this column.

Locked data type in a shared group

When an object is locked against changes by the Fixed-Point Tool, the Fixed-Point Tool does not propose a new data type for the object. If one of the results in a group of results that must share the same data type is locked, the Fixed-Point Tool proposes data types for all other objects in the group except for the locked object. If the data type proposed for the group is not compatible with the locked data type, a propagation error results.

To avoid incompatible data type proposals, perform one of the following.

- Lock all objects in the group against changes by the Fixed-Point Tool.
- Unlock the object in the group with the locked data type.

The **ProposedDT** column of the Fixed-Point Tool displays **Locked** for all results that are locked against changes by the Fixed-Point Tool.

Part of a Shared Data Type Group is Out of Scope

When results that are in a shared data type group share a data type from outside the scope of the system under design, the Fixed-Point Tool is not able to propose a data type.

To get a data type proposal, perform one of the following.

- Ensure that objects inside the system under design do not share their output data type with an object outside the selected system. One way to ensure that objects inside your system under design do not share their data type with objects outside the system, is by inserting Data Type Conversion blocks at the system boundaries.

- Ensure that all objects that must share a data type are within the scope of the system under design.

Model Reference Blocks

Systems that share data types across model reference boundaries may get data type propagation errors.

To avoid data type propagation errors, consider the following.

- Do not use the same signal object across model reference boundaries.
- Insert Data Type Conversion blocks at model reference boundaries.

Block Constraints

Certain blocks have constraints on which data types it can support. For example, the Merge block requires that all inputs use the same data type.

- Certain blocks in the Communications Toolbox, DSP System Toolbox, and Computer Vision Toolbox libraries have data type constraints. The Fixed-Point Tool is not aware of this requirement and does not use it for automatic data typing. Therefore, the tool might propose a data type that is inconsistent with the block requirements. In this case, manually edit the proposed data type such that it complies with block constraints.

Visit the individual block reference pages for more information on these constraints.

Internal Block Rules

Sum Blocks

Sum blocks have both an output data type as well as an accumulator data type. Under certain conditions, when the accumulator data type is set to `Inherit: Inherit via internal rule`, a data type propagation error can result.

To get a compatible data type proposal, perform one of the following.

- Change the accumulator data type to something other than `Inherit: Inherit via internal rule` and repropose data types for your model to get compatible data type proposals.
- Lock the block against changes by the fixed-point tools.

See Also

More About

- “Models That Might Cause Data Type Propagation Errors” on page 42-7

Resolve Range Analysis Issues

Issue

Different types of range analysis issues can occur in the Fixed-Point Tool depending on the specifics of your model.

Possible Solutions

Replace Unsupported Blocks

If the model contains blocks that are not supported for fixed-point conversion, range analysis will fail and the Fixed-Point Tool generates an error. Review the error message information and replace the unsupported blocks. For more information, see “Model Compatibility with Range Analysis” on page 43-4.

Fix Design Range Conflicts

If the model contains conflicting design range information, the analysis cannot derive range data and the Fixed-Point Tool generates an error. Examine the design ranges specified in the model to identify inconsistent design specifications and modify them to be consistent. For more information, see “Fix Design Range Conflicts” on page 43-22.

Specify Additional Design Range Information

If there is insufficient design range information specified, the analysis cannot derive range data and the Fixed-Point Tool highlights the results. Examine the model to determine which design range information is missing. Specify additional range information or reconfigure the code so that the Fixed-Point Tool can derive ranges for the model. For more information, see “Insufficient Design Range Information” on page 43-14 and “Troubleshoot Range Analysis of System Objects” on page 43-16.

See Also

More About

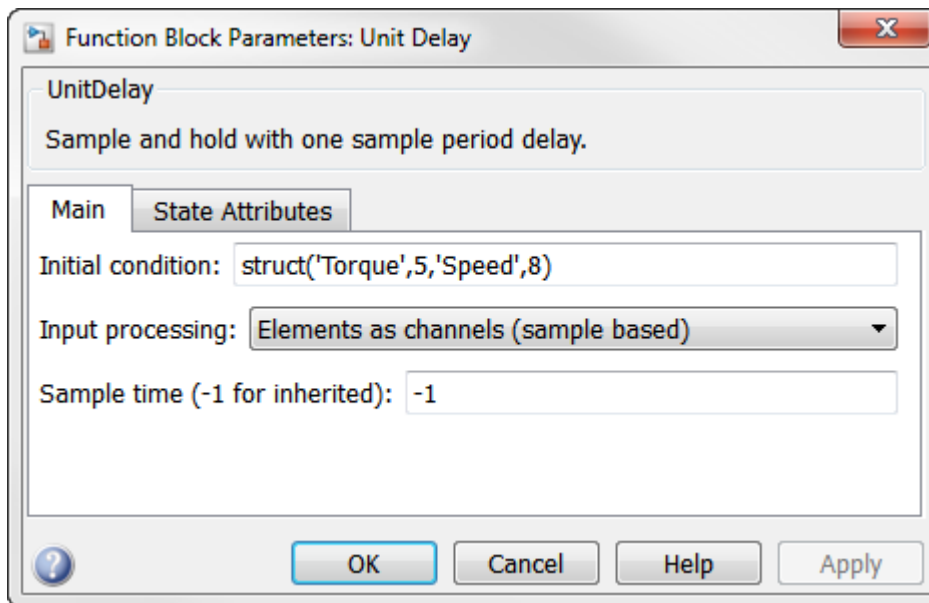
- “How Range Analysis Works” on page 43-2
- “Blocks That Do Not Support Fixed-Point Data Types” on page 49-19

Data Type Mismatch and Structure Initial Conditions

Specify Bus Signal Initial Conditions Using Simulink.Parameter Objects

This example shows how to replace a structure initial condition with a `Simulink.Parameter` object. This approach allows the structure to maintain its tunability.

- 1 Double-click the Unit Delay block to view the block parameters. The Unit Delay block uses a structure initial condition.



- 2 Define a `Simulink.Parameter` object at the MATLAB command line. Set the data type of the parameter object to the bus object `SensorData`. Set the value of the parameter object to the specified structure. To maintain tunability, set the `StorageClass` property to `ExportedGlobal`.

```
P = Simulink.Parameter;
P.DataType = 'Bus: SensorData';
P.Value = struct('Torque',5,'Speed',8);
P.StorageClass = 'ExportedGlobal';
```

- 3 In the Unit Delay block dialog box, set **Initial condition** to `P`, the `Simulink.Parameter` object you defined. The structure defined in the `Simulink.Parameter` object remains tunable.

For more information on generating code for bus signals that use tunable initial condition structures, see “Control Signal and State Initialization in the Generated Code”.

Data Type Mismatch and Masked Atomic Subsystems

A data type mismatch occurs when a structure initial condition drives a bus signal that you specified using a masked atomic subsystem.

Change the subsystem to non atomic, or specify the structure parameter using a `Simulink.Parameter` object (as described in “Specify Bus Signal Initial Conditions Using Simulink.Parameter Objects” on page 49-28) to avoid the data type mismatch error.

See Also

Related Examples

- “Bus Objects in the Fixed-Point Workflow” on page 46-2

Reconversion Using the Fixed-Point Tool

After you simulate your model using embedded types and compare the floating point and fixed-point behavior of your system, determine if the new behavior is satisfactory. If the behavior of the system using the newly applied fixed-point data types is not acceptable, you can iterate through the process of editing your proposal settings, proposing and applying data types, and comparing the results until you find settings that work for your system. You do not need to perform the range collection step again.

See Also

More About

- “Iterative Fixed-Point Conversion Using the Fixed-Point Tool” on page 42-9
- “Explore Multiple Floating-Point to Fixed-Point Conversions” on page 40-11

Data Type Optimization Not Successful

Issue

You can use the `fxpopt` function or the Optimized Fixed-Point Conversion workflow in the Fixed-Point Tool to optimize the data types of a model or subsystem. Sometimes, the optimization is not successful. The following sections describe how to troubleshoot these cases.

Possible Solutions

Unable to Model Problem — No Constraints Specified

To determine if the behavior of a new fixed-point implementation is acceptable, the optimization requires well-defined behavioral constraints. Use the `addTolerance` method of the `fxpOptimizationOptions` class to specify numerical constraints for the optimized design. Alternatively, use blocks from the Model Verification library. For more information, see “Specify Behavioral Constraints” on page 42-18.

Unable to Model Problem — Model is not Supported

The model containing the system you want to optimize must have the following characteristics:

- All blocks in the model must support fixed-point data types.
- The design ranges specified on blocks in the model must be consistent with the simulation ranges.
- If the model contains a MATLAB Function block, it must use MATLAB language features supported for fixed-point conversion. For more information, see “MATLAB Language Features Supported for Automated Fixed-Point Conversion” on page 7-35.
- The data logging format of the model must be set to `Dataset`.

To configure this setting, in the Configuration Parameters, in the **Data Import/Export** pane, set **Format** to `Dataset`.

- The model must have finite simulation stop time.

Data Type Conversion Blocks Were Ignored by Optimization

When the **Input and output to have equal** parameter of a Data Type Conversion block is set to `Stored Integer (SI)`, the Data Type Conversion block will be ignored by the optimization.

Unable to Find a Fixed-Point Implementation that Met the Tolerances

If the optimization cannot find a feasible solution, try these solutions:

- Relax signal tolerances.
- Allow larger word lengths to expand the search space.
- Consider using time windows when specifying signal tolerances. For more information, see “Tolerance Computation”.
- Instead of specifying low-level tolerances on individual signals, consider specifying high-level behavioral constraints using blocks from the “Model Verification” library. For more information, see “Specify Behavioral Constraints” on page 42-18.

Unable to Explore Results

When the optimization is not able to find a new valid result, the `fxpopt` function does not produce an `OptimizationResult` output. Invalid results are most often the result of using a model that is not supported for optimization. For more information, see “Unable to Model Problem — No Constraints Specified Unable to Model Problem — Model is not Supported” on page 49-31.

When the optimization is successful, you can explore several different implementations of your design that were found during the optimization process. Do not save the model until you are satisfied with the new design. Saving the model disables you from continuing to explore the other implementations.

Resolve Error: The RowNames property must be a string array or a cell array, with each name containing one or more characters

This error may occur if `clear all` is used during fixed-point conversion workflows in the Fixed-Point Tool. `clear all` is currently not supported by fixed-point conversion workflows. Do not use `clear all` in initialization functions (`InitFcn`), or at the MATLAB Command Window when using the Fixed-Point Tool.

Derived Range Analysis Does Not Work for Accumulator Data Type

Only block output signals participate in derived range analysis. If a block has additional data type controls, such as for the accumulator or intermediate results, ranges are not derived for these elements. As a result, when optimization considers both simulation ranges and derived ranges, only simulation range information is used to optimize accumulator data types. Therefore, the optimized accumulator data type and output data type for a given block may differ. For more information, see “How Range Analysis Works” on page 43-2.

See Also

Classes

`fxpOptimizationOptions` | `OptimizationResult` | `OptimizationSolution`

Functions

`addTolerance` | `showTolerances` | `explore` | `fxpopt`

More About

- “Optimize Fixed-Point Data Types for a System” on page 40-14

Compile-Time Recursion Limit Reached

Issue

You see a message such as:

```
Compile-time recursion limit reached. Size or type of
input #1 of function 'foo' may change at every call.
```

```
Compile-time recursion limit reached. Value of input #1
of function 'foo' may change at every call.
```

Cause

With compile-time recursion, the code generator produces multiple versions of the recursive function instead of producing a recursive function in the generated code. These versions are known as function specializations. The code generator is unable to use compile-time recursion for a recursive function in your MATLAB code because the number of function specializations exceeds the limit.

Solutions

To address the issue, try one of these solutions:

- “Force Run-Time Recursion” on page 49-33
- “Increase the Compile-Time Recursion Limit” on page 49-35

Force Run-Time Recursion

- For this message:

```
Compile-time recursion limit reached. Value of input #1
of function 'foo' may change at every call.
```

Use this solution:

“Force Run-Time Recursion by Treating the Input Value as Nonconstant” on page 49-33.

- For this message:

```
Compile-time recursion limit reached. Size or type of
input #1 of function 'foo' may change at every call.
```

In the code generation report, look at the function specializations. If you can see that the size of an argument is changing for each function specialization, then try this solution:

“Force Run-Time Recursion by Making the Input Variable-Size” on page 49-34.

Force Run-Time Recursion by Treating the Input Value as Nonconstant

Consider this function:

```
function y = call_recfcn(n)
A = ones(1,n);
```

```

x = 100;
y = recfcn(A,x);
end

function y = recfcn(A,x)
if size(A,2) == 1 || x == 1
    y = A(1);
else
    y = A(1)+recfcn(A(2:end),x-1);
end
end

```

The second input to `recfcn` has the constant value 100. The code generator determines that the number of recursive calls is finite and tries to produce 100 copies of `recfcn`. This number of specializations exceeds the compile-time recursion limit. To force run-time recursion, instruct the code generator to treat the second input as a nonconstant value by using `coder.ignoreConst`.

```

function y = call_recfcn(n)
A = ones(1,n);
x = coder.ignoreConst(100);
y = recfcn(A,x);
end

function y = recfcn(A,x)
if size(A,2) == 1 || x == 1
    y = A(1);
else
    y = A(1)+recfcn(A(2:end),x-1);
end
end

```

If the code generator cannot determine that the number of recursive calls is finite, it produces a run-time recursive function.

Force Run-Time Recursion by Making the Input Variable-Size

Consider this function:

```

function z = call_mysum(A)
%#codegen
z = mysum(A);
end

function y = mysum(A)
coder.inline('never');
if size(A,2) == 1
    y = A(1);
else
    y = A(1)+mysum(A(2:end));
end
end

```

If the input to `mysum` is fixed-size, the code generator uses compile-time recursion. If `A` is large enough, the number of function specializations exceeds the compile-time limit. To cause the code generator to use run-time conversion, make the input to `mysum` variable-size by using `coder. varsizes`.


```

function z = call_mysum(A)
    %#codegen
    B = A;
    coder.varsize('B');
    z = mysum(B);
end

function y = mysum(A)
    coder.inline('never');
    if size(A,2) == 1
        y = A(1);
    else
        y = A(1)+ mysum(A(2:end));
    end
end

```

Increase the Compile-Time Recursion Limit

The default compile-time recursion limit of 50 is large enough for most recursive functions that require compile-time recursion. Usually, increasing the limit does not fix the issue. However, if you can determine the number of recursive calls and you want compile-time recursion, increase the limit. For example, consider this function:

```

function z = call_mysum()
    %#codegen
    B = 1:125;
    z = mysum(B);
end

function y = mysum(A)
    coder.inline('never');
    if size(A,2) == 1
        y = A(1);
    else
        y = A(1)+ mysum(A(2:end));
    end
end

```

You can determine that the code generator produces 125 copies of the `mysum` function. In this case, if you want compile-time recursion, increase the compile-time recursion limit to 125.

To increase the limit, in a code acceleration configuration object, increase the value of the `CompileTimeRecursionLimit` configuration parameter.

See Also

More About

- “Code Generation for Recursive Functions” on page 14-12
- “Set Up C Compiler and Compilation Options” on page 12-16

Output Variable Must Be Assigned Before Run-Time Recursive Call

Issue

You see this error message:

```
All outputs must be assigned before any run-time recursive call. Output 'y' is not assigned here.
```

Cause

Run-time recursion produces a recursive function in the generated code. The code generator is unable to use run-time recursion for a recursive function in your MATLAB code because an output is not assigned before the first recursive call.

Solution

Rewrite the code so that it assigns the output before the recursive call.

Direct Recursion Example

In the following code, the statement `y = A(1)` assigns a value to the output `y`. This statement occurs after the recursive call `y = A(1)+ mysum(A(2:end))`.

```
function z = call_mysum(A)
B = A;
coder.varsize('B');
z = mysum(B);
end

function y = mysum(A)
coder.inline('never');
if size(A,2) > 1
    y = A(1)+ mysum(A(2:end));

else
    y = A(1);
end
end
```

Rewrite the code so that assignment `y = A(1)` occurs in the `if` block and the recursive call occurs in the `else` block.

```
function z = call_mysum(A)
B = A;
coder.varsize('B');
z = mysum(B);
end

function y = mysum(A)
coder.inline('never');
```

```

if size(A,2) == 1
    y = A(1);
else
    y = A(1)+ mysum(A(2:end));
end
end

```

Alternatively, before the if block, add an assignment, for example, $y = 0$.

```

function z = call_mysum(A)
B = A;
coder.varsize('B');
z = mysum(B);
end

function y = mysum(A)
coder.inline('never');
y = 0;
if size(A,2) > 1
    y = A(1)+ mysum(A(2:end));

else
    y = A(1);
end
end

```

Indirect Recursion Example

In the following code, rec1 calls rec2 before the assignment $y = 0$.

```

function y = rec1(x)
%#codegen

if x >= 0
    y = rec2(x-1)+1;
else
    y = 0;
end
end

function y = rec2(x)
y = rec1(x-1)+2;
end

```

Rewrite this code so that in rec1, the assignment $y = 0$ occurs in the if block and the recursive call occurs in the else block.

```

function y = rec1(x)
%#codegen

if x < 0
    y = 0;
else
    y = rec2(x-1)+1;
end
end

function y = rec2(x)

```

```
y = rec1(x-1)+2;  
end
```

See Also

More About

- “Code Generation for Recursive Functions” on page 14-12

Unable to Determine That Every Element of Cell Array Is Assigned

Issue

You see one of these messages:

Unable to determine that every element of 'y' is assigned before this line.

Unable to determine that every element of 'y' is assigned before exiting the function.

Unable to determine that every element of 'y' is assigned before exiting the recursively called function.

Cause

For code generation, before you use a cell array element, you must assign a value to it. When you use `cell` to create a variable-size cell array, for example, `cell(1,n)`, MATLAB assigns an empty matrix to each element. However, for code generation, the elements are unassigned. For code generation, after you use `cell` to create a variable-size cell array, you must assign all elements of the cell array before any use of the cell array.

The code generator analyzes your code to determine whether all elements are assigned before the first use of the cell array. The code generator detects that all elements are assigned when the code follows this pattern:

```
function z = CellVarSize1D(n, j)
%#codegen
x = cell(1,n);
for i = 1:n
    x{i} = i;
end
z = x{j};
end
```

Here is the pattern for a multidimensional cell array:

```
function z = CellAssign3D(m,n,p)
%#codegen
x = cell(m,n,p);
for i = 1:m
    for j = 1:n
        for k = 1:p
            x{i,j,k} = i+j+k;
        end
    end
end
z = x{m,n,p};
end
```

If the code generator detects that some elements are not assigned, code generation fails. Sometimes, even though your code assigns all elements of the cell array, code generation fails because the analysis does not detect that all elements are assigned.

Here are examples where the code generator is unable to detect that elements are assigned:

- Elements are assigned in different loops

```
...
x = cell(1,n)
for i = 1:5
    x{i} = 5;
end
for i = 6:n
    x{i} = 7;
end
...
```

- The variable that defines the loop end value is not the same as the variable that defines the cell dimension.

```
...
x = cell(1,n);
m = n;
for i = 1:m
    x{i} = 2;
end
...
```

For more information, see “Definition of Variable-Size Cell Array by Using cell” on page 30-8.

Solution

Try one of these solutions:

- “Use recognized pattern for assigning elements” on page 49-40
- “Use repmat” on page 49-40
- “Use coder.nullcopy” on page 49-41

Use recognized pattern for assigning elements

If possible, rewrite your code to follow this pattern:

```
...
x = cell(1,n);
for i = 1:n
    x{i} = i;
end
z = x{j};
...
```

Use repmat

Sometimes, you can use repmat to define the variable-size cell array.

Consider this code that defines a variable-size cell array. It assigns the value 1 to odd elements and the value 2 to even elements.

```
function z = repDefine(n, j)
%#codegen
c =cell(1,n);
for i = 1:2:n-1
    c{i} = 1;
end
for i = 2:2:n
    c{i} = 2;
end
z = c{j};
```

Code generation does not allow this code because:

- More than one loop assigns the elements.
- The loop counter does not increment by 1.

Rewrite the code to first use `cell` to create a 1-by-2 cell array whose first element is 1 and whose second element is 2. Then, use `repmat` to create a variable-size cell array whose element values alternate between 1 and 2.

```
function z = repVarSize(n, j)
%#codegen
c = cell(1,2);
c{1} = 1;
c{2} = 2;
c1= repmat(c,1,n);
z = c1{j};
end
```

You can pass an initially empty, variable-size cell array into or out of a function by using `repmat`. Use the following pattern:

```
function x = emptyVarSizeCellArray
x = repmat({'abc'},0,0);
coder.ysize('x');
end
```

This code indicates that `x` is an empty, variable-size cell array of 1x3 characters that can be passed into or out of functions.

Use `coder.nullcopy`

As a last resort, you can use `coder.nullcopy` to indicate that the code generator can allocate the memory for your cell array without initializing the memory. For example:

```
function z = nulcpyCell(n, j)
%#codegen
c =cell(1,n);
c1 = coder.nullcopy(c);
for i = 1:4
    c1{i} = 1;
end
for i = 5:n
    c1{i} = 2;
end
z = c1{j};
end
```

Use `coder.nullcopy` with caution. If you access uninitialized memory, results are unpredictable.

See Also

`cell` | `repmat` | `coder.nullcopy`

More About

- “Cell Array Limitations for Code Generation” on page 30-7

Nonconstant Index into varargin or varargout in a for-Loop

Issue

Your MATLAB code contains a for-loop that indexes into varargin or varargout. When you generate code, you see this error message:

```
Non-constant expression or empty matrix. This expression
must be constant because its value determines the size
or class of some expression.
```

Cause

At code generation time, the code generator must be able to determine the value of an index into varargin or varargout. When varargin or varargout are indexed in a for-loop, the code generator determines the index value for each loop iteration by unrolling the loop. Loop unrolling makes a copy of the loop body for each loop iteration. In each iteration, the code generator determines the value of the index from the loop counter.

The code generator is unable to determine the value of an index into varargin or varargout when:

- The number of copies of the loop body exceeds the limit for loop unrolling.
- Heuristics fail to identify that loop unrolling is warranted for a particular for-loop. For example, consider the following function:

```
function [x,y,z] = fcn(a,b,c)
    %#codegen

    [x,y,z] = subfcn(a,b,c);

    function varargout = subfcn(varargin)
        j = 0;
        for i = 1:nargin
            j = j+1;
            varargout{j} = varargin{j};
        end
```

The heuristics do not detect the relationship between the index *j* and the loop counter *i*. Therefore, the code generator does not unroll the for-loop.

Solution

Use one of these solutions:

- “Force Loop Unrolling” on page 49-43
- “Rewrite the Code” on page 49-44

Force Loop Unrolling

Force loop unrolling by using `coder.unroll`. For example:

```
function [x,y,z] = fcn(a,b,c)
    %#codegen
```

```
[x,y,z] = subfcn(a,b,c);  
  
function varargout = subfcn(varargin)  
j = 0;  
  
coder.unroll();  
for i = 1:nargin  
    j = j + 1;  
    varargout{j} = varargin{j};  
end
```

Rewrite the Code

Rewrite the code so that the code generator can detect the relationship between the index and the loop counter. For example:

```
function [x,y,z] = fcn(a,b,c)  
%#codegen  
[x,y,z] = subfcn(a,b,c);  
  
function varargout = subfcn(varargin)  
for i = 1:nargin  
    varargout{i} = varargin{i};  
end
```

See Also

`coder.unroll`

More About

- “Code Generation for Variable Length Argument Lists” on page 17-2

Single-Precision Conversion in Simulink

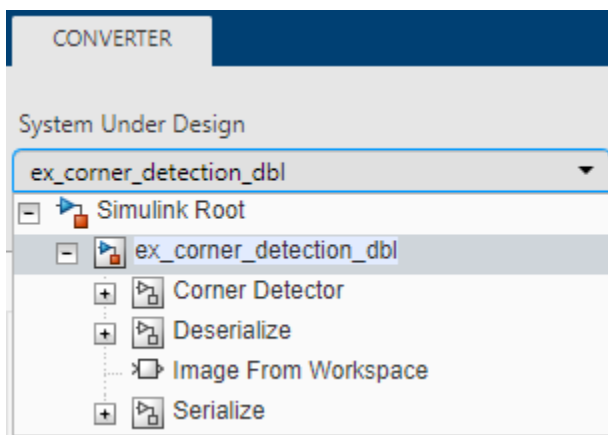
- “Getting Started with Single Precision Converter” on page 50-2
- “Convert a System to Single Precision” on page 50-5

Getting Started with Single Precision Converter

The Single Precision Converter converts your model or a system in your model from double precision to single precision. To open the Single Precision Converter, from the Simulink **Apps** tab, select **Single Precision Converter**.

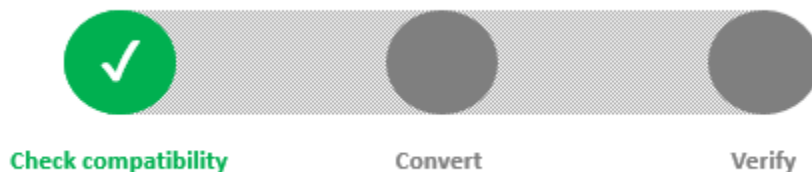
Select System Under Design

To begin, expand the **System Under Design** drop-down list and select the system to convert to single precision.



Check Compatibility

To start the conversion, click **Convert to Single**.



The Single Precision Converter performs these checks:

- Check that all blocks in the selected system support single precision.

The Single Precision Converter displays a list of blocks that do not support single precision or are locked against changes by the Fixed-Point Tools. To restart the conversion, replace the blocks that support only double precision and unlock the blocks that are locked against changes by the Fixed-Point Tools. Then click **Convert to Single**.

- Check that the system uses a library standard that supports single-precision designs.

To convert a system to single precision, the language standard must be set to C99 (ISO). If the specified language standard is not set to C99, the Single Precision Converter changes the math library.

- Check that the solver settings are set to fixed step.

Convert



Following the compatibility check, the Single Precision Converter converts the system to single-precision. The Converter makes these changes:

- Conversion of user-specified double-precision data types to single-precision data types (applies to block settings, Stateflow chart settings, signal objects, and bus objects).
- When the system under design contains a MATLAB Function block, the converter creates a variant subsystem containing a generated single-precision version of the MATLAB Function block and the original MATLAB Function block.
- Output signals and intermediate settings using inherited data types that compile to double-precision change to single-precision data types.

The converter does not change Boolean, built-in integer, or user-specified fixed-point data types. When the conversion is finished, the converter displays a table summarizing the compiled and proposed data types of the objects in the system under design.

Verify



Finally, the Single Precision Converter verifies that the model containing the converted system can successfully update the diagram. If the model is not able to update the diagram due to data type mismatch errors at the system boundaries, the Single Precision Converter displays a message.

To resolve the data type mismatch, insert Data Type Conversion blocks at the system boundaries. You can also resolve the data type mismatch errors by changing the output data type of the blocks feeding into the system to single or `Inherit: Inherit via back propagation`.

See Also

Related Examples

- “Convert a System to Single Precision” on page 50-5
- “Specify Single-Precision Data Type for Embedded Application”

Convert a System to Single Precision

In this section...

“Open Model” on page 50-5

“Convert to Single Precision” on page 50-5

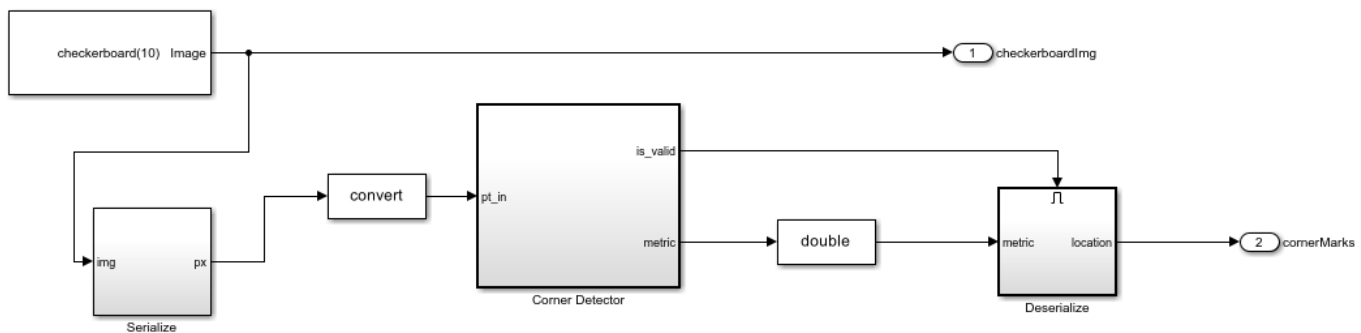
This example shows how to convert a system to single precision using the Single Precision Converter. This example converts a subsystem of a double-precision model to single precision. To convert a subsystem in a model to single precision, surround the subsystem under design with Data Type Conversion blocks before opening the Single Precision Converter.

Open Model

Open the `ex_corner_detection_double` model and set model parameters.

```
open_system("ex_corner_detection_double.slx")
```

```
R = 80; C = 80;  
g = fspecial('gaussian',[5 5],1.5);
```

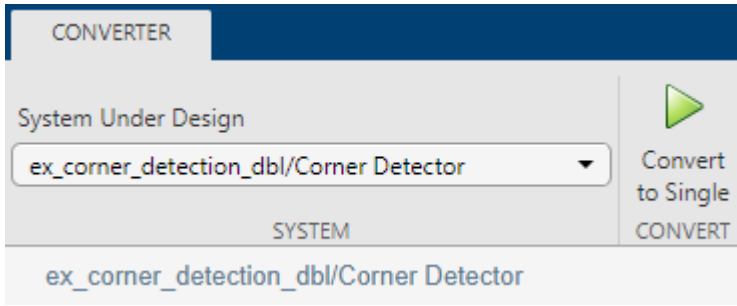


Copyright 2014-2017 The MathWorks, Inc.

The model uses a combination of double-precision, Boolean, and built-in integer data types.

Convert to Single Precision

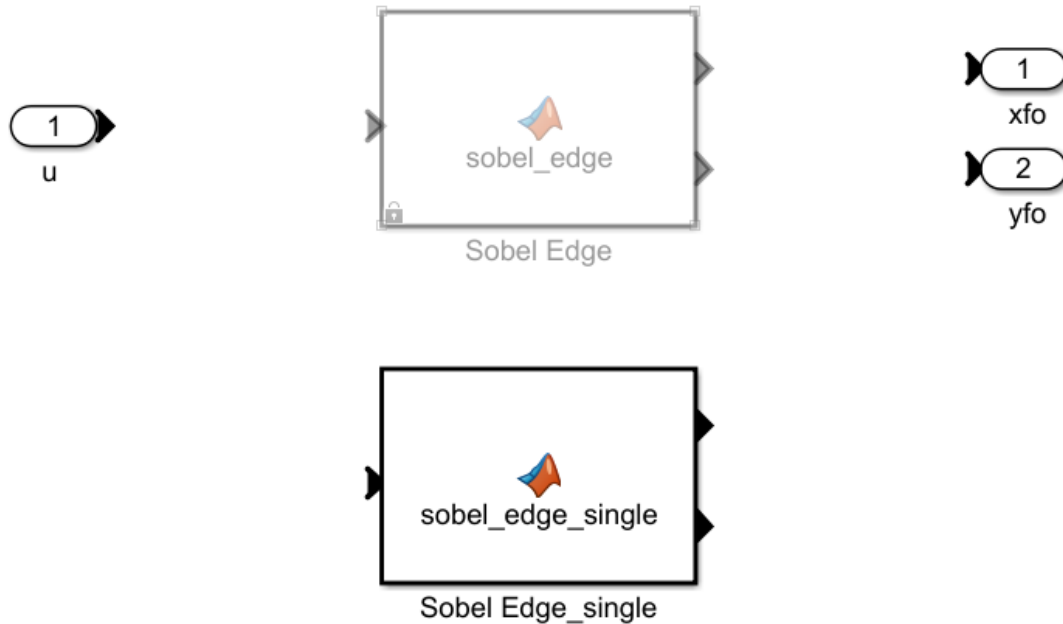
- 1 Open the Single Precision Converter. From the Simulink **Apps** tab, select **Single Precision Converter**.
- 2 Under **System Under Design**, select the system or subsystem to convert to single precision. For this example, select the **Corner Detector** subsystem. Click **Convert to Single**.



The converter first checks the system for compatibility with the conversion and changes any model settings that are incompatible. The language standard of the model must be set to C99 (ISO), and the model must use a fixed-step solver.

The converter converts the system and lists all converted data types. The converter changes only double-precision data types. It does not convert Boolean, fixed-point, or built-in integer types to single precision.

When the system under design contains a MATLAB Function block, the converter creates a variant subsystem containing a generated single-precision version of the MATLAB Function block and the original MATLAB Function block.



In the final stage of the conversion, the Converter verifies that the conversion was successful by updating the model.

Single Precision Converter

CONVERTER

System Under Design
 ex_corner_detection_db/Corner Detector

Convert to Single CONVERT

ex_corner_detection_db/Corner Detector

compatibility Convert Verify

Check for compatibility
 There are blocks that are locked against data type changes

- ex_corner_detection_db/Corner Detector/Corner Metric/Gx
- ex_corner_detection_db/Corner Detector/Corner Metric/Gy
- ex_corner_detection_db/Corner Detector/Corner Metric/Gaussian Filter/Gx
- ex_corner_detection_db/Corner Detector/Corner Metric/Gaussian Filter/Gy

Changes made to the model
 The standard math library of the following models was changed from c89 to c99

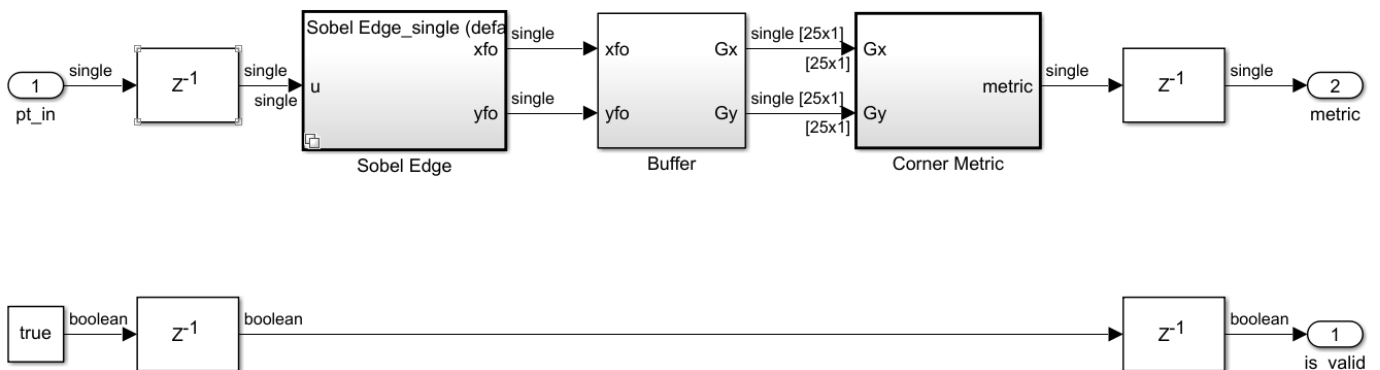
- ex_corner_detection_db

Convert system to single

| Block path | Compiled data type | Proposed data type |
|---|--------------------|--------------------|
| ex_corner_detection_db/Corner Detector/pt_in | double | single |
| ex_corner_detection_db/Corner Detector/Sobel Edge/u | double | single |
| ex_corner_detection_db/Corner Detector/Sobel Edge/Sobel Edge_sing | double | single |

Verify converted system
 The model updated successfully after conversion

3 Return to the model and update the diagram. The blocks inside the Corner Detector subsystem no longer use double-precision data types.



See Also

More About

- “Getting Started with Single Precision Converter” on page 50-2
- “Specify Single-Precision Data Type for Embedded Application”

Simulink Half Precision

- “The Half-Precision Data Type in Simulink” on page 51-2
- “Generate Native Half-Precision C Code from Simulink Models” on page 51-5
- “Half-Precision Field-Oriented Control Algorithm” on page 51-11
- “Image Quantization with Half-Precision Data Types” on page 51-14
- “Convert Single Precision Lookup Table to Half Precision” on page 51-15
- “Digit Classification with Half-Precision Data Types” on page 51-20

The Half-Precision Data Type in Simulink

Signals and block outputs can specify a half-precision data type. The half-precision data type is supported for simulation and code generation for parameters and a subset of blocks.

Math Operations in Half-Precision

In Simulink, half-precision inputs to blocks performing arithmetic operations, relational operations, and binary operations are always cast to single precision, and the operation is performed in single precision. If the output data type of the block is set to `half`, the output of the block is cast back to a half-precision data type.

In MATLAB, however, some functions perform arithmetic operations with half-precision inputs by emulating the half-precision floating-point math. For example, in MATLAB, the following code is performed using half-precision floating-point arithmetic.

```
y = mod(half(u1), half(u2))
```

In Simulink, using the `mod` function of the Math Function block, the same operation would be performed by casting the inputs to single precision and carrying out the operation in single-precision floating-point math. The result of the arithmetic operations is then cast back to half precision.

```
y = half(mod(single(half(u1)), single(half(u2))))
```

Software Features Supported for Half Precision

- The half-precision data type is supported for simulation in Normal, Accelerator, and Rapid Accelerator modes. The half-precision data type is also supported for SIL, PIL, and external modes.
- Half precision is supported for C/C++ code generation for `.ert` targets.

In the generated code, half-precision variables are stored in a class emulating the bit pattern of the value.

- For embedded hardware targets that natively support special types for half precision, native half-precision C code generation is supported. For more information, see “Generate Native Half-Precision C Code from Simulink Models” on page 51-5.
- HDL code generation using HDL Coder.

For more information, see “Getting Started with HDL Coder Native Floating-Point Support” (HDL Coder).

- MATLAB System block supports half-precision data type with real values.
- In Simulink, the half-precision data type only supports real values. Complex values cannot have a half-precision data type.

Blocks Supported for Half Precision

To view the blocks that support half precision, at the command line, type:

```
showblockdatatypetable
```

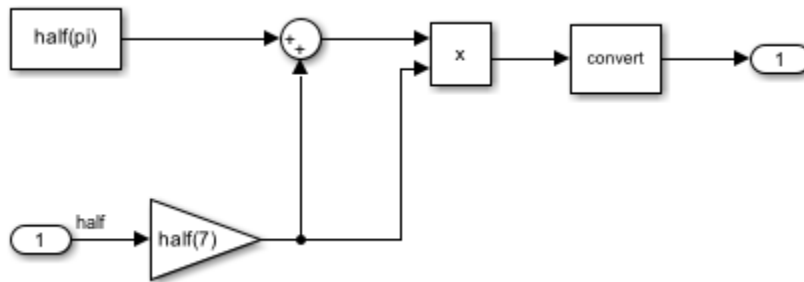
Blocks that support half precision display an X in the column labeled **Half**.

Generate Code for Half Precision Systems

You can generate C code targeting `.ert` targets for Simulink models using the half-precision data type. Code generation for `.ert` targets requires an Embedded Coder® license.

Open the `ex_half_arithmetic` model. The model performs several arithmetic operations. All parameter values and output data types specify a half-precision data type.

```
open_system('ex_half_arithmetic');
```



To generate C code for the model, press **Ctrl+B**. In the code generation report, open the `ex_half_arithmetic.c` file. Half-precision variables are types in the generated code as `real16_T`. For example, see the `rtb_Gain` variable.

```

33  /* Model step function */
34  void ex_half_arithmetic_step(void)
35  {
36  |   real16_T rtb_Gain;
37  |
  
```

In the generated code, half-precision variables are stored in a struct emulating the bit pattern of the value.

Half-precision input variables to arithmetic operations are cast to single precision, and the arithmetic operation is performed in single precision. If the output data type of the block is set to `half`, the result of the operation is cast back to a half-precision data type. For example, see the code computing the output of the Gain block.

```

38  /* Gain: '<Root>/Gain' incorporates:
39  *   Inport: '<Root>/Inport'
40  */
41  rtb_Gain = floatToHalf(7.0F * halfToFloat(ex_half_arithmetic_U.Inport));
42
  
```

See Also

`half`

More About

- “Floating-Point Numbers” on page 35-20
- “Generate Native Half-Precision C Code from Simulink Models” on page 51-5
- “Half-Precision Field-Oriented Control Algorithm” on page 51-11
- “Image Quantization with Half-Precision Data Types” on page 51-14
- “Digit Classification with Half-Precision Data Types” on page 51-20
- “Convert Single Precision Lookup Table to Half Precision” on page 51-15

Generate Native Half-Precision C Code from Simulink Models

Some embedded hardware targets natively support special types for half precision, such as `_Float16` and `_fp16` data types for ARM compilers. You can generate native half-precision C code for embedded hardware targets that natively support half precision floating-point data types. The process to generate native half C code is as follows:

- Register a new hardware target device that natively supports half precision using the `target` package.
- Set up the Simulink model for native half code generation with the new target hardware.
- Configure the Simulink model toolchain and set up the compiler for half precision.
- Generate native half type code.

Fixed-Point Designer includes preconfigured language implementations for Armclang and GCC compilers. For other hardware targets, you can specify a custom language implementation based on your hardware specifications.

Generate Native Half-Precision C Code for ARM Cortex-A with Armclang Compiler

In this example, an ARM Cortex-A processor is used as the hardware target. The model is configured to use this ARM target and the Armclang compiler toolchain. Refer to the ARM developer documentation for detailed information on using the half-precision data type on this hardware. See Arm Compiler armclang Reference Guide: Half-precision floating-point data types.

Register Target Hardware

Use the `target.create` function to create an ARM Cortex-A processor target that is compatible with half precision.

```
arm_half = target.create('Processor', ...
    'Manufacturer', "ARM Compatible", ...
    'Name', 'ARM Cortex-A Half-Precision');
```

Add the language implementation. The 32-bit version of Clang for ARM has half precision enabled by default. Use `target.get` to retrieve the target object from the internal database.

```
li = target.get('LanguageImplementation', "Clang ARM 32-bit");
```

Replace the default language implementation for ARM Cortex with Armclang.

```
arm_half.LanguageImplementations = li;
```

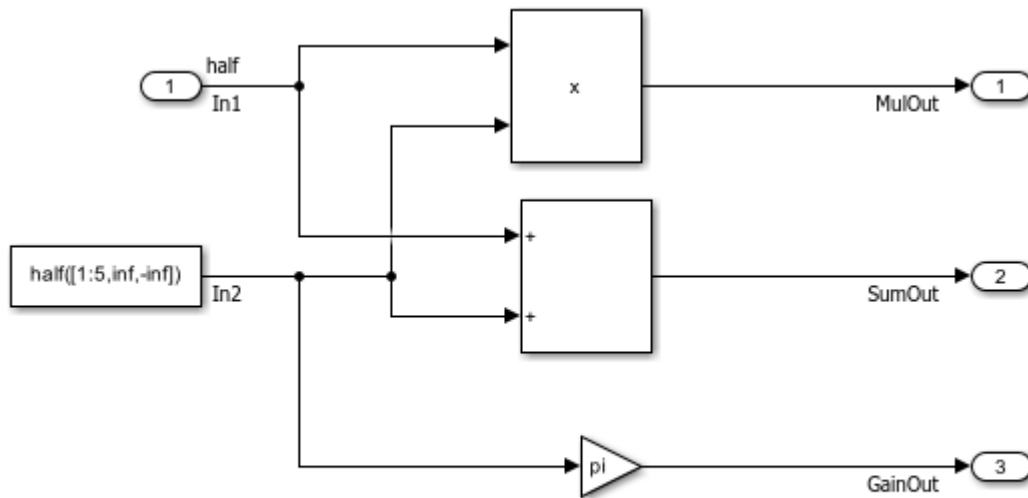
Use the `target.add` function to add the target object to the internal database.

```
target.add(arm_half);
```

Open Half-Precision Model

Open the model for which you want to generate native half-precision C code. The model should use half-precision data types and be configured for code generation.

```
model = 'simpleHalfModel';
open_system(model);
```



Copyright 2021 The MathWorks, Inc.

Set up Simulink Model with Target Hardware

Optionally, enable the model parameter “Inherit floating-point output type smaller than single precision” to propagate half precision.

```
set_param(model, 'InheritOutputTypeSmallerThanSingle', 'on');
```

Set the code generation target to the ARM target modified for half precision.

```
set_param(model, 'ProdEqTarget', 'on');
set_param(model, 'ProdHWDeviceType', ...
    'ARM Compatible->ARM Cortex-A Half-Precision');
```

Set up Simulink Model with Armclang Compiler Toolchain

Register the Armclang compiler toolchain according to Setup and Configure Armclang Compiler Toolchain for Code Generation.

Set up the Armclang compiler toolchain for code generation.

```
set_param(model, 'SystemTargetFile', 'ert.tlc');
set_param(model, 'GenCodeOnly', 'off');
set_param(model, 'Toolchain', 'Armclang Compiler');
set_param(model, 'BuildConfiguration', 'Specify');
```

Configure the Armclang compiler for custom C flags.

```
compilerOptions = get_param(model, 'CustomToolchainOptions');
compilerIdx = find(strcmp(compilerOptions, 'C Compiler'));
```

Add half-precision flags.


```
appendToCCompilerOptions = ' --target=arm-arm-none-eabi -mcpu=cortex-a75+fp16';
compilerOptions{compilerIdx+1} = strcat(compilerOptions{compilerIdx+1},appendToCCompilerOptions)
```

Apply the custom toolchain options to the model.

```
set_param(model, 'CustomToolchainOptions', compilerOptions );
```

You can confirm these settings in the **Code Generation** pane of the Configuration Parameters dialog box.

Target selection

System target file: ert.tlc
 Language: C
 Language standard: C89/C90 (ANSI)
 Description: Real-Time Workshop Embedded

Build process

Generate code only
 Package code and artifacts

Toolchain settings

Toolchain: Armclang Compiler
 Build configuration: Specify

| Tool | Options |
|------------|--|
| Assembler | --md \$(ASFLAGS) |
| C Compiler | -DMW_ARM_CLANG -c -O1 --target=arm-arm-none-eabi -mcpu=cortex-a75+fp16 |
| Linker | --list \$(PRODUCT_M --info sizes --info totals --info unused --info veneers --callgraph |

Generate Code

Once the model has been set up for the ARM Cortex-A target hardware and the Armclang compiler toolchain, you can generate C code. You can only run the model executable on an ARM target or emulator.

Use the `slbuild` function to build a standalone executable for the model.

```
slbuild(model);
```

You can inspect the code generation report to confirm that the custom header and type definitions are used.

```

13  * Target selection: ert.tlc
14  * Embedded hardware selection: ARM Compatible->ARM Cortex-A Half
15  * Code generation objectives: Unspecified
16  * Validation result: Not run
17  */
18
19  #ifndef HALF_TYPE_H
20  #define HALF_TYPE_H
21  #include "rtwtypes.h"
22
23  /* Custom Headers */
24  #include <arm_fp16.h>
25
26  /* Type definition */
27  typedef _Float16 real16_T;

```

The generated code uses the native half-precision data type.

```

40  for (i = 0; i < 7; i++) {
41      /* Output: '<Root>/Out1' incorporates:
42       * Constant: '<Root>/Constant'
43       * Inport: '<Root>/In1'
44       * Product: '<Root>/Product'
45       */
46      Y.Out1[i] = (real16_T)((real32_T)U.In1 * (real32_T)ConstWithInitP.Constant_Value[i]);
47
48      /* Sum: '<Root>/Sum' incorporates:
49       * Constant: '<Root>/Constant'
50       * Inport: '<Root>/In1'
51       */
52      tmp = (real32_T)(real16_T)((real32_T)U.In1 + (real32_T)ConstWithInitP.Constant_Value[i]);
53
54      /* Output: '<Root>/Out2' incorporates:
55       * Sum: '<Root>/Sum'
56       */
57      Y.Out2[i] = (real16_T)tmp;
58  }

```

Register ARM Cortex-A with GCC Compiler

In this example, an ARM Cortex-A processor is used as the hardware target. The model is configured to use this ARM target and the GCC compiler toolchain.

Use the `target.create` function to create an ARM Cortex-A processor target.

```

arm_half = target.create('Processor', ...
'Manufacturer', 'ARM Compatible', ...
'Name', 'ARM Cortex-A Half-Precision');

```

Add the language implementation. The 32-bit version of GCC for ARM has half precision enabled by default. Use the `target.get` function to retrieve the target object from the internal database.

```
li = target.get('LanguageImplementation', "GNU GCC ARM 32-bit");
```

Replace the default language implementation for ARM Cortex with Armclang.

```
arm_half.LanguageImplementations = li;
```

Use the `target.add` function to add the target object to the internal database.

```
target.add(arm_half);
```

Register ARM Target Hardware with Custom Language Implementation

In this example, you create a new custom language implementation with half precision for a compatible ARM target.

Use the `target.create` function to copy the ARM Compatible-ARM Cortex language implementation.

```
languageImplementation = target.create('LanguageImplementation', ...
    'Name', 'ARM with half', ...
    'Copy', 'ARM Compatible-ARM Cortex');
```

Specify custom half information and target specific headers, as given by your target hardware documentation. For more information, see “Register New Hardware Devices”.

```
customHalf = target.create('FloatingPointDataType', ...
    'Name', 'BCM2711 Half Type', ...
    'TypeName', 'float16_T', ...
    'LiteralSuffix', 'f16', ...
    'Size', 16, ...
    'SystemIncludes', "arm_fp16.h, arm_neon.h");
```

```
languageImplementation.DataTypes.NonStandardDataTypes = customHalf;
```

Provide information about your target processor.

```
% Broadcom BCM2711
% Quad core Cortex-A72 (ARM v8) 64-bit SoC
pi4a72 = target.create('Processor', ...
    'Manufacturer', 'Broadcom', ...
    'Name', 'BCM2711');
```

Add the custom half precision language implementation.

```
pi4a72.LanguageImplementations = languageImplementation;
```

Use the `target.add` function to add the target object to the internal database.

```
target.add(pi4a72);
```

Register Other Hardware Targets for Half Precision

This example uses a Renesas® RH850 to show you how to set up non-ARM targets for half precision.

Use the `target.create` function to create the target processor.

```
prh850 = target.create('Processor', ...  
    'Manufacturer', 'Renesas', ...  
    'Name', 'Renesas-RH850 With Half', ...  
    'Copy', 'Renesas-RH850');
```

Add the language implementation.

```
li = target.create('LanguageImplementation', ...  
    'Name', 'Renesas-RH850 with Half', ...  
    'Copy', 'Renesas-RH850');
```

Provide additional information the custom half-precision implementation for this hardware.

```
customHalf = target.create('FloatingPointDataType', ...  
    'Name', 'Renesas Half Type');  
customHalf.TypeName = '__fp16';  
customHalf.Size = 16;  
customHalf.LiteralSuffix = '';  
customHalf.SystemIncludes = '';
```

Add the custom half-precision data type.

```
li.DataTypes.NonStandardDataTypes = customHalf;  
prh850.LanguageImplementations = li;
```

Use the `target.add` function to add the target object to the internal database.

```
target.add(prh850);
```

See Also

`target.FloatingPointDataType` | `target.add` | `target.create` | `target.get` | `target.remove`

Related Examples

- “The Half-Precision Data Type in Simulink” on page 51-2
- “Register New Hardware Devices”

External Websites

- Clang Language Extensions for Half-Precision Floating Point
- Arm Compiler armclang Reference Guide: Half-precision floating-point data types
- GCC Half-Precision Floating Point
- Reduce the Program Data Size with Ease! Introducing Half-Precision Floating-Point Feature in Renesas Compiler Professional Edition

Half-Precision Field-Oriented Control Algorithm

This example shows how to implement a Field-Oriented Control (FOC) algorithm for a Permanent Magnet Synchronous Machine (PMSM). The example shows both a single-precision floating-point implementation and a half-precision floating-point implementation. When an algorithm contains large or unknown dynamic ranges (for example integrators in feedback loops) or when the algorithm uses operations that are difficult to design in fixed-point (for example, `atan2`), it can be advantageous to use floating-point representations. The half-precision data type occupies only 16 bits of memory, but its floating-point representation enables it to handle wider dynamic ranges than integer or fixed-point data types of the same size.

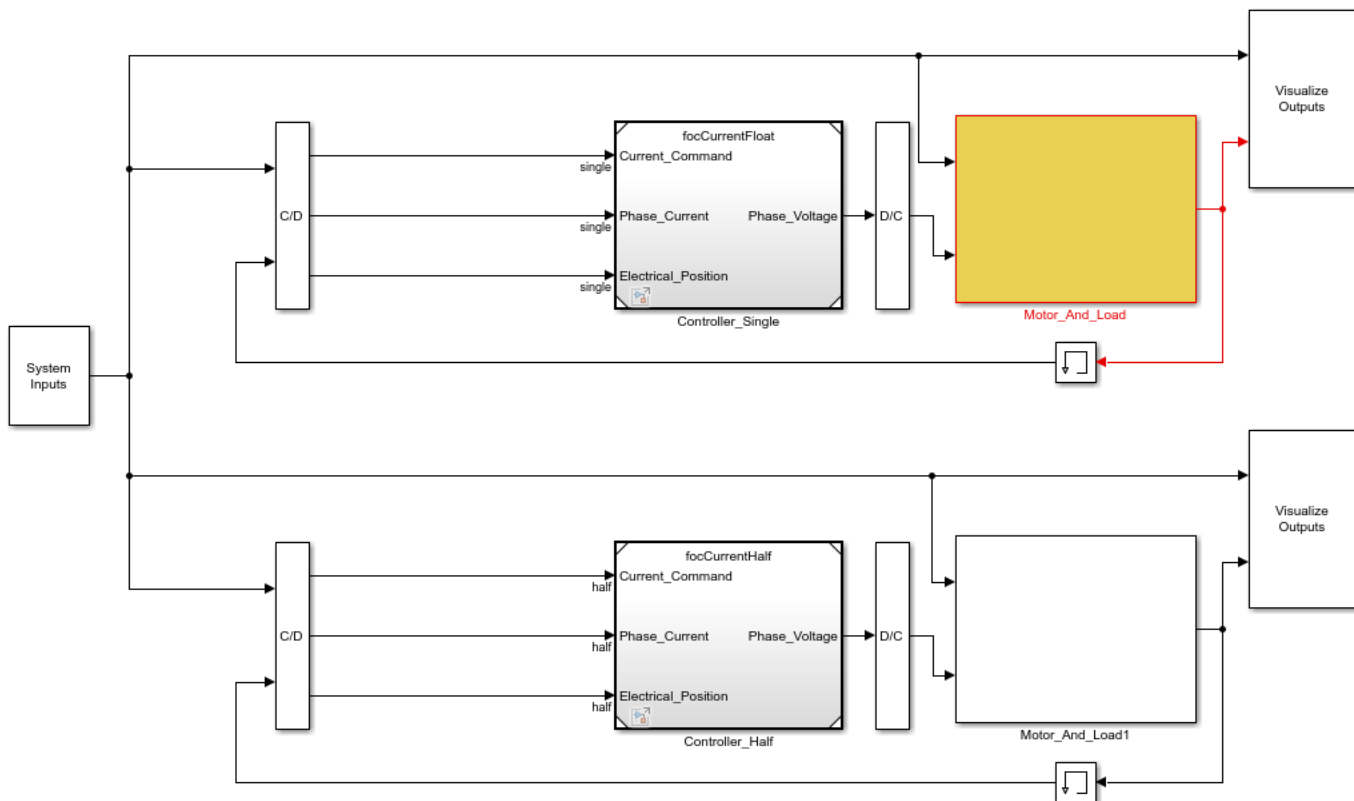
Load the data required to simulate the model.

```
focModelData
```

Open the `ex_foc_current` model. This model uses the same source block for two versions of a field-oriented control algorithm. The first version uses single-precision data types, while the second uses half-precision data types.

```
model = 'ex_foc_current.slx';
open_system(model)
```

Field-Oriented Control Current Control Test Bench

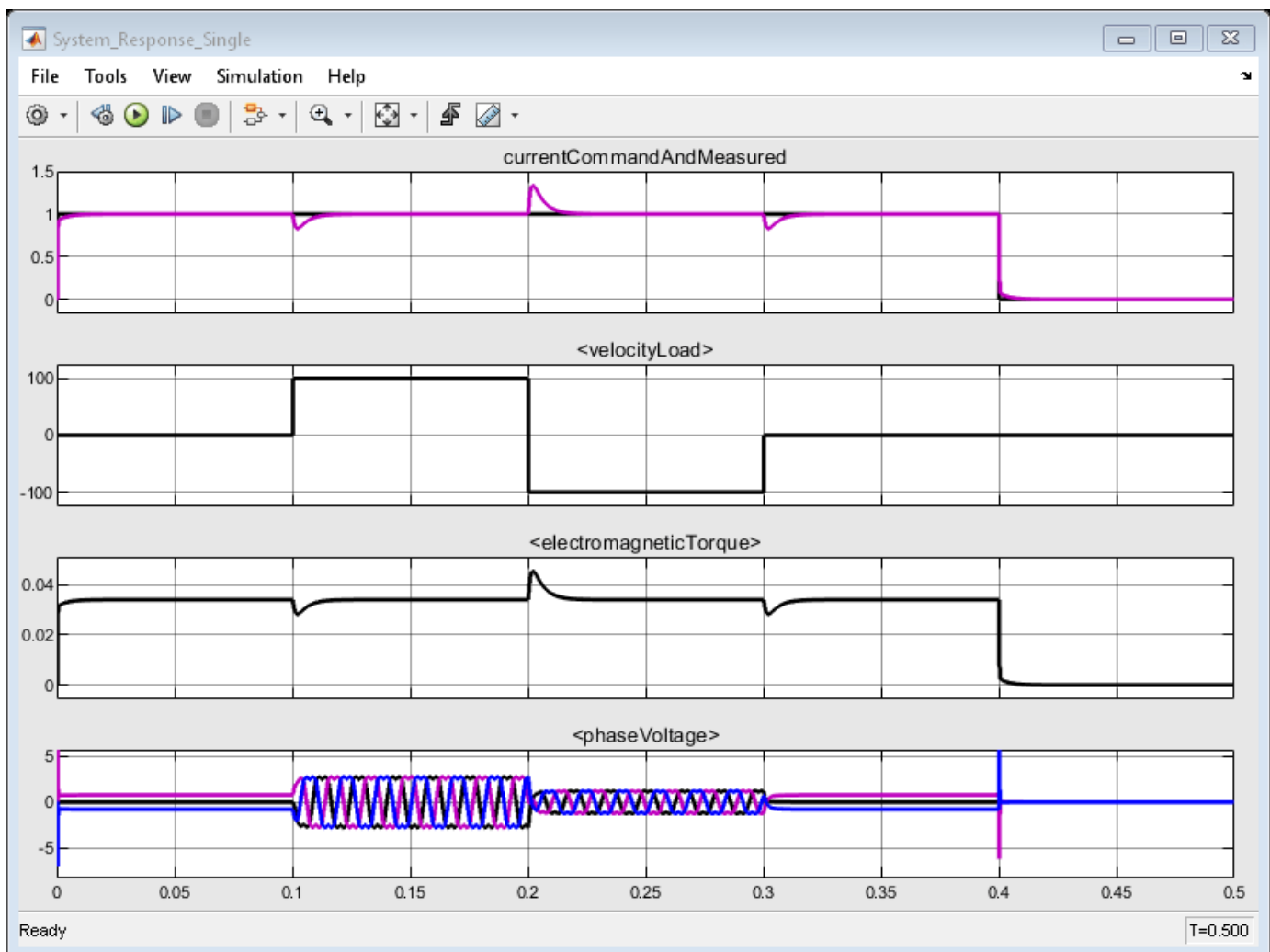


Simulate the model. You can see from the scope that the response of the single-precision implementation is identical to the response of the half-precision implementation.

```
sim(model)
```

```
ans =
```

```
Simulink.SimulationOutput:  
  logTestBench: [1x1 Simulink.SimulationData.Dataset]  
  
SimulationMetadata: [1x1 Simulink.SimulationMetadata]  
  ErrorMessage: [0x0 char]
```



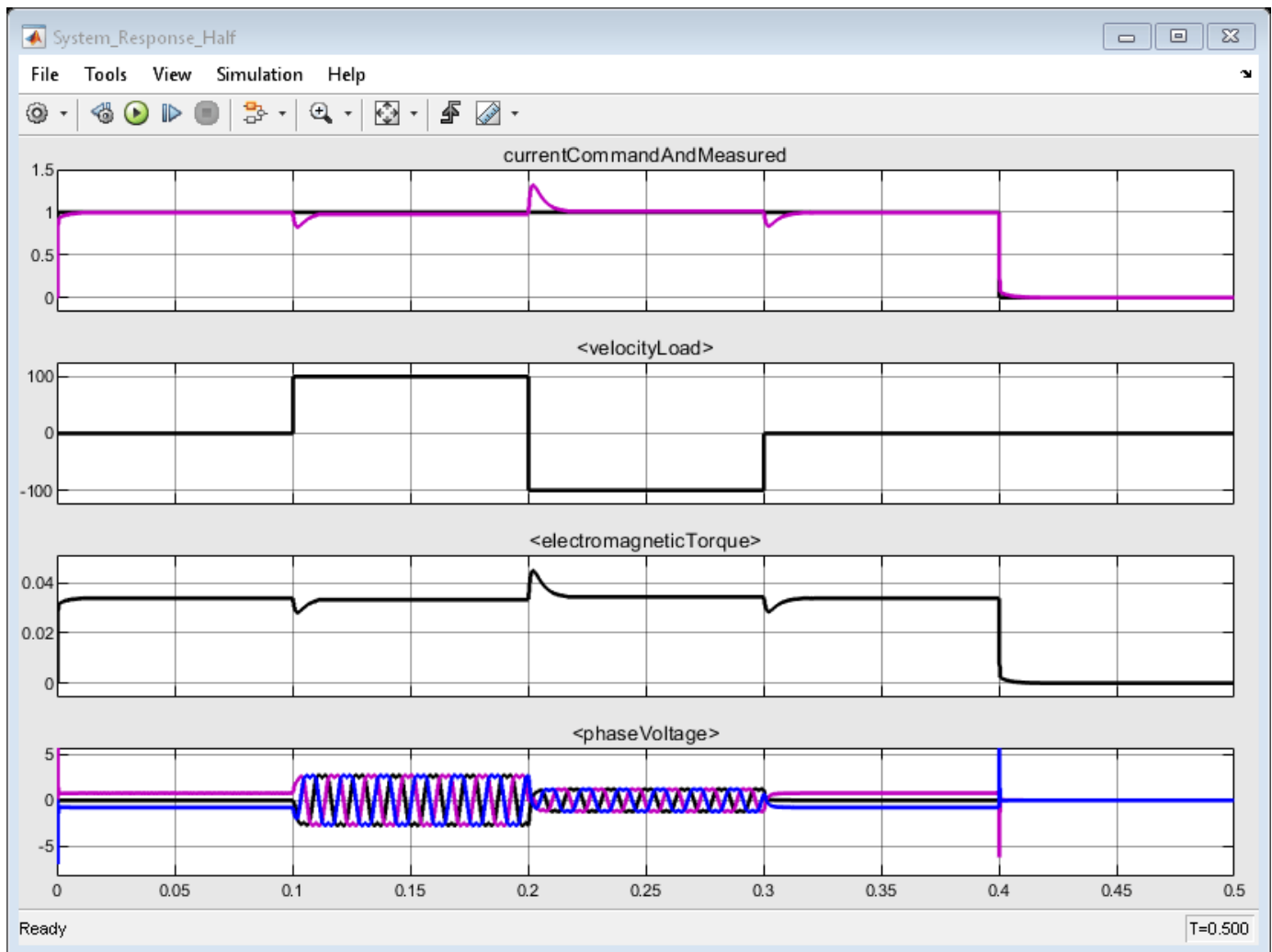
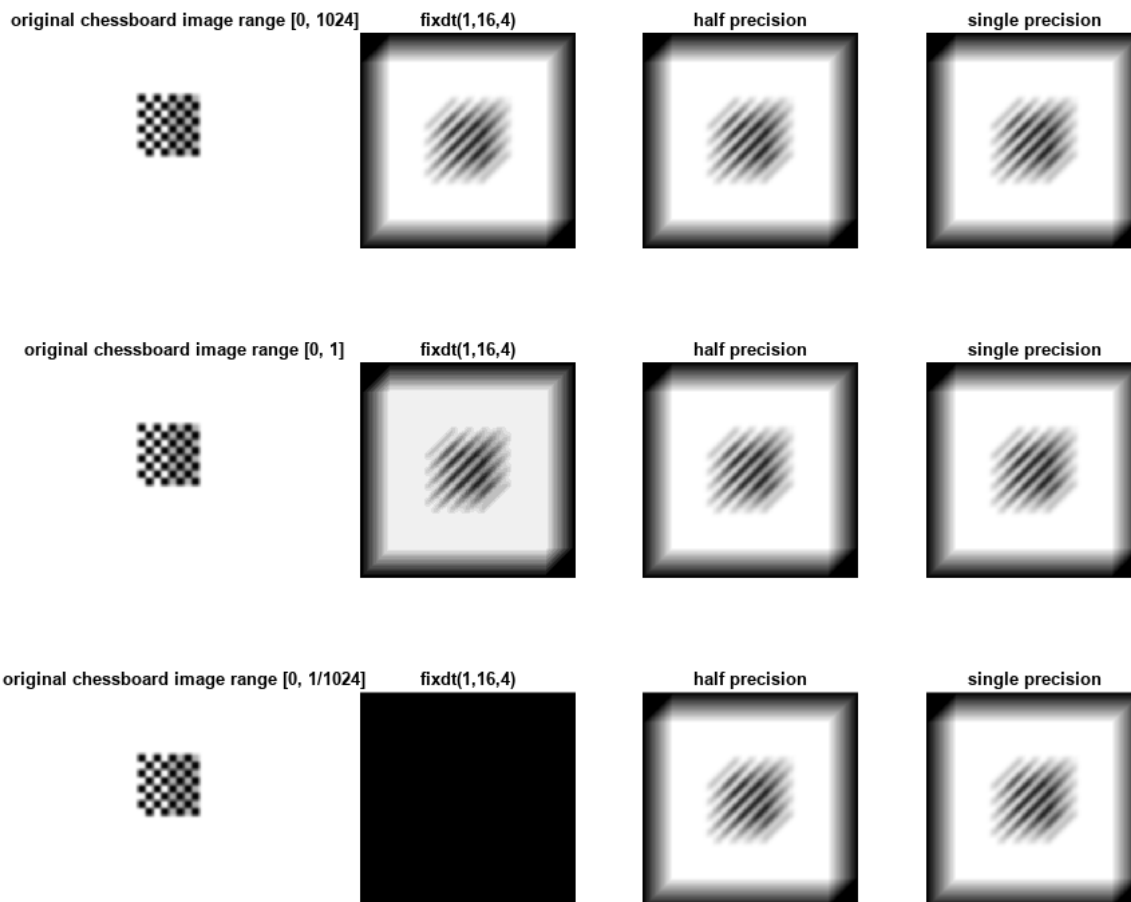


Image Quantization with Half-Precision Data Types

This example shows the effects of quantization on images. The `ex_imagequantization` model, computes the two-dimensional fourier transform of an image of a checkerboard. The original image is displayed in the left-most column, and the result is displayed with fixed-point, half-precision, and single-precision data types. You can see in the resulting images that, while the fixed-point data type does not always produce an acceptable result, the half-precision data type, which uses the same number of bits as the fixed-point data type, produces a result comparable to the single-precision result.

```
model = 'ex_image_quantization.slx';
open_system(model);
sim(model)
```



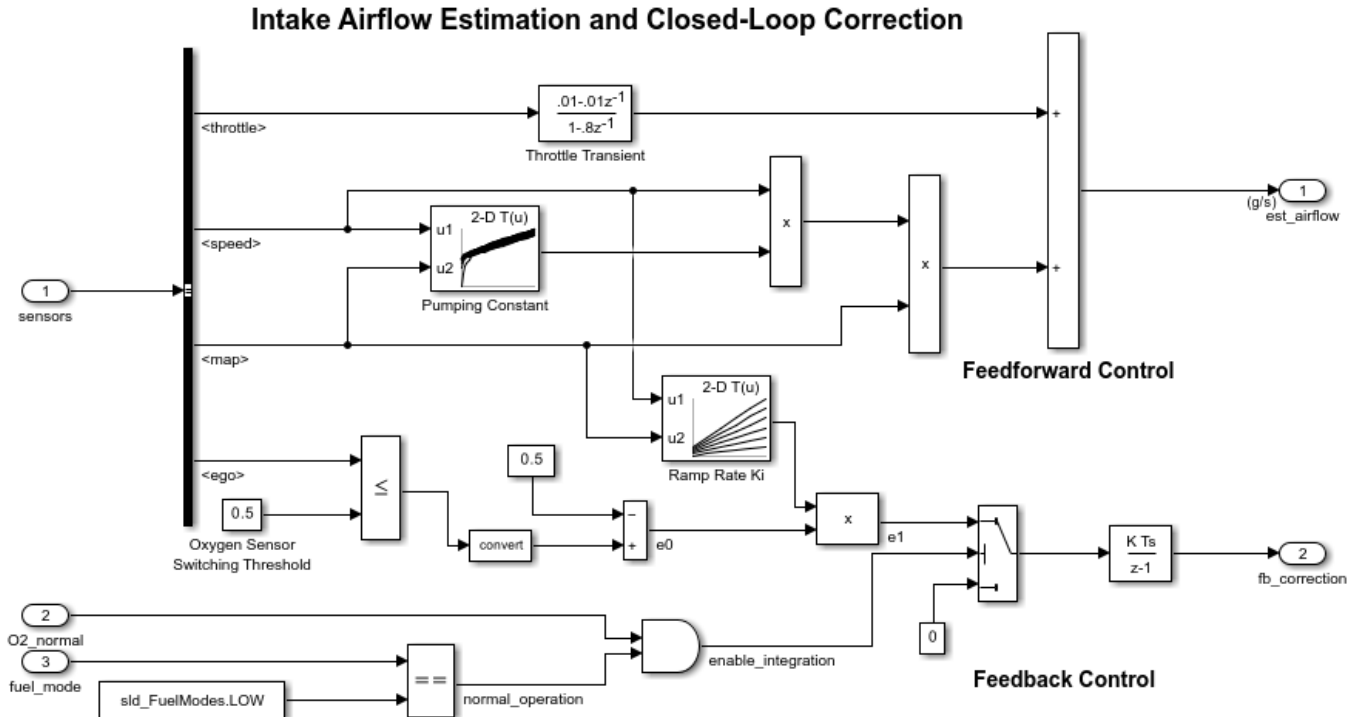
Convert Single Precision Lookup Table to Half Precision

This example shows how to convert a single-precision lookup table to use half precision. Half precision is the storage type; the lookup table computations continue to be performed using single precision. After the conversion, the examples halves the memory size of the Lookup Table blocks while keeping the desired system performance.

Task 1: Simulate and Obtain Baseline

1. Open the `airflow_calc` subsystem of the `sldemo_fuelsys` example model. It contains the Lookup Table blocks with the single precision table and breakpoint data to be converted to half precision.

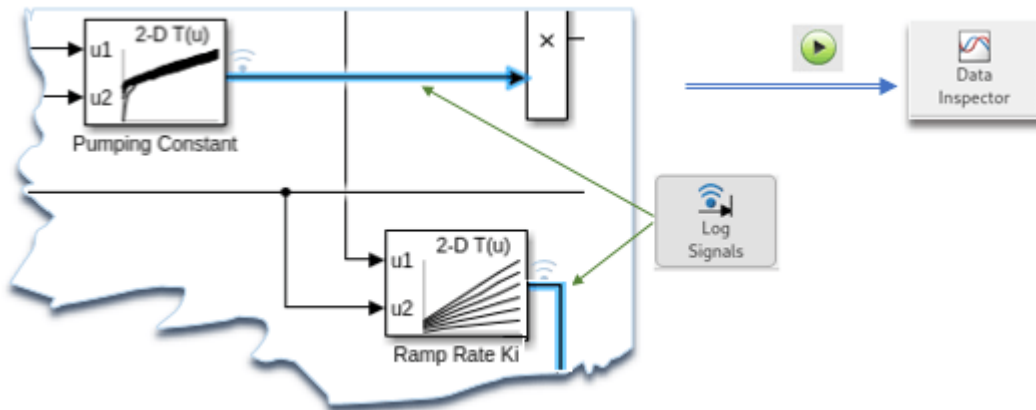
```
load_system('sldemo_fuelsys');
open_system('sldemo_fuelsys/fuel_rate_control/airflow_calc');
```



2. In MATLAB®, change your current folder to a writable folder.

3. Check that the table and breakpoints **Pumping Constant** and **Ramp Rate Ki** Lookup Table block data types are set as single.

4. Select output signals for both Lookup Table blocks. Mark them for logging. Simulate the model in normal mode and use the outputs of this run as a baseline.



Task 2: Analyze and Convert Data to Half

1. Get the table and breakpoint variable data from the model workspace for the specified Lookup Table block. For example, using the 'Pumping Constant' block:

```
mdlWks = get_param('sldemo_fuelsys', 'ModelWorkspace');
currentBlk = 'sldemo_fuelsys/fuel_rate_control/airflow_calc/Pumping Constant';
```

2. Analyze the table data and convert it to half precision if it can fit in the range of half precision type [-65504, 65504]. In Simulation Data Inspector, simulate the model and compare the logged output of the current block with the baseline. Verify that the output is within the specified absolute tolerance and relative tolerance. For example, this code sets both tolerances to $|0.01|$.

```
numBytesSaved(1) = analyzeDataConvertToHalf(mdlWks, currentBlk, 'Table');
```

3. Repeat step 2 of Task 2 for data of each of the breakpoints of the current block. The breakpoints data should remain monotonically increasing after converting from single to half precision.

```
numBytesSaved(2) = analyzeDataConvertToHalf(mdlWks, currentBlk, 'BreakpointsForDimension1');
numBytesSaved(3) = analyzeDataConvertToHalf(mdlWks, currentBlk, 'BreakpointsForDimension2');
totalBytesSaved = sum(numBytesSaved(1:3));
```

4. Repeat steps 2 and 3 of Task 2 for the remaining Lookup Table blocks in the same subsystem. Then, examine the number of bytes saved for the table and breakpoints of the n-D Lookup Table blocks.

```
currentBlk = 'sldemo_fuelsys/fuel_rate_control/airflow_calc/Ramp Rate Ki';
numBytesSaved(4) = analyzeDataConvertToHalf(mdlWks, currentBlk, 'Table');
numBytesSaved(5) = analyzeDataConvertToHalf(mdlWks, currentBlk, 'BreakpointsForDimension1');
numBytesSaved(6) = analyzeDataConvertToHalf(mdlWks, currentBlk, 'BreakpointsForDimension2');
```

```
numBytesSaved =
```

```
684    36    38    72    12    12
```

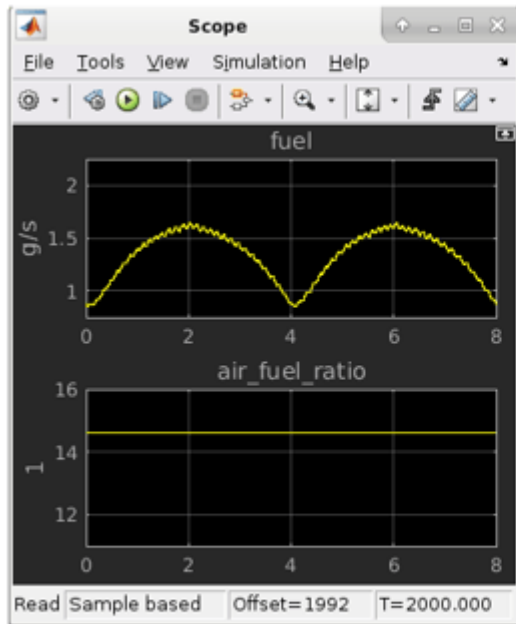
```
totalBytesSaved = totalBytesSaved + sum(numBytesSaved(4:6)),
```

```
totalBytesSaved =
```

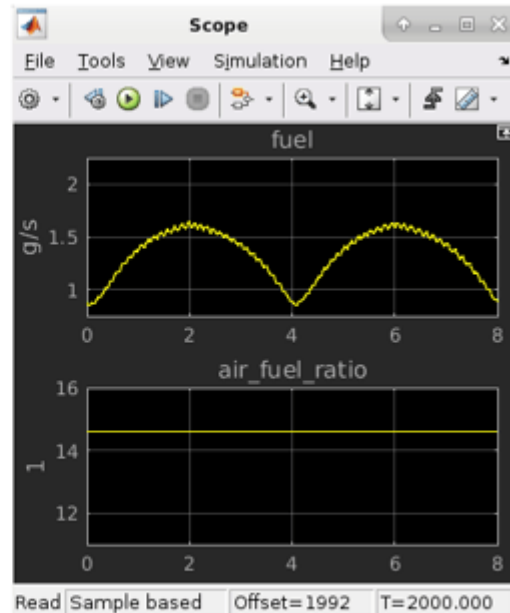
854

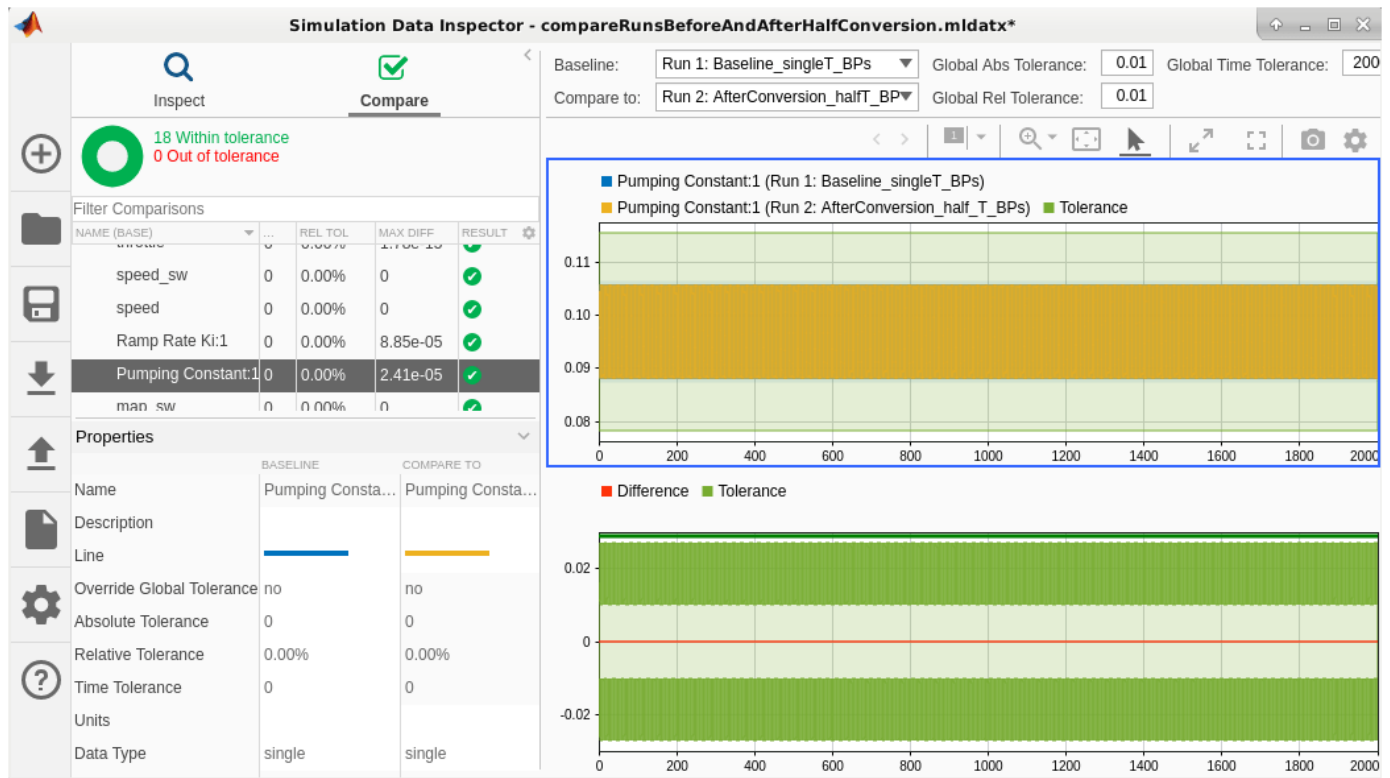
Task 3: Simulate and Compare

1. Simulate the converted model with half precision table and breakpoints.
2. Compare the 'fuel', 'air_fuel_ratio', etc. signals with the baseline in the Scope block and in the Simulation Data Inspector. Observe that the desired performance is still achieved.

**Single Precision Table & Breakpoints**

VS

**Half Precision Table & Breakpoints**



3. Save the model with a different name in the writable folder.

Task 4: Generate Code and Verify the Memory Optimization

1. Right-click the **fuel_rate_control** subsystem and select **C/C++ Code > Build This Subsystem**. To generate code, click **Build** in the **Build code for Subsystem** dialog box.
2. When the **build** finishes processing, a code generation report displays.
3. Click the `half_type.h` and `fuel_rate_control.c` files. Notice the half precision (`real16_T`) type definition, the `real16_T` type table, and breakpoint pointers in the `look2_ifbhlftHdIf_linlca` function call interface.

```
File: half\_type.h 27 /* C type definition */
28 typedef struct {
29     uint16_T bitPattern;
30 } half_t;
31
32 typedef half_t real16_T;
```

File: [fuel_rate_control.c](#)

```
70 real32_T look2_ifbhlftHdIf_linlca(real32_T u0, real32_T u1, const real16_T bp0[],
71     const real16_T bp1[], const real16_T table[], const uint32_T maxIndex[],
72     uint32_T stride)
```

4. Open the `fuel_rate_control.h` file and observe that **854 bytes** have been saved by using half precision as the storage type for table and breakpoints.

```

77 /* Constant parameters (default storage) */
78 typedef struct {
    ... ..
137 /* Computed Parameter: RampRateKi_tableData 152 /* Computed Parameter: PumpingConstant_tableData
138  * Referenced by: '<S2>/Ramp Rate Ki' 153  * Referenced by: '<S2>/Pumping Constant'
139  */ 154  */
140 real16_T RampRateKi_tableData[36]; 155 real16_T PumpingConstant_tableData[342];
141 156
142 /* Computed Parameter: RampRateKi_bp01Data 157 /* Computed Parameter: PumpingConstant_bp01Data
143  * Referenced by: '<S2>/Ramp Rate Ki' 158  * Referenced by: '<S2>/Pumping Constant'
144  */ 159  */
145 real16_T RampRateKi_bp01Data[6]; 160 real16_T PumpingConstant_bp01Data[18];
146 161
147 /* Computed Parameter: RampRateKi_bp02Data 162 /* Computed Parameter: PumpingConstant_bp02Data
148  * Referenced by: '<S2>/Ramp Rate Ki' 163  * Referenced by: '<S2>/Pumping Constant'
149  */ 164  */
150 real16_T RampRateKi_bp02Data[6]; 165 real16_T PumpingConstant_bp02Data[19];
151 166 } ConstParam;

```

```
close_system('sldemo_fuelsys',0);
```

Digit Classification with Half-Precision Data Types

This example compares the results of a trained neural network classification model in Simulink® in double precision and half precision. The model classifies images from the MNIST handwritten digit dataset.

To begin, load the data for the model, and specify the size of the test data set.

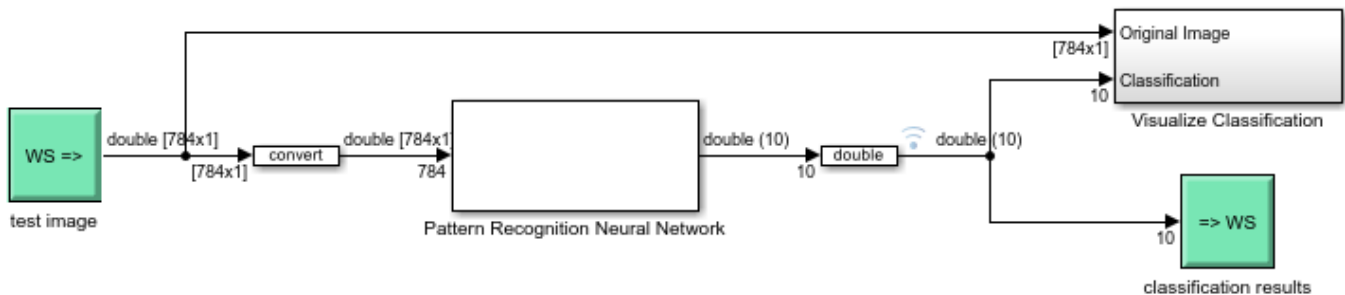
```
load trainImage.mat
testSetMaxIdx = 10;
```

Simulate the Model with Double-Precision Types

This model uses `numericType` objects to specify parameter, signal, and block output data types. To simulate the model using double-precision data types, define the `numericType` objects in the base workspace and set the data type of the objects to 'double'. Simulate the model.

```
floatType = numericType('double');
activationType = numericType('double');

model = 'ex_imagerecog_half.slx';
open_system(model);
sim(model);
```



The double-precision simulation results in 100% classification accuracy.

```
[c, ~] = confusion(ttestsubset(:,1:testSetMaxIdx), ytest.Data(2:testSetMaxIdx+1,:));
fprintf('Percentage Correct Classification : %f%%\n', 100*(1-c));
fprintf('Percentage Incorrect Classification : %f%%\n', 100*c);
```

```
Percentage Correct Classification : 100.000000%
Percentage Incorrect Classification : 0.000000%
```

Simulate the Model with Half-Precision Types

To simulate the model in half precision, redefine the `numericType` objects and set their data type to 'half'. Simulate the model.

```
floatType = numericType('half');
activationType = numericType('single');
sim(model);
```

In this example, there is no loss of accuracy when using the half-precision data type. The half-precision simulation also results in 100% classification accuracy.

```
[c, ~] = confusion(ttestsubset(:,1:testSetMaxIdx), ytest.Data(2:testSetMaxIdx+1,:));  
fprintf('Running Simulation with half precision :\n');  
fprintf('Percentage Correct Classification   : %f%\n', 100*(1-c));  
fprintf('Percentage Incorrect Classification : %f%\n', 100*c);
```

```
Running Simulation with half precision :  
Percentage Correct Classification   : 100.000000%  
Percentage Incorrect Classification : 0.000000%
```


Design Cost Estimation

- “Design Cost Model Metrics” on page 52-2
- “How to Collect Design Cost Metrics” on page 52-4

Design Cost Model Metrics

In this section...

“Data Segment Estimate” on page 52-2

“Operator Count” on page 52-3

Learn about design cost metrics that return metric data on the cost of implementing your Simulink design in embedded C code.

Data Segment Estimate

Metric ID: DataSegmentEstimate

Estimate the memory consumed by the data segment of generated code.

Description

Use this metric to estimate the amount of memory consumed, in bytes, by the data segment of code generated for the specified model unit. A data segment is a part of an object file or the corresponding address space of a program that holds initialized global variables and static local variables created during code generation. The size of the data segment is determined by the size of the values in the source code and does not change at run time.

To collect data for this metric:

- Use `getMetrics` with the metric identifier `DataSegmentEstimate`.

Collecting data for this metric requires a Fixed-Point Designer license.

Results

For this metric, instances of `metric.Result` return `Value` as an integer representing the total cost of the design in bytes. For a detailed breakdown of results, use the `generateReport` function.

Capabilities and Limitations

The metric:

- Analyzes one or more design units in a project, where a design unit represents a standalone Simulink or an entire model reference hierarchy.
- Requires that when the `execute` function specifies an `'ArtifactScope'`, then `scope` must refer to a top-level Simulink model.
- Metrics are collected only for designs that are code generation ready. If you collect metrics for a model reference hierarchy, each design within the hierarchy, including the top-level model, must be ready for code generation. Models that are not code generation ready will be ignored during metric execution and produce an error. You can use the `getArtifactErrors` function to see errors that occur during metric execution.

See Also

For an example of collecting metrics programmatically, see “How to Collect Design Cost Metrics” on page 52-4.

Operator Count

Metric ID: OperatorCount

Estimate the size of a design based on a weighted count of operators used in generated code. This metric is an abstraction of the actual size of generated code and is returned as a unitless value.

Description

Use this metric to estimate the size of generated code for the specified design unit or model reference hierarchy.

To collect data for this metric:

- Use `getMetrics` with the metric identifier `OperatorCount`.

Collecting data for this metric requires a Fixed-Point Designer license.

Results

For this metric, instances of `metric.Result` return `Value` as an integer representing the total cost of the design. For a detailed breakdown of results, use the `generateReport` function.

Capabilities and Limitations

The metric:

- Analyzes one or more design units in a project, where a design unit represents a standalone Simulink or an entire model reference hierarchy.
- Requires that when the `execute` function specifies an `'ArtifactScope'`, then `scope` must refer to a top-level Simulink model.
- Metrics are collected only for designs that are code generation ready. If you collect metrics for a model reference hierarchy, each design within the hierarchy, including the top-level model, must be ready for code generation. Models that are not code generation ready will be ignored during metric execution and produce an error. You can use the `getArtifactErrors` function to see errors that occur during metric execution.

See Also

For an example of collecting metrics programmatically, see “How to Collect Design Cost Metrics” on page 52-4.

See Also

`metric.Engine` | `metric.Result`

Related Examples

- “How to Collect Design Cost Metrics” on page 52-4

How to Collect Design Cost Metrics

In this section...

“Open Project” on page 52-4

“Collect Metric Results” on page 52-4

“Access High-Level Results Programmatically” on page 52-5

“Generate Report to Access Detailed Results” on page 52-6

“Operator Count” on page 52-6

“Data Segment Table” on page 52-8

Learn how to programmatically assess the cost of implementing your design in embedded C code. Design cost metrics analyze your model and report detailed cost data that can be traced back to the blocks in the Simulink model. You can use design cost metrics to:

- Estimate the program size based on a weighted count of operators in the generated code.
- Estimate the total size, in bytes, of all global variables and static local variables used during code generation.

Design cost metrics analyze one or more design units in a project, where a design unit represents a standalone Simulink, or an entire model reference hierarchy. After collecting metrics, you can programmatically access high-level results or generate a report for detailed cost breakdown information. By running a script that collects these metrics, you can automatically analyze the cost of your design to, for example, assess the impact of different design alternatives before deploying to hardware.

Open Project

To open the project that includes the models, in the MATLAB Command Window, enter this command.

```
dashboardCCProjectStart
```

This project contains models and requirements and test cases for the models. Some of the requirements have traceability links to the models and test cases, which help to verify that the functionality of a model meets the requirements.

Note To collect design cost metrics, all design files must be in a project. Your design must also be ready for code generation, including all models contained in a model reference hierarchy. Models that are not code generation ready will be ignored during metric execution and produce an error. You can use the `getArtifactErrors` function to see errors that occur during metric execution.

Collect Metric Results

Create a `metric.Engine` object.

```
metric_engine = metric.Engine();
```

A `metric.Engine` object represents the metric engine that you can execute with the `execute` object function to collect metric data on your design.

If models in your project have changed, update the trace information for `metric_engine` to reflect any pending artifact changes.

```
updateArtifacts(metric_engine)
```

Create an array of metric identifiers for the metrics you want to collect. For this example, create a list of the metric identifiers used for design cost estimation.

```
metric_Ids = {'OperatorCount', 'DataSegmentEstimate'}
```

You can use the `getAvailableMetricIds` function to create a full list of all available metric identifiers. For a list of design cost metrics and their identifiers, see “Design Cost Model Metrics” on page 52-2. For additional model testing metrics, see “Model Testing Metrics” (Simulink Check).

Collect Results for One Design Unit

You can collect metric results for one design unit in the project, where a design unit represents a standalone Simulink or the top-level model of a model reference hierarchy. When you collect and view results for a design unit, the metrics return data only for the artifacts that trace to the specified design unit.

Collect the metric results for the `cc_CruiseControl` model.

Create an array that identifies the path to the model file in the project and the name of the model.

```
unit = {fullfile(pwd, 'models', 'cc_CruiseControl.slx'), 'cc_CruiseControl'};
```

Execute the engine. Use `'ArtifactScope'` to specify the unit for which you want to collect results. The engine runs the metrics for only the artifacts that trace to the model that you specify. Collecting results for the design cost metrics requires a Fixed-Point Designer license.

```
execute(metric_engine, metric_Ids, 'ArtifactScope', unit)
```

Collect Results for Each Design Unit in Project

To collect the results for each design unit in the project, execute the engine without the argument for `ArtifactScope`.

```
execute(metric_engine, metric_Ids)
```

The project contains `cc_CruiseControl`, which is the top-level model in a model reference hierarchy. This model reference hierarchy represents one design unit.

Access High-Level Results Programmatically

To access the results programmatically, use the `getMetrics` function. The function returns the `metric.Result` objects that contain the result data for the specified unit and metrics. For this example, store the results for the `OperatorCount` and `DataSegmentEstimate` metrics in corresponding arrays.

```
results_OperatorCount = getMetrics(metric_engine, 'OperatorCount');
results_DataSegmentEstimate = getMetrics(metric_engine, 'DataSegmentEstimate');
```

The `OperatorCount` metric returns a unitless estimate of the program size based on operator count. Use the `disp` function to display the total cost of the design.

```
disp(['Unit: ', results_OperatorCount.Artifacts.Name])
disp(['Total Cost: ', num2str(results_OperatorCount.Value)])
```

```
Unit: cc_CruiseControl  
Total Cost: 333
```

This result shows that for `cc_CruiseControl`, the total cost of the design is 333. This is a unitless value based on a weighted count of operators in the generated code, and is an abstraction of the total program size.

The metric `DataSegmentEstimate` returns the estimated size of the data segment of the program. This value represents the total size, in bytes, of all global variables and static local variables used during code generation. Use the `disp` function to display the total data segment size.

```
disp(['Unit: ', results_DataSegmentEstimate.Artifacts.Name])  
disp(['Data Segment Size (bytes): ', num2str(results_DataSegmentEstimate.Value)])
```

```
Unit: cc_CruiseControl  
Data Segment Size (bytes): 79
```

This result shows that for `cc_CruiseControl`, the total data segment size estimate is 79 bytes.

Generate Report to Access Detailed Results

Generate a report that contains a detailed breakdown of design cost metric results. For this example, specify the HTML file format, use `pwd` to provide the path to the current folder. Name the report `'MetricResultsReport.html'`.

```
reportLocation = fullfile(pwd, 'MetricResultsReport.html');  
generateReport(metric_engine, 'App', 'DesignCostEstimation', ...  
    'Type', 'html-file', 'Location', reportLocation);
```

Note Report generation requires that the metric collection be executed in the current session. To recollect design cost metrics, first use the `deleteMetrics` function to delete the `metric.Result`. Then, use the `execute` function to collect metrics.

Open the HTML report. The report is in the current folder, at the root of the project.

```
web("MetricResultsReport.html")
```

To open the table of contents and navigate to results for each unit, click the menu icon in the top-left corner of the report. The report contains sections for each metric collected.

Operator Count

The Operator Count section of the report contains information on the cost of the design. The High Level Statistics section gives the total cost of the design and a breakdown of the top five most expensive blocks in the `cc_CruiseControl` model reference hierarchy.

Chapter 1. cc_CruiseControl

1.1. Operator Count

1.1.1. High Level Statistics

Total Cost of Design is 333

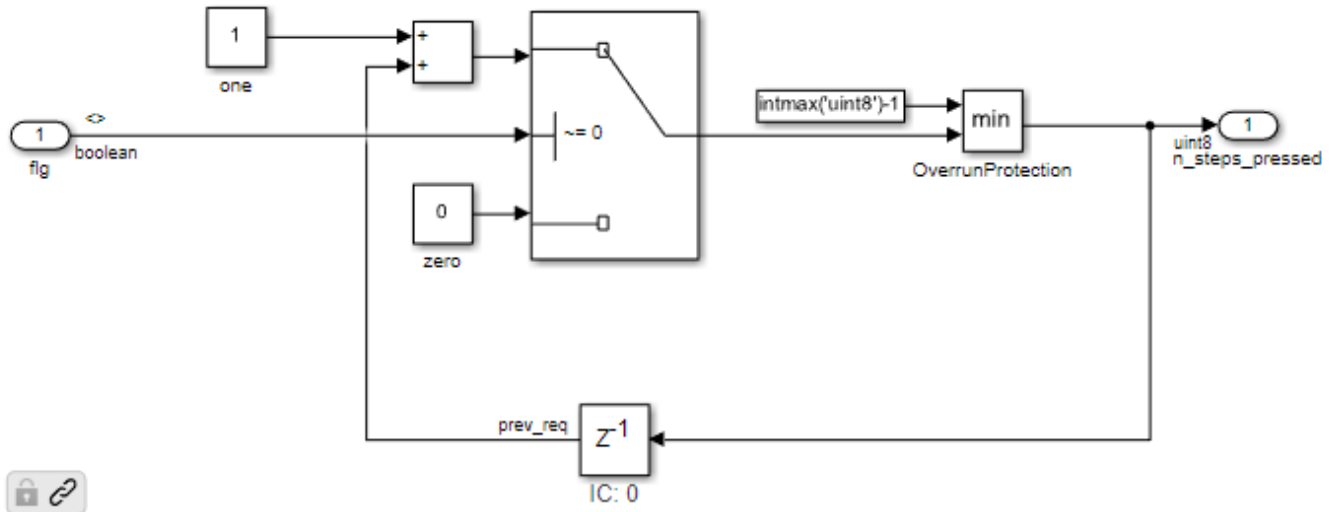
Table 1.1. Most Expensive Blocks

| Name | Cost |
|--|------|
| "cc_CruiseControl/DriverSwitchRequest" | 157 |
| "cc_CruiseControl/ControlMode" | 127 |
| "cc_CruiseControl/TargetSpeedThrottle" | 34 |
| "cc_CruiseControl/Model" | 14 |

1.1.2. Cost Breakdown Details

1.1.2.1. cc_CruiseControl

The Cost Breakdown Details section contains an image of each system, along with a table containing the cost for all blocks within each system.

1.1.2.57. *time_pressed*Figure 1.57. *time_pressed*Table 1.21. Block Cost Table - *time_pressed*

| BlockName | SelfCost | TotalCost |
|---------------------|----------|-----------|
| "OverrunProtection" | 8 | 12 |
| "Add" | 8 | 8 |
| "Switch" | 6 | 6 |
| "Delay" | 5 | 5 |

For each block within the system, the Block Cost Table displays SelfCost and TotalCost. For most blocks these cost metrics will be the same. However, for some blocks, the difference between SelfCost and TotalCost can highlight important results in the design. For example, if a subsystem contains two Lookup Table blocks with similar configuration, the blocks may share some utilities in common. In this case, the TotalCost of one Lookup Table block is equal to the SelfCost plus the cost of any shared utilities. For another example, a Subsystem block has a SelfCost of zero, and the TotalCost is equal to the cost of everything contained within the subsystem.

Data Segment Table

The data segment table displays the following information:

- Identifier — Actual global or local static variable in the generated code.
- Size — Estimated size, in bytes, of the memory required by the variable.
- Parent — Name of the struct to which the Identifier belongs. This field is empty for standalone variables.
- Source Location — List of the blocks that use this variable.

1.2. Data Segment Table

Table 1.28. Data Segment Variables

| Variable | Size (bytes) | Model Location |
|---|---------------|--|
| | | <ul style="list-style-type: none"> • cc_DriverSwRequest DetectRaiseAndOverride1/Switch • cc_ControlMode Control_Mode_StateMachine/Transitions_From_ACTIVE/throttle_pressed/Compare • cc_DriverSwRequest DetectResume/Switch • cc_ControlMode Control_Mode_StateMachine/Transitions_From_ENABLED/Compare To Constant3/Compare • cc_DriverSwRequest DetectRaiseAndOverride2/Switch • cc_DriverSwRequest DecHoldDetector/LONG • cc_DriverSwRequest DetectRaiseAndOverride/Switch • cc_DriverSwRequest IncHoldDetector/LONG • cc_DriverSwRequest DecHoldDetector/SHORT • cc_DriverSwRequest IncHoldDetector/SHORT • cc_ControlMode Control_Mode_StateMachine/Transitions_From_ACTIVE/Switch1 • cc_ControlMode Control_Mode_StateMachine/wasActiveOnce/Delay1 • cc_DriverSwRequest DecHoldDetector DetectDecrease/Delay Input1 • cc_DriverSwRequest DetectResume Detect Increase/Delay Input1 • cc_DriverSwRequest DetectRaiseAndOverride Detect Increase/Delay Input1 • cc_DriverSwRequest DetectRaiseAndOverride1 Detect Increase/Delay Input1 • cc_DriverSwRequest DetectRaiseAndOverride2 Detect Increase/Delay Input1 |
| | | cc_ControlMode Control_Mode_StateMachine/Transitions_From_ACTIVE/Subsys1 |
| "cc_DriverSwRequest_DW.If.ActiveSubsystem" | 1,000 | "cc_DriverSwRequest DecHoldDetector.If" |
| "cc_DriverSwRequest_DW.If.ActiveSubsystem_rvnm" | 1,000 | "cc_DriverSwRequest IncHoldDetector.If" |
| "cc_LightControl_B.Switch1" | 1,000 | <ul style="list-style-type: none"> • cc_LightControl/Enabled Subsystem/Switch1 • cc_DriverSwRequest DetectRaiseAndOverride1 Detect Increase/Delay Input1 • cc_LightControl/OR • cc_DriverSwRequest DetectResume Detect Increase/Delay Input1 • cc_DriverSwRequest IncHoldDetector time_pressed/Delay |
| "cc_ControlMode_DW.SwitchCase.ActiveSubsystem" | 1,000 | <ul style="list-style-type: none"> • cc_ControlMode/Target_Speed_Calculator/Switch Case • cc_ThrottleController/Switch Case • cc_DriverSwRequest DecHoldDetector DetectDecrease/FixPt Relational Operator |
| "cc_LightControl_DW.icLoad" | 1,000 | <ul style="list-style-type: none"> • cc_LightControl/Enabled Subsystem/Delay2 • cc_DriverSwRequest DetectRaiseAndOverride/Switch • cc_DriverSwRequest DetectRaiseAndOverride1/Switch |
| "cc_LightControl_DW.prev_key_position_DSTATE" | 1,000 | <ul style="list-style-type: none"> • cc_LightControl/prev_key_position • cc_DriverSwRequest DetectRaiseAndOverride/Switch1 |
| Total | 79,000 | |

See Also

"Design Cost Model Metrics" on page 52-2 | "Model Testing Metrics" (Simulink Check) | metric.Engine | execute | generateReport | getAvailableMetricIds | updateArtifacts

Fixed-Point HDL-Optimized Blocks

Choose a Block for HDL-Optimized Fixed-Point Matrix Operations

In this section...

“Define the Problem to Solve” on page 53-2

“Choose an Architecture” on page 53-3

“Linear System Solvers: Select Synchronous or Asynchronous Operation” on page 53-4

“Data Complexity” on page 53-5

“Hardware Control Signals” on page 53-5

You can use the Fixed-Point Designer HDL Support library of blocks to perform fixed-point matrix operations and generate efficient HDL code. These blocks model design patterns for systems of linear equations and core matrix operations, such as QR decomposition and singular value decomposition, for hardware-efficient implementation on FPGAs. For an introduction to these concepts, see “Factorizations” and “Singular Values”.

This topic discusses how to choose an appropriate block from the Fixed-Point Designer HDL Support library for your application.

Define the Problem to Solve

First, define the math problem that you need to solve and the algorithm to use.

Linear System Solvers

Use the Linear System Solver library of blocks to solve these systems of linear equations.

| Operation | Blocks | Description |
|------------|---|---|
| $Ax = B$ | <i>Matrix Solve Using QR Decomposition</i> blocks | Use QR decomposition to solve the system of linear equations $Ax = B$. To compute $x = A^{-1}$, set B to be the identity matrix. |
| $A'AX = B$ | <i>Matrix Solve Using Q-less QR Decomposition</i> blocks | Solve the system of linear equations $A'AX = B$ using QR decomposition, without computing Q . |
| $A'AX = B$ | <i>Matrix Solve Using Q-less QR Decomposition with Forgetting Factor</i> blocks | Solve the system of linear equations $A'AX = B$ using QR decomposition, without computing Q . A is an infinitely tall matrix representing streaming data. |

Matrix Factorizations

Use the Matrix Factorizations library of blocks to perform QR decomposition, also known as QR factorization.

| Operation | Blocks | Description |
|---|---|---|
| QR decomposition | QR Decomposition blocks | Use QR decomposition to compute R and $C=Q'B$, where $QR=A$, where A and B are your input matrices. The least-squares solution to $Ax=B$ is $x=R \setminus C$. R is an upper-triangular matrix and Q is an orthogonal matrix. To compute $C=Q'$, set B to be the identity matrix. |
| QR decomposition without computing Q | Q-less QR Decomposition blocks | Use Q-less QR decomposition to compute the economy size upper-triangular R factor of the QR decomposition $A = QR$, without computing Q . The solution to $A'Ax = B$ is $x = R \setminus R'b$. |
| QR decomposition without computing Q and an infinite number of rows | Q-less QR Decomposition with Forgetting Factor blocks | Use Q-less QR decomposition to compute the economy size upper-triangular R factor of the QR decomposition $A = QR$, without computing Q . A is an infinitely tall matrix representing streaming data. |
| Singular value decomposition | Square Jacobi SVD HDL Optimized block | Use Square Jacobi SVD to compute the singular value decomposition of a square matrix A using the two-sided Jacobi algorithm. |

Choose an Architecture

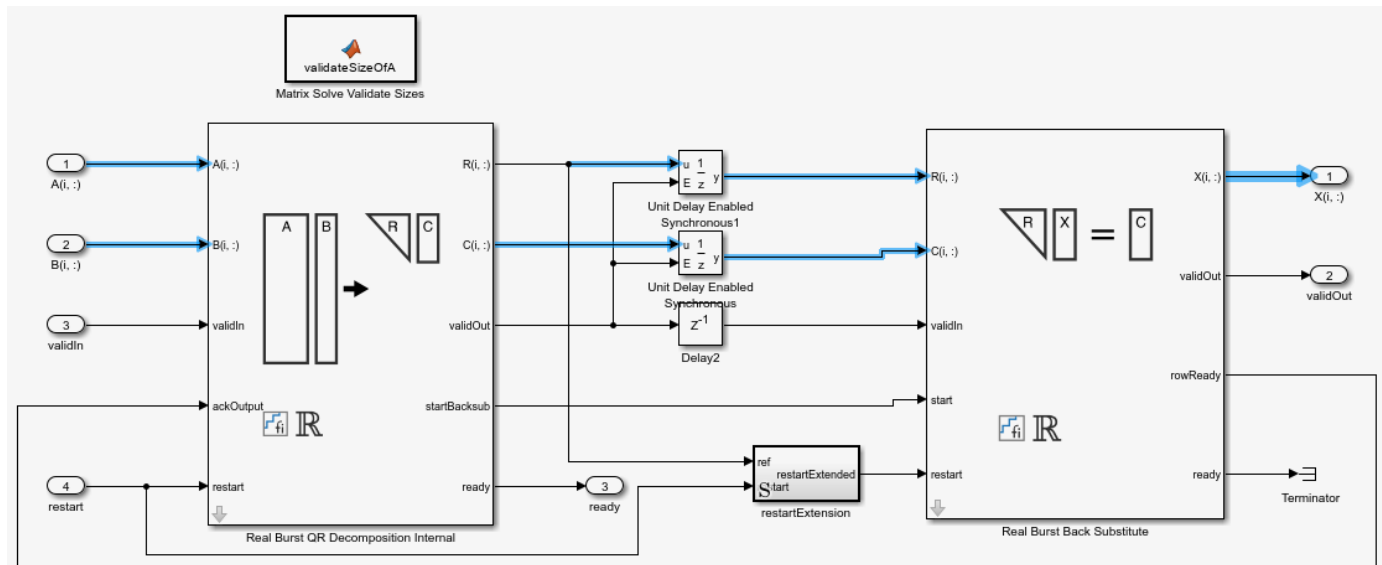
Blocks in the **Fixed-Point Designer HDL Support > Matrices and Linear Algebra** library are available in burst and partial-systolic implementations. Partial-systolic implementations prioritize speed of computations over space constraints, while burst implementations prioritize space constraints at the expense of speed of the operations. The following table illustrates the tradeoffs between the implementations available for matrix decompositions and solving systems of linear equations.

| Implementation | Ready | Latency | Area |
|------------------|--------|-----------|----------|
| Partial-Systolic | C | $O(m)$ | $O(n^2)$ |
| Burst | $O(n)$ | $O(mn^2)$ | $O(n)$ |

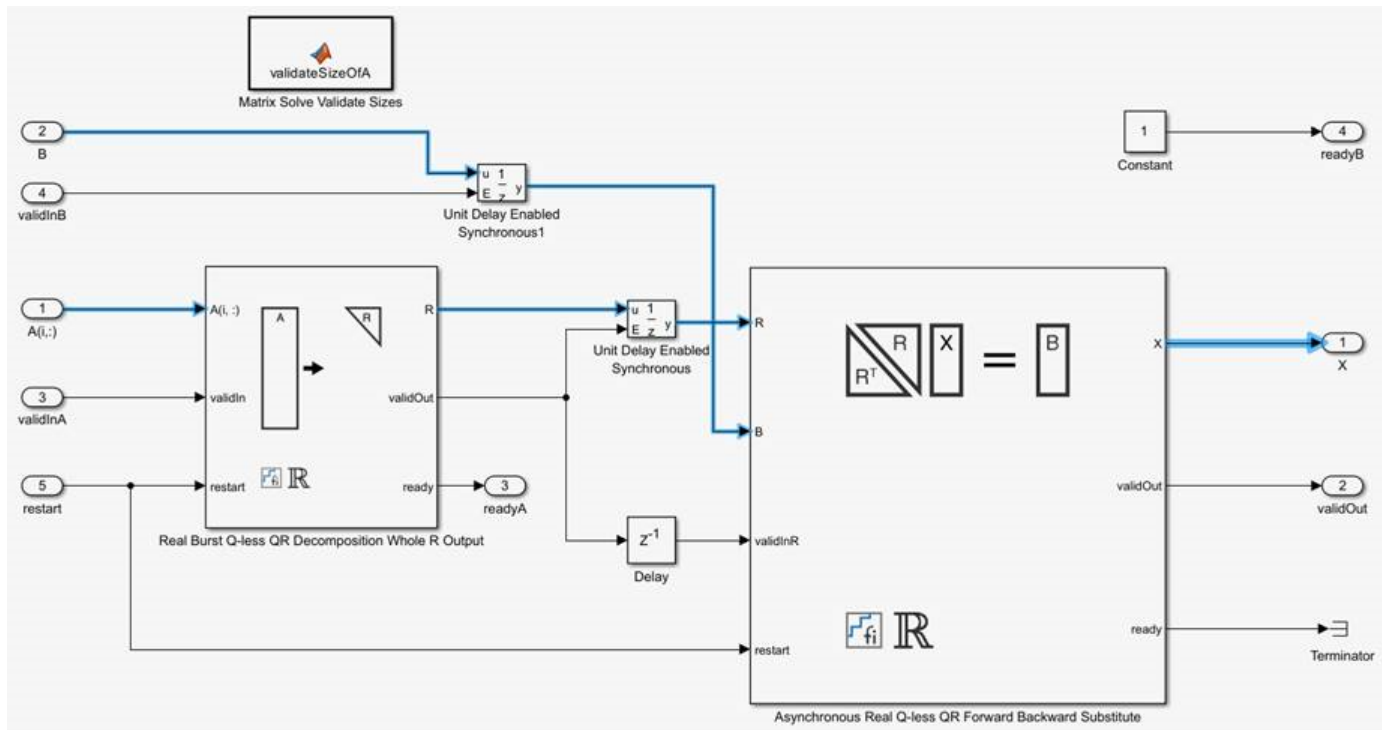
Where C is a constant proportional to the word length of the data, m is the number of rows in matrix A , and n is the number of columns in matrix A .

Linear System Solvers: Select Synchronous or Asynchronous Operation

The *Matrix Solve Using QR Decomposition* blocks operate synchronously. These blocks first decompose the input A and B matrices into R and C matrices using a QR decomposition block. Then, a back substitute block computes $RX = C$. The input A and B matrices propagate through the system in parallel, in a synchronized way.



The *Matrix Solve Using Q-less QR Decomposition* blocks operate asynchronously. First, Q-less QR decomposition is performed on the input A matrix and the resulting R matrix is put into a buffer. Then, a forward backward substitution block uses the input B matrix and the buffered R matrix to compute $R'RX = B$. Because the R and B matrices are stored separately in buffers, the upstream Q-less QR decomposition block and the downstream Forward Backward Substitute block can run independently. The Forward Backward Substitute block starts processing when the first R and B matrices are available. Then it runs continuously using the latest buffered R and B matrices, regardless of the status of the Q-less QR Decomposition block. For example, if the upstream block stops providing A and B matrices, the Forward Backward Substitute block continues to generate the same output using the last pair of R and B matrices.



The *Burst (Asynchronous) Matrix Solve Using Q-less QR Decomposition* blocks are available in both synchronous and asynchronous operation variants, as denoted by the block name.

Data Complexity

All blocks in the **Fixed-Point Designer HDL Support > Matrices and Linear Algebra** library are available in real and complex variants. Choose the real or complex variant of the block based on the complexity of your data.

Hardware Control Signals

Restart Signal

Some blocks in the **Fixed-Point Designer HDL Support > Matrices and Linear Algebra** library provide an input reset signal that clears internal states.

AMBA AXI Handshake Process

Blocks in the **Fixed-Point Designer HDL Support > Matrices and Linear Algebra** library use the AMBA AXI handshake protocol [1]. The *valid/ready* handshake process is used to transfer data and control information. This two-way control mechanism allows both the manager and subordinate to control the rate at which information moves between manager and subordinate. A *valid* signal indicates when data is available. The *ready* signal indicates that the block can accept the data. Transfer of data occurs only when both the *valid* and *ready* signals are high.

References

- [1] "AMBA AXI and ACE Protocol Specification Version E." <https://developer.arm.com/documentation/ih0022/e/AMBA-AXI3-and-AXI4-Protocol-Specification/Single-Interface-Requirements/Basic-read-and-write-transactions/Handshake-process>

See Also

Blocks

Real Burst Matrix Solve Using QR Decomposition | Real Partial-Systolic Matrix Solve Using Q-less QR Decomposition | Complex Burst Q-less QR Decomposition Whole R Output | Complex Partial-Systolic Q-less QR Decomposition with Forgetting Factor | Square Jacobi SVD HDL Optimized

Related Examples

- "Implement Hardware-Efficient Real Burst Matrix Solve Using QR Decomposition" on page 48-85
- "Implement Hardware-Efficient Real Burst Matrix Solve Using Q-less QR Decomposition with Tikhonov Regularization" on page 48-245
- "Implement Hardware-Efficient Complex Partial-Systolic QR Decomposition" on page 48-114
- "Implement Hardware-Efficient Real Burst Q-less QR with Forgetting Factor" on page 48-274
- "Algorithms to Determine Fixed-Point Types for Real Least-Squares Matrix Solve $AX=B$ " on page 48-183
- "Determine Fixed-Point Types for Real Least-Squares Matrix Solve $AX=B$ " on page 48-194

Writing Fixed-Point S-Functions

This appendix discusses the API for user-written fixed-point S-functions, which enables you to write Simulink C S-functions that directly handle fixed-point data types. Note that the API also provides support for standard floating-point and integer data types. You can find the files and examples associated with this API in the following locations:

- `matlabroot/simulink/include/`
- `matlabroot/toolbox/simulink/fixedandfloat/fixpdemos/`

Data Type Support

| In this section... |
|---|
| “Supported Data Types” on page A-2 |
| “The Treatment of Integers” on page A-2 |
| “Data Type Override” on page A-3 |

Supported Data Types

The API for user-written fixed-point S-functions provides support for a variety of Simulink and Fixed-Point Designer data types, including

- Built-in Simulink data types
 - `single`
 - `double`
 - `uint8`
 - `int8`
 - `uint16`
 - `int16`
 - `uint32`
 - `int32`
 - `uint64`
 - `int64`
- Fixed-point Simulink data types, such as
 - `sfix16_En15`
 - `ufix32_En16`
 - `ufix128`
 - `sfix37_S3_B5`
- Data types resulting from a data type override with Scaled `double`, such as
 - `flts16`
 - `flts16_En15`
 - `fltu32_S3_B5`

For more information, see “Fixed-Point Data Type and Scaling Notation” on page 35-13.

The Treatment of Integers

The API treats integers as fixed-point numbers with trivial scaling. In [Slope Bias] representation, fixed-point numbers are represented as

$$\text{real-world value} = (\text{slope} \times \text{integer}) + \text{bias}.$$

In the trivial case, $\text{slope} = 1$ and $\text{bias} = 0$.

In terms of binary-point-only scaling, the binary point is to the right of the least significant bit for trivial scaling, meaning that the fraction length is zero:

$$\text{real-world value} = \text{integer} \times 2^{-\text{fraction length}} = \text{integer} \times 2^0.$$

In either case, trivial scaling means that the real-world value is equal to the stored integer value:

$$\text{real-world value} = \text{integer}.$$

All integers, including Simulink built-in integers such as `uint8`, are treated as fixed-point numbers with trivial scaling by this API. However, Simulink built-in integers are different in that their use does not cause a Fixed-Point Designer software license to be checked out.

Data Type Override

The Fixed-Point Tool enables you to perform various data type overrides on fixed-point signals in your simulations. This API can handle signals whose data types have been overridden in this way:

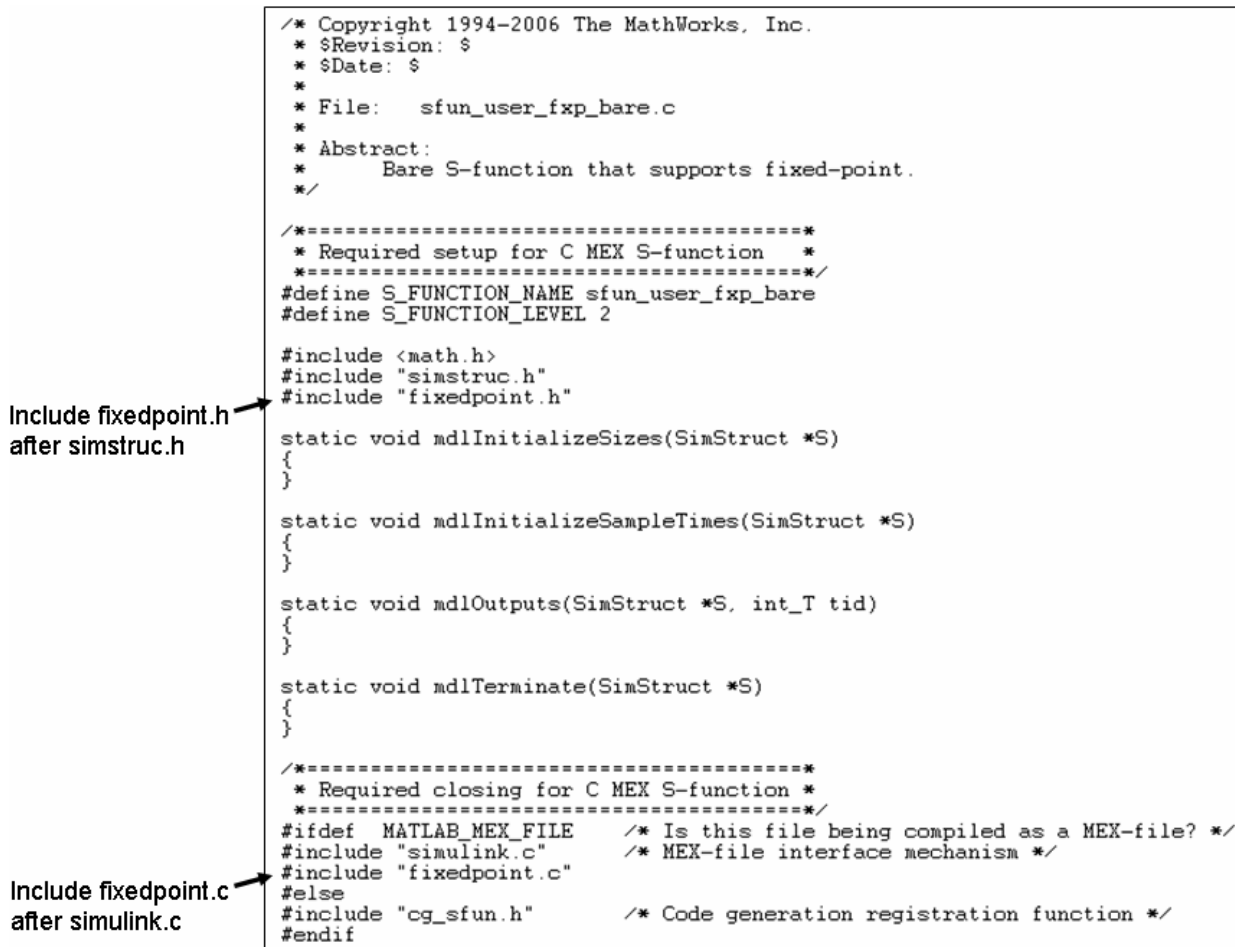
- A signal that has been overridden with `Single` is treated as a Simulink built-in `single`.
- A signal that has been overridden with `Double` is treated as a Simulink built-in `double`.
- A signal that has been overridden with `Scaled double` is treated as being of data type `ScaledDouble`.

`ScaledDouble` signals are a hybrid between floating-point and fixed-point signals, in that they are stored as `doubles` with the scaling, sign, and word length information retained. The value is stored as a floating-point `double`, but as with a fixed-point number, the distinction between the stored integer value and the real-world value remains. The scaling information is applied to the stored integer `double` to obtain the real-world value. By storing the value in a `double`, overflow and precision issues are almost always eliminated. Refer to any individual API function reference page at the end of this appendix to learn how that function treats `ScaledDouble` signals.

For more information about the Fixed-Point Tool and data type override, see **Fixed-Point Tool**.

Structure of the S-Function

The following diagram shows the basic structure of an S-function that directly handles fixed-point data types.



The callouts in the diagram alert you to the fact that you must include `fixedpoint.h` and `fixedpoint.c` at the appropriate places in the S-function. The other elements of the S-function displayed in the diagram follow the standard requirements for S-functions.

To learn how to create a MEX-file for your user-written fixed-point S-function, see “Create MEX-Files” on page A-16.

Storage Containers

In this section...

“Introduction” on page A-5

“Storage Containers in Simulation” on page A-5

“Storage Containers in Code Generation” on page A-7

Introduction

While coding with the API for user-written fixed-point S-functions, it is important to keep in mind the difference between storage container size, storage container word length, and signal word length. The sections that follow discuss the containers used by the API to store signals in simulation and code generation.

Storage Containers in Simulation

In simulation, signals are stored in one of several types of containers of a specific size.

Storage Container Categories

During simulation, fixed-point signals are held in one of the types of storage containers, as shown in the following table. In many cases, signals are represented in containers with more bits than their specified word length.

Fixed-Point Storage Containers

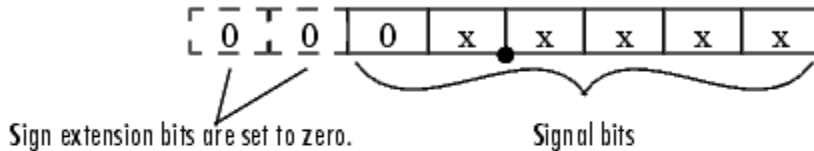
| Container Category | Signal Word Length | Container Word Length | Container Size |
|---|--|---|---|
| FXP_STORAGE_INT8 (signed) FXP_STORAGE_UINT8 (unsigned) | 1 to 8 bits | 8 bits | 1 byte |
| FXP_STORAGE_INT16 (signed) FXP_STORAGE_UINT16 (unsigned) | 9 to 16 bits | 16 bits | 2 bytes |
| FXP_STORAGE_INT32 (signed) FXP_STORAGE_UINT32 (unsigned) | 17 to 32 bits | 32 bits | 4 bytes |
| FXP_STORAGE_OTHER_SINGLE_WORD | 33 to word length of long data type | Length of long data type | Length of long data type |
| FXP_STORAGE_MULTIWORD | Greater than the word length of long data type to 128 bits | Multiples of length of long data type to 128 bits | Multiples of length of long data type to 128 bits |

When the number of bits in the signal word length is less than the size of the container, the word length bits are always stored in the least significant bits of the container. The remaining container bits must be sign extended:

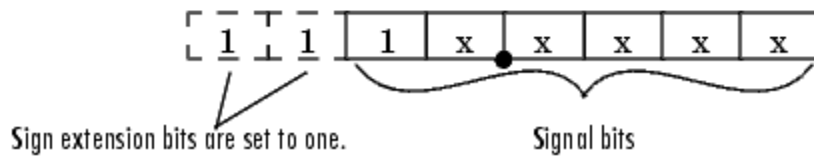
- If the data type is unsigned, the sign extension bits must be cleared to zero.
- If the data type is signed, the sign extension bits must be set to one for strictly negative numbers, and cleared to zero otherwise.

For example, a signal of data type `sfixed64` is held in a `FXP_STORAGE_INT8` container. The signal is held in the six least significant bits. The remaining two bits are set to zero when the signal is positive or zero, and to one when it is negative.

8-bit container for a signed, 6-bit signal that is positive or zero

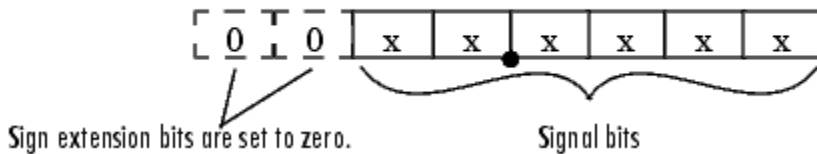


8-bit container for a signed, 6-bit signal that is negative



A signal of data type `ufixed64` is held in a `FXP_STORAGE_UINT8` container. The signal is held in the six least significant bits. The remaining two bits are always cleared to zero.

8-bit container for an unsigned, 6-bit signal



The signal and storage container word lengths are returned by the `ssGetDataTypeInfoWordLength` and `ssGetDataTypeInfoContainerWordLen` functions, respectively. The storage container size is returned by the `ssGetDataTypeInfoStorageContainerSize` function. The container category is returned by the `ssGetDataTypeInfoStorageContainerCat` function, which in addition to those in the table above, can also return the following values.

Other Storage Containers

| Container Category | Description |
|---------------------------------------|--|
| <code>FXP_STORAGE_UNKNOWN</code> | Returned if the storage container category is unknown |
| <code>FXP_STORAGE_SINGLE</code> | The container type for a Simulink single |
| <code>FXP_STORAGE_DOUBLE</code> | The container type for a Simulink double |
| <code>FXP_STORAGE_SCALEDDOUBLE</code> | The container type for a data type that has been overridden with Scaled double |

Storage Containers in Simulation Example

An `sfix24_En10` data type has a word length of 24, but is actually stored in 32 bits during simulation. For this signal,

- `ssGetDataTypeInfoStorageContainerCat` returns `FXP_STORAGE_INT32`.
- `ssGetDataTypeInfoStorageContainerSize` or `sizeof()` returns 4, which is the storage container size in bytes.
- `ssGetDataTypeInfoFxpContainerWordLen` returns 32, which is the storage container word length in bits.
- `ssGetDataTypeInfoFxpWordLength` returns 24, which is the data type word length in bits.

Storage Containers in Code Generation

The storage containers used by this API for code generation are not always the same as those used for simulation. During code generation, a native C data type is always used. Floating-point data types are held in C `double` or `float`. Fixed-point data types are held in C signed and unsigned `char`, `short`, `int`, or `long`.

Emulation

Because it is valuable for rapid prototyping and hardware-in-the-loop testing, the emulation of smaller signals inside larger containers is supported in code generation. For example, a 29-bit signal is supported in code generation if there is a C data type available that has at least 32 bits. The rules for placing a smaller signal into a larger container, and for dealing with the extra container bits, are the same in code generation as for simulation.

If a smaller signal is emulated inside a larger storage container in simulation, it is not necessarily emulated in code generation. For example, a 24-bit signal is emulated in a 32-bit storage container in simulation. However, some DSP chips have native support for 24-bit quantities. On such a target, the C compiler can define an `int` or a `long` to be exactly 24 bits. In this case, the 24-bit signal is held in a 32-bit container in simulation, and in a 24-bit container in code generation.

Conversely, a signal that was not emulated in simulation might need to be emulated in code generation. For example, some DSP chips have minimal support for integers. On such chips, `char`, `short`, `int`, and `long` might all be defined to 32 bits. In that case, it is necessary to emulate 8- and 16-bit fixed-point data types in code generation.

Storage Container TLC Functions

Since the mapping of storage containers in simulation to storage containers in code generation is not one-to-one, the Target Language Compiler (TLC) functions for storage containers are different from those in simulation:

- `FixPt_DataTypeInfoNativeType`
- `FixPt_DataTypeInfoStorageDouble`
- `FixPt_DataTypeInfoStorageSingle`
- `FixPt_DataTypeInfoStorageScaledDouble`
- `FixPt_DataTypeInfoStorageSInt`
- `FixPt_DataTypeInfoStorageUInt`

- `FixPt_DataTypeStorageSLong`
- `FixPt_DataTypeStorageULong`
- `FixPt_DataTypeStorageSShort`
- `FixPt_DataTypeStorageUShort`
- `FixPt_DataTypeStorageMultiword`

The first of these TLC functions, `FixPt_DataTypeNativeType`, is the closest analogue to `ssGetDataTypeStorageContainCat` in simulation. `FixPt_DataTypeNativeType` returns a TLC string that specifies the type of the storage container, and the Simulink Coder product automatically inserts a `typedef` that maps the string to a native C data type in the generated code.

For example, consider a fixed-data type that is held in `FXP_STORAGE_INT8` in simulation. `FixPt_DataTypeNativeType` will return `int8_T`. The `int8_T` will be `typedef`'d to a `char`, `short`, `int`, or `long` in the generated code, depending upon what is appropriate for the target compiler.

The remaining TLC functions listed above return `TRUE` or `FALSE` depending on whether a particular standard C data type is used to hold a given API-registered data type. Note that these functions do not necessarily give mutually exclusive answers for a given registered data type, due to the fact that C data types can potentially overlap in size. In C,

`sizeof(char) ≤ sizeof(short) ≤ sizeof(int) ≤ sizeof(long)`.

One or more of these C data types can be, and very often are, the same size.

Data Type IDs

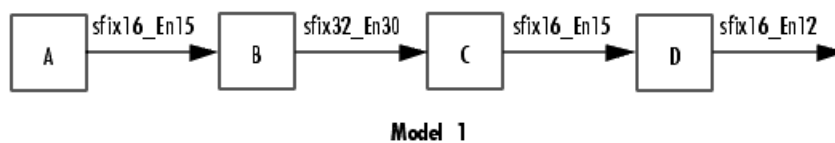
In this section...

“The Assignment of Data Type IDs” on page A-9
 “Registering Data Types” on page A-10
 “Setting and Getting Data Types” on page A-11
 “Getting Information About Data Types” on page A-11
 “Converting Data Types” on page A-13

The Assignment of Data Type IDs

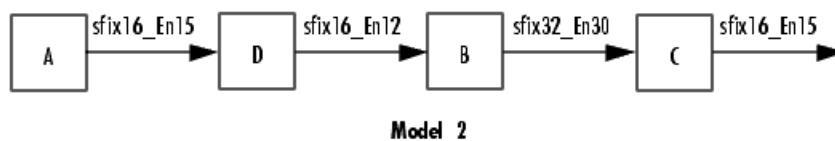
Each data type used in your S-function is assigned a data type ID. You should always use data type IDs to get and set information about data types in your S-function.

In general, the Simulink software assigns data type IDs during model initialization on a “first come, first served” basis. For example, consider the generalized schema of a block diagram below.



The Simulink software assigns a data type ID for each output data type in the diagram in the order it is requested. For simplicity, assume that the order of request occurs from left to right. Therefore, the output of block A may be assigned data type ID 13, and the output of block B may be assigned data type ID 14. The output data type of block C is the same as that of block A, so the data type ID assigned to the output of block C is also 13. The output of block D is assigned data type ID 15.

Now if the blocks in the model are rearranged,



The Simulink software still assigns the data type IDs in the order in which they are used. Therefore each data type might end up with a different data type ID. The output of block A is still assigned data type ID 13. The output of block D is now next in line and is assigned data type ID 14. The output of block B is assigned data type ID 15. The output data type of block C is still the same as that of block A, so it is also assigned data type ID 13.

This table summarizes the two cases described above.

| Block | Data Type ID in Model_1 | Data Type ID in Model_2 |
|-------|-------------------------|-------------------------|
| A | 13 | 13 |

| Block | Data Type ID in Model_1 | Data Type ID in Model_2 |
|-------|-------------------------|-------------------------|
| B | 14 | 15 |
| C | 13 | 13 |
| D | 15 | 14 |

This example illustrates that there is no strict relationship between the attributes of a data type and the value of its data type ID. In other words, the data type ID is not assigned based on the characteristics of the data type it is representing, but rather on when that data type is first needed.

Note Because of the nature of the assignment of data type IDs, you should always use API functions to extract information from a data type ID about a data type in your S-function.

Registering Data Types

The functions in the following table are available in the API for user-written fixed-point S-functions for registering data types in simulation. Each of these functions will return a data type ID. To see an example of a function being used, go to the file and line indicated in the table.

Data Type Registration Functions

| Function | Description | Example of Use |
|--|--|---|
| <code>ssRegisterDataTypeFxpBinaryPoint</code> | Register a fixed-point data type with binary-point-only scaling and return its data type ID | <code>sfun_user_fxp_asr.c</code> Line 252 |
| <code>ssRegisterDataTypeFxpFSlopeFixExpBias</code> | Register a fixed-point data type with [Slope Bias] scaling specified in terms of fractional slope, fixed exponent, and bias, and return its data type ID | Not Available |
| <code>ssRegisterDataTypeFxpScaledDouble</code> | Register a scaled double data type with [Slope Bias] scaling specified in terms of fractional slope, fixed exponent, and bias, and return its data type ID | Not Available |
| <code>ssRegisterDataTypeFxpSlopeBias</code> | Register a data type with [Slope Bias] scaling and return its data type ID | <code>sfun_user_fxp_dtprop.c</code> Line 319 |

Preassigned Data Type IDs

The Simulink software registers its built-in data types, and those data types always have preassigned data type IDs. The built-in data type IDs are given by the following tokens:

- `SS_DOUBLE`
- `SS_SINGLE`

- SS_INT8
- SS_UINT8
- SS_INT16
- SS_UINT16
- SS_INT32
- SS_UINT32
- SS_BOOLEAN

You do not need to register these data types. If you attempt to register a built-in data type, the registration function simply returns the preassigned data type ID.

Setting and Getting Data Types

Data type IDs are used to specify the data types of input and output ports, run-time parameters, and DWork states. To set fixed-point data types for quantities in your S-function, the procedure is as follows:

- 1 Register a data type using one of the functions listed in the table Data Type Registration Functions. A data type ID is returned to you.

Alternately, you can use one of the preassigned data type IDs of the Simulink built-in data types.

- 2 Use the data type ID to set the data type for an input or output port, run-time parameter, or DWork state using one of the following functions:

- `ssSetInputPortDataType`
- `ssSetOutputPortDataType`
- `ssSetRunTimeParamInfo`
- `ssSetDWorkDataType`

To get the data type ID of an input or output port, run-time parameter, or DWork state, use one of the following functions:

- `ssGetInputPortDataType`
- `ssGetOutputPortDataType`
- `ssGetSFcnParamDataType` or `ssGetRunTimeParamInfo`
- `ssGetDWorkDataType`

Getting Information About Data Types

You can use data type IDs with functions to get information about the built-in and registered data types in your S-function. The functions in the following tables are available in the API for extracting information about registered data types. To see an example of a function being used, go to the file and line indicated in the table. Note that data type IDs can also be used with all the standard data type access methods in `simstruc.h`, such as `ssGetDataTypeSize`.

Storage Container Information Functions

| Function | Description | Example of Use |
|-----------------------------------|---|---|
| ssGetDataTypeFxpContainWordLen | Return the word length of the storage container of a registered data type | sfun_user_fxp_ContainWordLenProbe.c Line 181 |
| ssGetDataTypeStorageContainCat | Return the storage container category of a registered data type | sfun_user_fxp_asr.c Line 294 |
| ssGetDataTypeStorageContainerSize | Return the storage container size of a registered data type | sfun_user_fxp_StorageContainSizeProbe.c Line 171 |

Signal Data Type Information Functions

| Function | Description | Example of Use |
|--------------------------------|---|---|
| ssGetDataTypeFxpIsSigned | Determine whether a fixed-point registered data type is signed or unsigned | sfun_user_fxp_asr.c Line 254 |
| ssGetDataTypeFxpWordLength | Return the word length of a fixed-point registered data type | sfun_user_fxp_asr.c Line 255 |
| ssGetDataTypeIsFixedPoint | Determine whether a registered data type is a fixed-point data type | sfun_user_fxp_const.c Line 127 |
| ssGetDataTypeIsFloatingPoint | Determine whether a registered data type is a floating-point data type | sfun_user_fxp_IsFloatingPointProbe.c Line 176 |
| ssGetDataTypeIsFxpFltApiCompat | Determine whether a registered data type is supported by the API for user-written fixed-point S-functions | sfun_user_fxp_asr.c Line 184 |
| ssGetDataTypeIsScalingPow2 | Determine whether a registered data type has power-of-two scaling | sfun_user_fxp_asr.c Line 203 |
| ssGetDataTypeIsScalingTrivial | Determine whether the scaling of a registered data type is slope = 1, bias = 0 | sfun_user_fxp_IsScalingTrivialProbe.c Line 171 |

Signal Scaling Information Functions

| Function | Description | Example of Use |
|-----------------------------|--|------------------------------------|
| ssGetDataTypeBias | Return the bias of a registered data type | sfun_user_fxp_dtprop.c Line 243 |
| ssGetDataTypeFixedExponent | Return the exponent of the slope of a registered data type | sfun_user_fxp_dtprop.c Line 237 |
| ssGetDataTypeFracSlope | Return the fractional slope of a registered data type | sfun_user_fxp_dtprop.c Line 234 |
| ssGetDataTypeFractionLength | Return the fraction length of a registered data type with power-of-two scaling | sfun_user_fxp_asr.c Line 256 |
| ssGetDataTypeTotalSlope | Return the total slope of the scaling of a registered data type | sfun_user_fxp_dtprop.c Line 240 |

Converting Data Types

The functions in the following table allow you to convert values between registered data types in your fixed-point S-function.

Data Type Conversion Functions

| Function | Description | Example of Use |
|--------------------------------|--|----------------|
| ssFxpConvert | Convert a value from one data type to another data type. | Not Available |
| ssFxpConvertFromRealWorldValue | Convert a value of data type <code>double</code> to another data type. | Not Available |
| ssFxpConvertToRealWorldValue | Convert a value of any data type to a <code>double</code> . | Not Available |

Overflow Handling and Rounding Methods

| |
|--|
| In this section... |
| “Tokens for Overflow Handling and Rounding Methods” on page A-14 |
| “Overflow Logging Structure” on page A-14 |

Tokens for Overflow Handling and Rounding Methods

The API for user-written fixed-point S-functions provides functions for some mathematical operations, such as conversions. When these operations are performed, a loss of precision or overflow may occur. The tokens in the following tables allow you to control the way an API function handles precision loss and overflow. The data type of the overflow handling methods is `fxpModeOverflow`. The data type of the rounding modes is `fxpModeRounding`.

Overflow Handling Tokens

| Token | Description |
|------------------------------------|--------------------|
| <code>FXP_OVERFLOW_SATURATE</code> | Saturate overflows |
| <code>FXP_OVERFLOW_WRAP</code> | Wrap overflows |

Rounding Method Tokens

| Token | Description |
|-----------------------------------|---|
| <code>FXP_ROUND_CEIL</code> | Round to the closest representable number in the direction of positive infinity |
| <code>FXP_ROUND_CONVERGENT</code> | Round toward nearest integer with ties rounding to nearest even integer |
| <code>FXP_ROUND_FLOOR</code> | Round to the closest representable number in the direction of negative infinity |
| <code>FXP_ROUND_NEAR</code> | Round to the closest representable number, with the exact midpoint rounded in the direction of positive infinity |
| <code>FXP_ROUND_NEAR_ML</code> | Round toward nearest. Ties round toward negative infinity for negative numbers, and toward positive infinity for positive numbers |
| <code>FXP_ROUND_SIMPLEST</code> | Automatically chooses between round toward floor and round toward zero to produce generated code that is as efficient as possible |
| <code>FXP_ROUND_ZERO</code> | Round to the closest representable number in the direction of zero |

Overflow Logging Structure

Math functions of the API, such as `ssFxpConvert`, can encounter overflows when carrying out an operation. These functions provide a mechanism to log the occurrence of overflows and to report that log back to the caller.

You can use a fixed-point overflow logging structure in your S-function by defining a variable of data type `fxpOverflowLogs`. Some API functions, such as `ssFxpConvert`, accept a pointer to this

structure as an argument. The function initializes the logging structure and maintains a count of each the following events that occur while the function is being performed:

- Overflows
- Saturations
- Divide-by-zeros

When a function that accepts a pointer to the logging structure is invoked, the function initializes the event counts of the structure to zero. The requested math operations are then carried out. Each time an event is detected, the appropriate event count is incremented by one.

The following fields contain the event-count information of the structure:

- `OverflowOccurred`
- `SaturationOccurred`
- `DivisionByZeroOccurred`

Create MEX-Files

To create a MEX-file for a user-written fixed-point C S-function on either Windows or UNIX® systems:

- In your S-function, include `fixedpoint.c` and `fixedpoint.h`. For more information, see “Structure of the S-Function” on page A-4.
- Pass an extra argument, `-lfixedpoint`, to the `mex` command. For example,

```
mex('sfun_user_fxp_asr.c', '-lfixedpoint')
```


Fixed-Point S-Function Examples

The following files in *matlabroot/toolbox/simulink/fixedandfloat/fixpdemos/* are examples of S-functions written with the API for user-written fixed-point S-functions:

- `sfun_user_fxp_asr.c`
- `sfun_user_fxp_BiasProbe.c`
- `sfun_user_fxp_const.c`
- `sfun_user_fxp_ContainWordLenProbe.c`
- `sfun_user_fxp_dtprop.c`
- `sfun_user_fxp_FixedExponentProbe.c`
- `sfun_user_fxp_FracLengthProbe.c`
- `sfun_user_fxp_FracSlopeProbe.c`
- `sfun_user_fxp_IsFixedPointProbe.c`
- `sfun_user_fxp_IsFloatingPointProbe.c`
- `sfun_user_fxp_IsFxpFltApiCompatProbe.c`
- `sfun_user_fxp_IsScalingPow2Probe.c`
- `sfun_user_fxp_IsScalingTrivialProbe.c`
- `sfun_user_fxp_IsSignedProbe.c`
- `sfun_user_fxp_prodsum.c`
- `sfun_user_fxp_StorageContainCatProbe.c`
- `sfun_user_fxp_StorageContainSizeProbe.c`
- `sfun_user_fxp_TotalSlopeProbe.c`
- `sfun_user_fxp_U32BitRegion.c`
- `sfun_user_fxp_WordLengthProbe.c`

See Also

Related Examples

- “Get the Input Port Data Type” on page A-18
- “Set the Output Port Data Type” on page A-20
- “Interpret an Input Value” on page A-21
- “Write an Output Value” on page A-23
- “Determine Output Type Using the Input Type” on page A-25

Get the Input Port Data Type

Within your S-function, you might need to know the data types of different ports, run-time parameters, and DWorks. In each case, you will need to get the data type ID of the data type, and then use functions from this API to extract information about the data type.

For example, suppose you need to know the data type of your input port. To do this,

- 1 Use `ssGetInputPortDataType`. The data type ID of the input port is returned.
- 2 Use API functions to extract information about the data type.

The following lines of example code are from `sfun_user_fxp_dtprop.c`.

In lines 191 and 192, `ssGetInputPortDataType` is used to get the data type ID for the two input ports of the S-function:

```
dataTypeIdU0 = ssGetInputPortDataType( S, 0 );
dataTypeIdU1 = ssGetInputPortDataType( S, 1 );
```

Further on in the file, the data type IDs are used with API functions to get information about the input port data types. In lines 205 through 226, a check is made to see whether the input port data types are `single` or `double`:

```
storageContainerU0 = ssGetDataTypeInfoStorageContainer( S,
dataTypeIdU0 );
storageContainerU1 = ssGetDataTypeInfoStorageContainer( S,
dataTypeIdU1 );
if ( storageContainerU0 == FXP_STORAGE_DOUBLE ||
storageContainerU1 == FXP_STORAGE_DOUBLE )
{
/* Doubles take priority over all other rules.
* If either of first two inputs is double,
* then third input is set to double.
*/
dataTypeIdU2Desired = SS_DOUBLE;
}
else if ( storageContainerU0 == FXP_STORAGE_SINGLE ||
storageContainerU1 == FXP_STORAGE_SINGLE )
{
/* Singles take priority over all other rules,
* except doubles.
* If either of first two inputs is single
* then third input is set to single.
*/
dataTypeIdU2Desired = SS_SINGLE;
}
else
```

In lines 227 through 244, additional API functions are used to get information about the data types if they are neither `single` nor `double`:

```
{
isSignedU0 = ssGetDataTypeInfoIsSigned( S, dataTypeIdU0 );
isSignedU1 = ssGetDataTypeInfoIsSigned( S, dataTypeIdU1 );

wordLengthU0 = ssGetDataTypeInfoWordLength( S, dataTypeIdU0 );
wordLengthU1 = ssGetDataTypeInfoWordLength( S, dataTypeIdU1 );
```

```
fracSlopeU0 = ssGetDataTypeInfoFracSlope( S, dataTypeIdU0 );
fracSlopeU1 = ssGetDataTypeInfoFracSlope( S, dataTypeIdU1 );

fixedExponentU0 = ssGetDataTypeInfoFixedExponent( S, dataTypeIdU0 );
fixedExponentU1 = ssGetDataTypeInfoFixedExponent( S, dataTypeIdU1 );

totalSlopeU0 = ssGetDataTypeInfoTotalSlope( S, dataTypeIdU0 );
totalSlopeU1 = ssGetDataTypeInfoTotalSlope( S, dataTypeIdU1 );

biasU0 = ssGetDataTypeInfoBias( S, dataTypeIdU0 );
biasU1 = ssGetDataTypeInfoBias( S, dataTypeIdU1 );
}
```

The functions used above return whether the data types are signed or unsigned, as well as their word lengths, fractional slopes, exponents, total slopes, and biases. Together, these quantities give full information about the fixed-point data types of the input ports.

See Also

Related Examples

- “Set the Output Port Data Type” on page A-20

Set the Output Port Data Type

You may want to set the data type of various ports, run-time parameters, or DWorks in your S-function.

For example, suppose you want to set the output port data type of your S-function. To do this,

- 1 Register a data type by using one of the functions listed in the table Data Type Registration Functions. A data type ID is returned.

Alternately, you can use one of the predefined data type IDs of the Simulink built-in data types.

- 2 Use `ssSetOutputPortDataType` with the data type ID from Step 1 to set the output port to the desired data type.

In the example below from lines 336 - 352 of `sfun_user_fxp_const.c`, `ssRegisterDataTypeFxpBinaryPoint` is used to register the data type. `ssSetOutputPortDataType` then sets the output data type either to the given data type ID, or to be dynamically typed:

```
/* Register data type
   */
if ( notSizesOnlyCall )
{
    DTypeId DataTypeId = ssRegisterDataTypeFxpBinaryPoint(
        S,
        V_ISSIGNED,
        V_WORDLENGTH,
        V_FRACTIONLENGTH,
        1 /* true means obey data type override setting for
           this subsystem */ );

    ssSetOutputPortDataType( S, 0, DataTypeId );
}
else
{
    ssSetOutputPortDataType( S, 0, DYNAMICALLY_TYPED );
}
```

See Also

Related Examples

- “Interpret an Input Value” on page A-21

Interpret an Input Value

Suppose you need to get the value of the signal on your input port to use in your S-function. You should write your code so that the pointer to the input value is properly typed, so that the values read from the input port are interpreted correctly. To do this, you can use these steps, which are shown in the example code below:

- 1 Create a void pointer to the value of the input signal.
- 2 Get the data type ID of the input port using `ssGetInputPortDataType`.
- 3 Use the data type ID to get the storage container type of the input.
- 4 Have a case for each input storage container type you want to handle. Within each case, you will need to perform the following in some way:
 - Create a pointer of the correct type according to the storage container, and cast the original void pointer into the new fully typed pointer (see **a** and **c**).
 - You can now store and use the value by dereferencing the new, fully typed pointer (see **b** and **d**).

For example,

```
static void mdlOutputs(SimStruct *S, int_T tid)
{
    const void *pVoidIn =
        (const void *)ssGetInputPortSignal( S, 0 ); (1)

    DTypeId dataTypeIdU0 = ssGetInputPortDataType( S, 0 ); (2)

    fxpStorageContainerCategory storageContainerU0 =
        ssGetDataTypeIdStorageContainerCat( S, dataTypeIdU0 ); (3)

    switch ( storageContainerU0 )
    {
        case FXP_STORAGE_UINT8: (4)
        {
            const uint8_T *pU8_Properly_Typed_Pointer_To_U0; (a)

            uint8_T u8_Stored_Integer_U0; (b)

            pU8_Properly_Typed_Pointer_To_U0 =
                (const uint8_T *)pVoidIn; (c)

            u8_Stored_Integer_U0 =
                *pU8_Properly_Typed_Pointer_To_U0; (d)

            <snip: code that uses input when it's in a uint8_T>
        }
        break;

        case FXP_STORAGE_INT8: (4)
        {
            const int8_T *pS8_Properly_Typed_Pointer_To_U0; (a)

            int8_T s8_Stored_Integer_U0; (b)

            pS8_Properly_Typed_Pointer_To_U0 =
```

```
        (const int8_T *)pVoidIn; (c)
s8_Stored_Integer_U0 =
    *pS8_Properly_Typed_Pointer_To_U0; (d)
    <snip: code that uses input when it's in a int8_T>
}
break;
```

See Also

Related Examples

- “Write an Output Value” on page A-23

Write an Output Value

Suppose you need to write the value of the output signal to the output port in your S-function. You should write your code so that the pointer to the output value is properly typed. To do this, you can use these steps, which are followed in the example code below:

- 1 Create a void pointer to the value of the output signal.
- 2 Get the data type ID of the output port using `ssGetOutputPortDataType`.
- 3 Use the data type ID to get the storage container type of the output.
- 4 Have a case for each output storage container type you want to handle. Within each case, you will need to perform the following in some way:
 - Create a pointer of the correct type according to the storage container, and cast the original void pointer into the new fully typed pointer (see **a** and **c**).
 - You can now write the value by dereferencing the new, fully typed pointer (see **b** and **d**).

For example,

```
static void mdlOutputs(SimStruct *S, int_T tid)
{
    <snip>

    void *pVoidOut = ssGetOutputPortSignal( S, 0 ); (1)

    DTypeId dataTypeIdY0 = ssGetOutputPortDataType( S, 0 ); (2)

    fxpStorageContainerCategory storageContainerY0 =
        ssGetDataTypeIdStorageContainerCat( S,
        dataTypeIdY0 ); (3)

    switch ( storageContainerY0 )
    {
        case FXP_STORAGE_UINT8: (4)
        {
            const uint8_T *pU8_Properly_Typed_Pointer_To_Y0; (a)

            uint8_T u8_Stored_Integer_Y0; (b)

            <snip: code that puts the desired output stored integer
            value in to temporary variable u8_Stored_Integer_Y0>

            pU8_Properly_Typed_Pointer_To_Y0 =
                (const uint8_T *)pVoidOut; (c)

            *pU8_Properly_Typed_Pointer_To_Y0 =
                u8_Stored_Integer_Y0; (d)

        }
        break;

        case FXP_STORAGE_INT8: (4)
        {
            const int8_T *pS8_Properly_Typed_Pointer_To_Y0; (a)

            int8_T s8_Stored_Integer_Y0; (b)

            <snip: code that puts the desired output stored integer
            value in to temporary variable s8_Stored_Integer_Y0>

            pS8_Properly_Typed_Pointer_To_Y0 =
                (const int8_T *)pVoidOut; (c)

            *pS8_Properly_Typed_Pointer_To_Y0 =
                s8_Stored_Integer_Y0; (d)

        }
    }
}
```

```
    }  
    break;  
<snip>
```

See Also

Related Examples

- “Determine Output Type Using the Input Type” on page A-25

Determine Output Type Using the Input Type

The following sample code from lines 243 through 261 of `sfun_user_fxp_asr.c` gives an example of using the data type of the input to your S-function to calculate the output data type. Notice that in this code

- The output is signed or unsigned to match the input **(a)**.
- The output is the same word length as the input **(b)**.
- The fraction length of the output depends on the input fraction length and the number of shifts **(c)**.

```
#define MDL_SET_INPUT_PORT_DATA_TYPE
static void mdlSetInputPortDataType(SimStruct *S, int port,
                                   DTypeId dataTypeIdInput)
{
    if ( isDataTypeSupported( S, dataTypeIdInput ) )
    {
        DTypeId dataTypeIdOutput;

        ssSetInputPortDataType( S, port, dataTypeIdInput );

        dataTypeIdOutput = ssRegisterDataTypeFxpBinaryPoint(
            S,
            ssGetDataTypeFxpIsSigned( S, dataTypeIdInput ), (a)
            ssGetDataTypeFxpWordLength( S, dataTypeIdInput ), (b)
            ssGetDataTypeFractionLength( S, dataTypeIdInput )
            - V_NUM_BITS_TO_SHIFT_RGHT, (c)
            0 /* false means do NOT obey data type override
               setting for this subsystem */ );

        ssSetOutputPortDataType( S, 0, dataTypeIdOutput );
    }
}
```

API Function Reference

ssFxpConvert

Convert value from one data type to another

Syntax

```
extern void ssFxpConvert (SimStruct *S,  
                          void *pVoidDest,  
                          size_t sizeofDest,  
                          DTypeId dataTypeIdDest,  
                          const void *pVoidSrc,  
                          size_t sizeofSrc,  
                          DTypeId dataTypeIdSrc,  
                          fxpModeRounding roundMode,  
                          fxpModeOverflow overflowMode,  
                          fxpOverflowLogs *pFxpOverflowLogs)
```

Arguments

S

SimStruct representing an S-function block.

pVoidDest

Pointer to the converted value.

sizeofDest

Size in memory of the converted value.

dataTypeIdDest

Data type ID of the converted value.

pVoidSrc

Pointer to the value you want to convert.

sizeofSrc

Size in memory of the value you want to convert.

dataTypeIdSrc

Data type ID of the value you want to convert.

roundMode

Rounding mode you want to use if a loss of precision is necessary during the conversion. Possible values are FXP_ROUND_CEIL, FXP_ROUND_CONVERGENT, FXP_ROUND_FLOOR, FXP_ROUND_NEAR, FXP_ROUND_NEAR_ML, FXP_ROUND_SIMPLEST and FXP_ROUND_ZERO.

overflowMode

Overflow mode you want to use if overflow occurs during the conversion. Possible values are FXP_OVERFLOW_SATURATE and FXP_OVERFLOW_WRAP.

pFxpOverflowLogs

Pointer to the fixed-point overflow logging structure.

Description

This function converts a value of any registered built-in or fixed-point data type to any other registered built-in or fixed-point data type.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

None

See Also

`ssFxpConvertFromRealWorldValue`, `ssFxpConvertToRealWorldValue`

Version History

Introduced before R2006a

ssFxpConvertFromRealWorldValue

Convert value of data type double to another data type

Syntax

```
extern void ssFxpConvertFromRealWorldValue
    (SimStruct *S,
     void *pVoidDest,
     size_t sizeofDest,
     DTypeId dataTypeIdDest,
     double dblRealWorldValue,
     fxpModeRounding roundMode,
     fxpModeOverflow overflowMode,
     fxpOverflowLogs *pFxpOverflowLogs)
```

Arguments

S

SimStruct representing an S-function block.

pVoidDest

Pointer to the converted value.

sizeofDest

Size in memory of the converted value.

dataTypeIdDest

Data type ID of the converted value.

dblRealWorldValue

Double value you want to convert.

roundMode

Rounding mode you want to use if a loss of precision is necessary during the conversion. Possible values are FXP_ROUND_CEIL, FXP_ROUND_CONVERGENT, FXP_ROUND_FLOOR, FXP_ROUND_NEAR, FXP_ROUND_NEAR_ML, FXP_ROUND_SIMPLEST and FXP_ROUND_ZERO.

overflowMode

Overflow mode you want to use if overflow occurs during the conversion. Possible values are FXP_OVERFLOW_SATURATE and FXP_OVERFLOW_WRAP.

pFxpOverflowLogs

Pointer to the fixed-point overflow logging structure.

Description

This function converts a double value to any registered built-in or fixed-point data type.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

None

See Also

`ssFxpConvert`, `ssFxpConvertToRealWorldValue`

Version History

Introduced before R2006a

ssFxpConvertToRealWorldValue

Convert value of any data type to double

Syntax

```
extern double ssFxpConvertToRealWorldValue (SimStruct *S,  
                                             const void *pVoidSrc,  
                                             size_t sizeofSrc,  
                                             DTypeId dataTypeIdSrc)
```

Arguments

S

SimStruct representing an S-function block.

pVoidSrc

Pointer to the value you want to convert.

sizeofSrc

Size in memory of the value you want to convert.

dataTypeIdSrc

Data type ID of the value you want to convert.

Description

This function converts a value of any registered built-in or fixed-point data type to a double.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

None

See Also

ssFxpConvert, ssFxpConvertFromRealWorldValue

Version History

Introduced before R2006a

ssFxpGetU32BitRegion

Return stored integer value for 32-bit region of real, scalar signal element

Syntax

```
extern uint32 ssFxpGetU32BitRegion(SimStruct *S,
                                   const void *pVoid
                                   DTypeId dataTypeId
                                   unsigned int regionIndex)
```

Arguments

S

SimStruct representing an S-function block.

pVoid

Pointer to the storage container of the real, scalar signal element in which the 32-bit region of interest resides.

dataTypeId

Data type ID of the registered data type corresponding to the signal.

regionIndex

Index of the 32-bit region whose stored integer value you want to retrieve, where 0 accesses the least significant 32-bit region.

Description

This function returns the stored integer value in the 32-bit region specified by `regionIndex`, associated with the fixed-point data type designated by `dataTypeId`. You can use this function with any fixed-point data type, including those with word sizes less than 32 bits. If the fixed-point word size is less than 32 bits, the remaining bits are sign extended.

This function generates an error if `dataTypeId` represents a floating-point data type.

To view an example model whose S-functions use the `ssFxpGetU32BitRegion` function, at the MATLAB prompt, enter `fxpdemo_sfun_user_U32BitRegion`.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see "Structure of the S-Function" on page A-4.

Languages

C

See Also

`ssFxpSetU32BitRegion`

Version History

Introduced in R2007b

ssFxpGetU32BitRegionCompliant

Determine whether S-function is compliant with the U32 bit region interface

Note The `ssFxpGetU32BitRegionCompliant` function can be ignored. This function no longer has any impact on the memory layout for inputs and outputs. The memory layout introduced in R2008a is always used.

Syntax

```
extern ssFxpSGetU32BitRegionCompliant(SimStruct *S,  
                                     int *result)
```

Arguments

S

SimStruct representing an S-function block.

result

- 1 if S-function calls `ssFxpSetU32BitRegionCompliant` to declare compliance with memory footprint for fixed-point data types with 33 or more bits
- 0 if S-function does not call `ssFxpSetU32BitRegionCompliant`

Description

This function checks whether the S-function calls `ssFxpSetU32BitRegionCompliant` to declare compliance with the memory footprint for fixed-point data types with 33 or more bits. Before calling any other Fixed-Point Designer API function on data with 33 or more bits, you must call `ssFxpSetU32BitRegionCompliant` as follows:

```
ssFxpSetU32BitRegionCompliant(S,1);
```

Note The Fixed-Point Designer software assumes that S-functions that use fixed-point data types with 33 or more bits without calling `ssFxpSetU32BitRegionCompliant` are using the obsolete memory footprint that existed until R2007b. Either redesign these S-functions or isolate them using the library `fixpt_legacy_sfun_support`.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

See Also

`ssFxpSetU32BitRegionCompliant`

Version History

Introduced in R2009a

ssFxpSetU32BitRegion

Set stored integer value for 32-bit region of real, scalar signal element

Note The `ssFxpSetU32BitRegionCompliant` function can be ignored. This function no longer has any impact on the memory layout for inputs and outputs. The memory layout introduced in R2008a is always used.

Syntax

```
extern ssFxpSetU32BitRegion(SimStruct *S,
                           void *pVoid
                           DTypeId dataTypeId
                           uint32 regionValue
                           unsigned int regionIndex)
```

Arguments

S

SimStruct representing an S-function block.

pVoid

Pointer to the storage container of the real, scalar signal element in which the 32-bit region of interest resides.

dataTypeId

Data type ID of the registered data type corresponding to the signal.

regionValue

Stored integer value that you want to assign to a 32-bit region.

regionIndex

Index of the 32-bit region whose stored integer value you want to set, where 0 accesses the least significant 32-bit region.

Description

This function sets `regionValue` as the stored integer value of the 32-bit region specified by `regionIndex`, associated with the fixed-point data type designated by `dataTypeId`. You can use this function with any fixed-point data type, including those with word sizes less than 32 bits. If the fixed-point word size is less than 32 bits, ensure that the remaining bits are sign extended.

This function generates an error if `dataTypeId` represents a floating-point data type, or if the stored integer value that you set is invalid.

To view an example model whose S-functions use the `ssFxpSetU32BitRegion` function, at the MATLAB prompt, enter `fxpdemo_sfun_user_U32BitRegion`.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

See Also

`ssFxpGetU32BitRegion`

Version History

Introduced in R2007b

ssFxpSetU32BitRegionCompliant

Declare compliance with the U32 bit region interface for fixed-point data types with 33 or more bits

Note The `ssFxpSetU32BitRegionCompliant` function can be ignored. This function no longer has any impact on the memory layout for inputs and outputs. The memory layout introduced in R2008a is always used.

Syntax

```
extern ssFxpSetU32BitRegionCompliant(SimStruct *S,
                                     int Value)
```

Arguments

S

SimStruct representing an S-function block.

Value

- 1 declare compliance with memory footprint for fixed-point data types with 33 or more bits.

Description

This function declares compliance with the Fixed-Point Designer bit region interface for data types with 33 or more bits. The memory footprint for data types with 33 or more bits varies between MATLAB host platforms and might change between software releases. To make an S-function robust to memory footprint changes, use the U32 bit region interface. You can use identical source code on different MATLAB host platforms and with any software release from R2008b. If the memory footprint changes between releases, you do not have to recompile U32 bit region compliant S-functions. To make an S-function U32 bit region compliant, before calling any other Fixed-Point Designer API function on data with 33 or more bits, you must call this function as follows:

```
ssFxpSetU32BitRegionCompliant(S,1);
```

If an S-function block contains a fixed-point data type with 33 or more bits, call this function in `mdlInitializeSizes()`.

Note The Fixed-Point Designer software assumes that S-functions that use fixed-point data types with 33 or more bits without calling `ssFxpSetU32BitRegionCompliant` are using the obsolete memory footprint that existed until R2007b. Either redesign these S-functions or isolate them using the library `fixpt_legacy_sfun_support`.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

See Also

`ssFxpGetU32BitRegionCompliant`

Version History

Introduced in R2009a

ssGetDataTypeBias

Return bias of registered data type

Syntax

```
extern double ssGetDataTypeBias(SimStruct *S, DTypeId
                               dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered data type for which you want to know the bias.

Description

Fixed-point numbers can be represented as

real-world value = (slope × integer) + bias.

This function returns the bias of a registered data type:

- For both trivial scaling and power-of-two scaling, 0 is returned.
- If the registered data type is `ScaledDouble`, the bias returned is that of the nonoverridden data type.

This function errors out when `ssGetDataTypeIsFxpFltApiCompat` returns `FALSE`.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

FixPt_DataTypeBias

See Also

ssGetDataTypeFixedExponent, ssGetDataTypeFracSlope, ssGetDataTypeTotalSlope

Version History

Introduced before R2006a

ssGetDataTypeFixedExponent

Return exponent of slope of registered data type

Syntax

```
extern int ssGetDataTypeFixedExponent (SimStruct *S, DTypeId
                                       dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered data type for which you want to know the exponent.

Description

Fixed-point numbers can be represented as

real-world value = (slope × integer) + bias,

where the slope can be expressed as

slope = fractional slope × 2^{exponent}.

This function returns the exponent of a registered fixed-point data type:

- For power-of-two scaling, the exponent is the negative of the fraction length.
- If the data type has trivial scaling, including for data types `single` and `double`, the exponent is 0.
- If the registered data type is `ScaledDouble`, the exponent returned is that of the nonoverridden data type.

This function errors out when `ssGetDataTypeIsFxpFltApiCompat` returns `FALSE`.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

FixPt_DataTypeFixedExponent

See Also

`ssGetDataTypeInfoBias`, `ssGetDataTypeInfoFracSlope`, `ssGetDataTypeInfoTotalSlope`

Version History

Introduced before R2006a

ssGetDataTypeFracSlope

Return fractional slope of registered data type

Syntax

```
extern double ssGetDataTypeFracSlope(SimStruct *S, DTypeId
                                     dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered data type for which you want to know the fractional slope.

Description

Fixed-point numbers can be represented as

real-world value = (slope × integer) + bias,

where the slope can be expressed as

slope = fractional slope × 2^{exponent}.

This function returns the fractional slope of a registered fixed-point data type. To get the total slope, use `ssGetDataTypeTotalSlope`:

- For power-of-two scaling, the fractional slope is 1.
- If the data type has trivial scaling, including data types `single` and `double`, the fractional slope is 1.
- If the registered data type is `ScaledDouble`, the fractional slope returned is that of the nonoverridden data type.

This function errors out when `ssGetDataTypeIsFxpFltApiCompat` returns `FALSE`.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

`FixPt_DataTypeFracSlope`

See Also

`ssGetDataTypeBias`, `ssGetDataTypeFixedExponent`, `ssGetDataTypeTotalSlope`

Version History

Introduced before R2006a

ssGetDataTypeFractionLength

Return fraction length of registered data type with power-of-two scaling

Syntax

```
extern int ssGetDataTypeFractionLength (SimStruct *S, DTypeId
                                       dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered data type for which you want to know the fraction length.

Description

This function returns the fraction length, or the number of bits to the right of the binary point, of the data type designated by dataTypeId.

This function errors out when ssGetDataTypeIsScalingPow2 returns FALSE.

This function also errors out when ssGetDataTypeIsFxpFltApiCompat returns FALSE.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

FixPt_DataTypeFractionLength

See Also

ssGetDataTypeFxpWordLength

Version History

Introduced before R2006a

ssGetDataTypeFxpContainWordLen

Return word length of storage container of registered data type

Syntax

```
extern int ssGetDataTypeFxpContainWordLen (SimStruct *S,  
                                           DTypeId dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered data type for which you want to know the container word length.

Description

This function returns the word length, in bits, of the storage container of the fixed-point data type designated by `dataTypeId`. This function does not return the size of the storage container or the word length of the data type. To get the storage container size, use `ssGetDataTypeStorageContainerSize`. To get the data type word length, use `ssGetDataTypeFxpWordLength`.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

Examples

An `sfix24_En10` data type has a word length of 24, but is actually stored in 32 bits during simulation. For this signal,

- `ssGetDataTypeFxpContainWordLen` returns 32, which is the storage container word length in bits.
- `ssGetDataTypeFxpWordLength` returns 24, which is the data type word length in bits.
- `ssGetDataTypeStorageContainerSize` or `sizeof()` returns 4, which is the storage container size in bytes.

TLC Functions

FixPt_DataTypeFxpContainWordLen

See Also

ssGetDataTypeFxpWordLength, ssGetDataTypeStorageContainCat,
ssGetDataTypeStorageContainerSize

Version History

Introduced before R2006a

ssGetDataTypeFxpIsSigned

Determine whether fixed-point registered data type is signed or unsigned

Syntax

```
extern int ssGetDataTypeFxpIsSigned (SimStruct *S, DTypeId  
                                     dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered fixed-point data type for which you want to know whether it is signed.

Description

This function determines whether a registered fixed-point data type is signed:

- If the fixed-point data type is signed, the function returns TRUE. If the fixed-point data type is unsigned, the function returns FALSE.
- If the registered data type is `ScaledDouble`, the function returns TRUE or FALSE according to the signedness of the nonoverridden data type.
- If the registered data type is `single` or `double`, this function errors out.

This function errors out when `ssGetDataTypeIsFxpFltApiCompat` returns FALSE.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

`FixPt_DataTypeFxpIsSigned`

Version History

Introduced before R2006a

ssGetDataTypeFxpWordLength

Return word length of fixed-point registered data type

Syntax

```
extern int ssGetDataTypeFxpWordLength (SimStruct *S, DTypeId
                                       dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered fixed-point data type for which you want to know the word length.

Description

This function returns the word length of the fixed-point data type designated by `dataTypeId`. This function does not return the word length of the container of the data type. To get the container word length, use `ssGetDataTypeFxpContainWordLen`:

- If the registered data type is fixed point, this function returns the total word length including any sign bits, integer bits, and fractional bits.
- If the registered data type is `ScaledDouble`, this function returns the word length of the nonoverridden data type.
- If registered data type is `single` or `double`, this function errors out.

This function errors out when `ssGetDataTypeIsFxpFltApiCompat` returns `FALSE`.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see "Structure of the S-Function" on page A-4.

Languages

C

Examples

An `sfix24_En10` data type has a word length of 24, but is actually stored in 32 bits during simulation. For this signal,

- `ssGetDataTypeFxpWordLength` returns 24, which is the data type word length in bits.
- `ssGetDataTypeFxpContainWordLen` returns 32, which is the storage container word length in bits.

- `ssGetDataTypeStorageContainerSize` or `sizeof()` returns 4, which is the storage container size in bytes.

TLC Functions

`FixPt_DataTypeFxpWordLength`

See Also

`ssGetDataTypeFxpContainWordLen`, `ssGetDataTypeFractionLength`,
`ssGetDataTypeStorageContainerSize`

Version History

Introduced before R2006a

ssGetDataTypesDoubleSingleOrHalf

Determine whether registered data type is double, single, or half-precision data type

Syntax

```
extern int ssGetDataTypesDoubleSingleOrHalf
        (SimStruct *S,
         DTypeId dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered data type for which you want to determine if the data type is double, single, or half-precision.

Description

This function determines whether a registered data type is a double, single, or half-precision data type.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

FixPt_DataTypesDoubleSingleOrHalf

See Also

`ssGetDataTypesHalfPrecision`, `ssGetDataTypesFloatingPoint`, `ssRegisterDataTypesHalfPrecision`

Version History

Introduced in R2020b

ssGetDataTypesFixedPoint

Determine whether registered data type is fixed-point data type

Syntax

```
extern int ssGetDataTypesFixedPoint(SimStruct *S, DTypeId  
                                   dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered data type for which you want to know whether it is fixed-point.

Description

This function determines whether a registered data type is a fixed-point data type:

- This function returns `TRUE` if the registered data type is fixed-point, and `FALSE` otherwise.
- If the registered data type is a pure Simulink integer, such as `int8`, this function returns `TRUE`.
- If the registered data type is `ScaledDouble`, this function returns `FALSE`.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

`FixPt_DataTypeIsFixedPoint`

See Also

`ssGetDataTypesFloatingPoint`

Version History

Introduced before R2006a

ssGetDataTypesFloatingPoint

Determine whether registered data type is floating-point data type

Syntax

```
extern int ssGetDataTypesFloatingPoint (SimStruct *S, DTypeId  
                                       dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered data type for which you want to know whether it is floating-point.

Description

This function determines whether a registered data type is single or double:

- If the registered data type is either single or double, this function returns TRUE, and FALSE is returned otherwise.
- If the registered data type is ScaledDouble, this function returns FALSE.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

FixPt_DataTypeIsFloatingPoint

See Also

ssGetDataTypesFixedPoint

Version History

Introduced before R2006a

ssGetDataTypesFxpFltApiCompat

Determine whether registered data type is supported by API for user-written fixed-point S-functions

Syntax

```
extern int ssGetDataTypesFxpFltApiCompat(SimStruct *S, DTypeId  
                                         dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered data type for which you want to determine compatibility with the API for user-written fixed-point S-functions.

Description

This function determines whether the registered data type is supported by the API for user-written fixed-point S-functions. The supported data types are all standard Simulink data types, all fixed-point data types, and data type override data types.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

None. Checking for API-compatible data types is done in simulation. Checking for API-compatible data types is not supported in TLC.

Version History

Introduced before R2006a

ssGetDataTypesHalfPrecision

Determine whether registered data type is half-precision data type

Syntax

```
extern int ssGetDataTypesHalfPrecision  
                (SimStruct *S,  
                DTypeId dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered data type for which you want to determine if the data type is half-precision.

Description

This function determines whether a registered data type is a half-precision data type as defined in “Half-Precision Format” on page 35-22.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

FixPt_DataTypesHalfPrecision

See Also

ssGetDataTypesDoubleSingleorHalf, ssRegisterDataTypesHalfPrecision

Version History

Introduced in R2020b

ssGetDataTypesScalingPow2

Determine whether registered data type has power-of-two scaling

Syntax

```
extern int ssGetDataTypesScalingPow2 (SimStruct *S, DTypeId
                                     dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered data type for which you want to know whether the scaling is strictly power-of-two.

Description

This function determines whether the registered data type is scaled strictly by a power of two. Fixed-point numbers can be represented as

$$\text{real-world value} = (\text{slope} \times \text{integer}) + \text{bias},$$

where the slope can be expressed as

$$\text{slope} = \text{fractional slope} \times 2^{\text{exponent}}.$$

When $\text{bias} = 0$ and $\text{fractional slope} = 1$, the only scaling factor that remains is a power of two:

$$\text{real-world value} = (2^{\text{exponent}} \times \text{integer}) = (2^{-\text{fraction length}} \times \text{integer}).$$

Trivial scaling is considered a case of power-of-two scaling, with the exponent being equal to zero.

Note Many fixed-point algorithms are designed to accept only power-of-two scaling. For these algorithms, you can call `ssGetDataTypesScalingPow2` in `mdlSetInputPortDataType` and `mdlSetOutputPortDataType`, to prevent unsupported data types from being accepted.

This function errors out when `ssGetDataTypesIsFxpFltApiCompat` returns `FALSE`.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

FixPt_DataTypeIsScalingPow2

See Also

ssGetDataTypeIsScalingTrivial

Version History

Introduced before R2006a

ssGetDataTypesScalingTrivial

Determine whether scaling of registered data type is slope = 1, bias = 0

Syntax

```
extern int ssGetDataTypesScalingTrivial (SimStruct *S, DTypeId
                                         dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered data type for which you want to know whether the scaling is trivial.

Description

This function determines whether the scaling of a registered data type is trivial. In [Slope Bias] representation, fixed-point numbers can be represented as

real-world value = (*slope* × *integer*) + *bias*.

In the trivial case, *slope* = 1 and *bias* = 0.

In terms of binary-point-only scaling, the binary point is to the right of the least significant bit for trivial scaling, meaning that the fraction length is zero:

real-world value = *integer* × 2^{-fraction length} = *integer* × 2⁰.

In either case, trivial scaling means that the real-world value is simply equal to the stored integer value:

real-world value = *integer*.

Scaling is always trivial for pure integers, such as `int8`, and also for the true floating-point types `single` and `double`.

This function errors out when `ssGetDataTypesFxpFltApiCompat` returns `FALSE`.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

FixPt_DataTypeIsScalingTrivial

See Also

ssGetDataTypeIsScalingPow2

Version History

Introduced before R2006a

ssGetDataTypeNumberOfChunks

Return number of chunks in multiword storage container of registered data type

Syntax

```
extern int ssGetDataTypeNumberOfChunks(SimStruct *S,  
                                       DTypeId dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered data type for which you want to know the number of chunks in its multiword storage container.

Description

This function returns the number of chunks in the multiword storage container of the fixed-point data type designated by `dataTypeId`. This function is valid only for a registered data type whose storage container uses a multiword representation. You can use the `ssGetDataTypeStorageContainCat` function to identify the storage container category; for multiword storage containers, the function returns the category `FXP_STORAGE_MULTIWORD`.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

See Also

`ssGetDataTypeStorageContainCat`

Version History

Introduced in R2007b

ssGetDataTypeStorageContainCat

Return storage container category of registered data type

Syntax

```
extern fxpStorageContainerCategory
ssGetDataTypeStorageContainCat(SimStruct *S, DTypeId dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered data type for which you want to know the container category.

Description

This function returns the storage container category of the data type designated by `dataTypeId`. The container category returned by this function is used to store input and output signals, run-time parameters, and DWorks during Simulink simulations.

During simulation, fixed-point signals are held in one of the types of containers shown in the following table. Therefore in many cases, signals are represented in containers with more bits than their actual word length.

Fixed-Point Storage Containers

| Container Category | Signal Word Length | Container Word Length | Container Size |
|---|--|---|---|
| FXP_STORAGE_INT8 (signed) FXP_STORAGE_UINT8 (unsigned) | 1 to 8 bits | 8 bits | 1 byte |
| FXP_STORAGE_INT16 (signed) FXP_STORAGE_UINT16 (unsigned) | 9 to 16 bits | 16 bits | 2 bytes |
| FXP_STORAGE_INT32 (signed) FXP_STORAGE_UINT32 (unsigned) | 17 to 32 bits | 32 bits | 4 bytes |
| FXP_STORAGE_OTHER_SINGLE_WORD | 33 to word length of long data type | Length of long data type | Length of long data type |
| FXP_STORAGE_MULTIWORD | Greater than the word length of long data type to 128 bits | Multiples of length of long data type to 128 bits | Multiples of length of long data type to 128 bits |

When the number of bits in the signal word length is less than the size of the container, the word length bits are always stored in the least significant bits of the container. The remaining container bits must be sign extended to fit the bits of the container:

- If the data type is unsigned, then the sign-extended bits must be cleared to zero.
- If the data type is signed, then the sign-extended bits must be set to one for strictly negative numbers, and cleared to zero otherwise.

The `ssGetDataTypeStorageContainCat` function can also return the following values.

Other Storage Containers

| Container Category | Description |
|---------------------------------------|--|
| <code>FXP_STORAGE_UNKNOWN</code> | Returned if the storage container category is unknown |
| <code>FXP_STORAGE_SINGLE</code> | Container type for a Simulink single |
| <code>FXP_STORAGE_DOUBLE</code> | Container type for a Simulink double |
| <code>FXP_STORAGE_SCALEDDOUBLE</code> | Container type for a data type that has been overridden with Scaled double |

This function errors out when `ssGetDataTypeIsFxpFltApiCompat` returns `FALSE`.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

Because the mapping of storage containers in simulation to storage containers in code generation is not one-to-one, the TLC functions for storage containers in TLC are different from those in simulation. Refer to “Storage Container TLC Functions” on page A-7 for more information:

- `FixPt_DataTypeNativeType`
- `FixPt_DataTypeStorageDouble`
- `FixPt_DataTypeStorageSingle`
- `FixPt_DataTypeStorageScaledDouble`
- `FixPt_DataTypeStorageSInt`
- `FixPt_DataTypeStorageUInt`
- `FixPt_DataTypeStorageSLong`
- `FixPt_DataTypeStorageULong`
- `FixPt_DataTypeStorageSShort`
- `FixPt_DataTypeStorageUShort`

See Also

`ssGetDataTypeStorageContainerSize`

Version History

Introduced before R2006a

ssGetDataTypeStorageContainerSize

Return storage container size of registered data type

Syntax

```
extern size_t ssGetDataTypeStorageContainerSize  
              (SimStruct *S, DTypeId  
              dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered data type for which you want to know the container size.

Description

This function returns the storage container size of the data type designated by `dataTypeId`. This function returns the same value as would the `sizeof()` function; it does not return the word length of either the storage container or the data type. To get the word length of the storage container, use `ssGetDataTypeFxpContainWordLen`. To get the word length of the data type, use `ssGetDataTypeFxpWordLength`.

The container of the size returned by this function stores input and output signals, run-time parameters, and DWorks during Simulink simulations. It is also the appropriate size measurement to pass to functions like `memcpy()`.

This function errors out when `ssGetDataTypeIsFxpFltApiCompat` returns `FALSE`.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

Examples

An `sfix24_En10` data type has a word length of 24, but is actually stored in 32 bits during simulation. For this signal,

- `ssGetDataTypeStorageContainerSize` or `sizeof()` returns 4, which is the storage container size in bytes.

- `ssGetDataTypeFxpContainWordLen` returns 32, which is the storage container word length in bits.
- `ssGetDataTypeFxpWordLength` returns 24, which is the data type word length in bits.

TLC Functions

`FixPt_GetDataTypeStorageContainerSize`

See Also

`ssGetDataTypeFxpContainWordLen`, `ssGetDataTypeFxpWordLength`,
`ssGetDataTypeStorageContainCat`

Version History

Introduced before R2006a

`ssGetDataTypeTotalSlope`

Return total slope of scaling of registered data type

Syntax

```
extern double ssGetDataTypeTotalSlope (SimStruct *S, DTypeId
                                       dataTypeId)
```

Arguments

S

SimStruct representing an S-function block.

dataTypeId

Data type ID of the registered data type for which you want to know the total slope.

Description

Fixed-point numbers can be represented as

real-world value = (slope × integer) + bias,

where the slope can be expressed as

slope = fractional slope × 2^{exponent}.

This function returns the total slope, rather than the fractional slope, of the data type designated by `dataTypeId`. To get the fractional slope, use `ssGetDataTypeFracSlope`:

- If the registered data type has trivial scaling, including `double` and `single` data types, the function returns a total slope of 1.
- If the registered data type is `ScaledDouble`, the function returns the total slope of the nonoverridden data type. Refer to the examples below.

This function errors out when `ssGetDataTypeIsFxpFltApiCompat` returns `FALSE`.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

Examples

The data type `sfix32_En4` becomes `flts32_En4` with data type override. The total slope returned by this function in either case is `0.0625` (2^{-4}).

The data type `ufix16_s7p98` becomes `flt16_s7p98` with data type override. The total slope returned by this function in either case is `7.98`.

TLC Functions

`FixPt_DataTypeTotalSlope`

See Also

`ssGetDataTypeBias`, `ssGetDataTypeFixedExponent`, `ssGetDataTypeFracSlope`

Version History

Introduced before R2006a

ssLogFixptInstrumentation

Record information collected during simulation

Syntax

```
extern void ssLogFixptInstrumentation
    (SimStruct *S,
     double minValue,
     double maxValue,
     int countOverflows,
     int countSaturations,
     int countDivisionsByZero,
     char *pStrName)
```

Arguments

S

SimStruct representing an S-function block.

minValue

Minimum output value that occurred during simulation.

maxValue

Maximum output value that occurred during simulation.

countOverflows

Number of overflows that occurred during simulation.

countSaturations

Number of saturations that occurred during simulation.

countDivisionsByZero

Number of divisions by zero that occurred during simulation.

***pStrName**

The string argument is currently unused.

Description

ssLogFixptInstrumentation records information collected during a simulation, such as output maximum and minimum, any overflows, saturations, and divisions by zero that occurred. The Fixed-Point Tool displays this information after a simulation.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

Version History

Introduced in R2008b

ssRegisterDataTypeFxpBinaryPoint

Register fixed-point data type with binary-point-only scaling and return its data type ID

Syntax

```
extern DTypeId ssRegisterDataTypeFxpBinaryPoint
    (SimStruct *S,
     int isSigned,
     int wordLength,
     int fractionLength,
     int obeyDataTypeOverride)
```

Arguments

S

SimStruct representing an S-function block.

isSigned

TRUE if the data type is signed.

FALSE if the data type is unsigned.

wordLength

Total number of bits in the data type, including any sign bit.

fractionLength

Number of bits in the data type to the right of the binary point.

obeyDataTypeOverride

TRUE indicates that the **Data Type Override** setting for the subsystem is to be obeyed. Depending on the value of **Data Type Override**, the resulting data type could be `Double`, `Single`, `Scaled double`, or the fixed-point data type specified by the other arguments of the function.

FALSE indicates that the **Data Type Override** setting is to be ignored.

Description

This function fully registers a fixed-point data type with the Simulink software and returns a data type ID. Note that unlike the standard Simulink function `ssRegisterDataType`, you do not need to take any additional registration steps. The data type ID can be used to specify the data types of input and output ports, run-time parameters, and DWork states. It can also be used with all the standard data type access methods in `simstruc.h`, such as `ssGetDataTypeSize`.

Use this function if you want to register a fixed-point data type with binary-point-only scaling. Alternatively, you can use one of the other fixed-point registration functions:

- Use `ssRegisterDataTypeFxpFSlopeFixExpBias` to register a data type with [Slope Bias] scaling by specifying the word length, fractional slope, fixed exponent, and bias.
- Use `ssRegisterDataTypeFxpScaledDouble` to register a scaled double.

- Use `ssRegisterDataTypeFxpSlopeBias` to register a data type with [Slope Bias] scaling.

If the registered data type is not one of the Simulink built-in data types, a Fixed-Point Designer software license is checked out. To prevent a Fixed-Point Designer software license from being checked out when you simply open or view a model, protect registration calls with

```
if (ssGetSimMode(S) != SS_SIMMODE_SIZES_CALL_ONLY )  
    ssRegisterDataType...
```

Note Because of the nature of the assignment of data type IDs, you should always use API functions to extract information from a data type ID about a data type in your S-function. For more information, refer to “Data Type IDs” on page A-9.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

None. Data types should be registered in the Simulink software. Registration of data types is not supported in TLC.

See Also

`ssRegisterDataTypeFxpFSlopeFixExpBias`, `ssRegisterDataTypeFxpScaledDouble`,
`ssRegisterDataTypeFxpSlopeBias`, `ssRegisterDataTypeHalfPrecision`

Version History

Introduced before R2006a

ssRegisterDataTypeFxpFSlopeFixExpBias

Register fixed-point data type with [Slope Bias] scaling specified in terms of fractional slope, fixed exponent, and bias, and return its data type ID

Syntax

```
extern DTypeId ssRegisterDataTypeFxpFSlopeFixExpBias
    (SimStruct *S,
     int isSigned,
     int wordLength,
     double fractionalSlope,
     int fixedExponent,
     double bias,
     int obeyDataTypeOverride)
```

Arguments

S

SimStruct representing an S-function block.

isSigned

TRUE if the data type is signed.

FALSE if the data type is unsigned.

wordLength

Total number of bits in the data type, including any sign bit.

fractionalSlope

Fractional slope of the data type.

fixedExponent

Exponent of the slope of the data type.

bias

Bias of the scaling of the data type.

obeyDataTypeOverride

TRUE indicates that the **Data Type Override** setting for the subsystem is to be obeyed.

Depending on the value of **Data Type Override**, the resulting data type could be **Double**, **Single**, **Scaled double**, or the fixed-point data type specified by the other arguments of the function.

FALSE indicates that the **Data Type Override** setting is to be ignored.

Description

This function fully registers a fixed-point data type with the Simulink software and returns a data type ID. Note that unlike the standard Simulink function `ssRegisterDataType`, you do not need to take any additional registration steps. The data type ID can be used to specify the data types of input and output ports, run-time parameters, and DWork states. It can also be used with all the standard data type access methods in `simstruc.h`, such as `ssGetDataTypeSize`.

Use this function if you want to register a fixed-point data type by specifying the word length, fractional slope, fixed exponent, and bias. Alternatively, you can use one of the other fixed-point registration functions:

- Use `ssRegisterDataTypeFxpBinaryPoint` to register a data type with binary-point-only scaling.
- Use `ssRegisterDataTypeFxpScaledDouble` to register a scaled double.
- Use `ssRegisterDataTypeFxpSlopeBias` to register a data type with [Slope Bias] scaling.

If the registered data type is not one of the Simulink built-in data types, a Fixed-Point Designer software license is checked out. To prevent a Fixed-Point Designer software license from being checked out when you simply open or view a model, protect registration calls with

```
if (ssGetSimMode(S) != SS_SIMMODE_SIZES_CALL_ONLY )
    ssRegisterDataType...
```

Note Because of the nature of the assignment of data type IDs, you should always use API functions to extract information from a data type ID about a data type in your S-function. For more information, refer to “Data Type IDs” on page A-9.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

None. Data types should be registered in the Simulink software. Registration of data types is not supported in TLC.

See Also

`ssRegisterDataTypeFxpBinaryPoint`, `ssRegisterDataTypeFxpScaledDouble`,
`ssRegisterDataTypeFxpSlopeBias`, `ssRegisterDataTypeHalfPrecision`

Version History

Introduced before R2006a

ssRegisterDataTypeFxpScaledDouble

Register scaled double data type with [Slope Bias] scaling specified in terms of fractional slope, fixed exponent, and bias, and return its data type ID

Syntax

```
extern DTypeId ssRegisterDataTypeFxpScaledDouble
    (SimStruct *S,
     int isSigned,
     int wordLength,
     double fractionalSlope,
     int fixedExponent,
     double bias,
     int obeyDataTypeOverride)
```

Arguments

S

SimStruct representing an S-function block.

isSigned

TRUE if the data type is signed.

FALSE if the data type is unsigned.

wordLength

Total number of bits in the data type, including any sign bit.

fractionalSlope

Fractional slope of the data type.

fixedExponent

Exponent of the slope of the data type.

bias

Bias of the scaling of the data type.

obeyDataTypeOverride

TRUE indicates that the **Data Type Override** setting for the subsystem is to be obeyed.

Depending on the value of **Data Type Override**, the resulting data type could be **Double**, **Single**, **Scaled double**, or the fixed-point data type specified by the other arguments of the function.

FALSE indicates that the **Data Type Override** setting is to be ignored.

Description

This function fully registers a fixed-point data type with the Simulink software and returns a data type ID. Note that unlike the standard Simulink function `ssRegisterDataType`, you do not need to take any additional registration steps. The data type ID can be used to specify the data types of input and output ports, run-time parameters, and DWork states. It can also be used with all the standard data type access methods in `simstruc.h`, such as `ssGetDataTypeSize`.

Use this function if you want to register a scaled double data type. Alternatively, you can use one of the other fixed-point registration functions:

- Use `ssRegisterDataTypeFxpBinaryPoint` to register a data type with binary-point-only scaling.
- Use `ssRegisterDataTypeFxpFSlopeFixExpBias` to register a data type with [Slope Bias] scaling by specifying the word length, fractional slope, fixed exponent, and bias.
- Use `ssRegisterDataTypeFxpSlopeBias` to register a data type with [Slope Bias] scaling.

If the registered data type is not one of the Simulink built-in data types, a Fixed-Point Designer software license is checked out. To prevent a Fixed-Point Designer software license from being checked out when you simply open or view a model, protect registration calls with

```
if (ssGetSimMode(S) != SS_SIMMODE_SIZES_CALL_ONLY )
    ssRegisterDataType...
```

Note Because of the nature of the assignment of data type IDs, you should always use API functions to extract information from a data type ID about a data type in your S-function. For more information, refer to “Data Type IDs” on page A-9.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

None. Data types should be registered in the Simulink software. Registration of data types is not supported in TLC.

See Also

`ssRegisterDataTypeFxpBinaryPoint`, `ssRegisterDataTypeFxpFSlopeFixExpBias`,
`ssRegisterDataTypeFxpSlopeBias`, `ssRegisterDataTypeHalfPrecision`

Version History

Introduced before R2006a

ssRegisterDataTypeFxpSlopeBias

Register data type with [Slope Bias] scaling and return its data type ID

Syntax

```
extern DTypeId ssRegisterDataTypeFxpSlopeBias
    (SimStruct *S,
     int isSigned,
     int wordLength,
     double totalSlope,
     double bias,
     int obeyDataTypeOverride)
```

Arguments

S

SimStruct representing an S-function block.

isSigned

TRUE if the data type is signed.

FALSE if the data type is unsigned.

wordLength

Total number of bits in the data type, including any sign bit.

totalSlope

Total slope of the scaling of the data type.

bias

Bias of the scaling of the data type.

obeyDataTypeOverride

TRUE indicates that the **Data Type Override** setting for the subsystem is to be obeyed. Depending on the value of **Data Type Override**, the resulting data type could be **Double**, **Single**, **Scaled double**, or the fixed-point data type specified by the other arguments of the function.

FALSE indicates that the **Data Type Override** setting is to be ignored.

Description

This function fully registers a fixed-point data type with the Simulink software and returns a data type ID. Note that unlike the standard Simulink function `ssRegisterDataType`, you do not need to take any additional registration steps. The data type ID can be used to specify the data types of input and output ports, run-time parameters, and DWork states. It can also be used with all the standard data type access methods in `simstruc.h`, such as `ssGetDataTypeSize`.

Use this function if you want to register a fixed-point data type with [Slope Bias] scaling. Alternately, you can use one of the other fixed-point registration functions:

- Use `ssRegisterDataTypeFxpBinaryPoint` to register a data type with binary-point-only scaling.
- Use `ssRegisterDataTypeFxpFSlopeFixExpBias` to register a data type with [Slope Bias] scaling by specifying the word length, fractional slope, fixed exponent, and bias.
- Use `ssRegisterDataTypeFxpScaledDouble` to register a scaled double.

If the registered data type is not one of the Simulink built-in data types, a Fixed-Point Designer software license is checked out. To prevent a Fixed-Point Designer software license from being checked out when you simply open or view a model, protect registration calls with

```
if (ssGetSimMode(S) != SS_SIMMODE_SIZES_CALL_ONLY )
    ssRegisterDataType...
```

Note Because of the nature of the assignment of data type IDs, you should always use API functions to extract information from a data type ID about a data type in your S-function. For more information, refer to “Data Type IDs” on page A-9.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

None.

See Also

`ssRegisterDataTypeFxpBinaryPoint`, `ssRegisterDataTypeFxpFSlopeFixExpBias`,
`ssRegisterDataTypeFxpScaledDouble`, `ssRegisterDataTypeHalfPrecision`

Version History

Introduced before R2006a

ssRegisterDataTypeHalfPrecision

Register half-precision data type and return its data type ID

Syntax

```
extern int ssRegisterDataTypeHalfPrecision
        (SimStruct *S,
         int obeyDataTypeOverride)
```

Arguments

S

SimStruct representing an S-function block.

obeyDataTypeOverride

TRUE indicates that the **Data Type Override** setting is to be obeyed. Depending on the value of **Data Type Override**, the resulting data type could be Double, Single, Scaled double, or the fixed-point data type specified by the other arguments of the function.

FALSE indicates that the **Data Type Override** setting is to be ignored.

Description

This function fully registers the half-precision data type with the Simulink software and returns a data type ID. Note that unlike the standard Simulink function `ssRegisterDataType`, you do not need to take any additional registration steps. The data type ID can be used to specify the data types of input and output ports, run-time parameters, and DWork states. It can also be used with all the standard data type access methods in `simstruc.h`, such as `ssGetDataTypeSize`.

Use this function if you want to register a half-precision data type. For more information on the supported half-precision format, see “Half-Precision Format” on page 35-22.

The registered data type is not one of the Simulink built-in data types, so a Fixed-Point Designer software license is checked out. To prevent a Fixed-Point Designer software license from being checked out when you open or view a model, protect registration calls with

```
if (ssGetSimMode(S) != SS_SIMMODE_SIZES_CALL_ONLY )
    ssRegisterDataType...
```

Note Because of the nature of the assignment of data type IDs, you should always use API functions to extract information from a data type ID about a data type in your S-function. For more information, refer to “Data Type IDs” on page A-9.

Requirement

To use this function, you must include `fixedpoint.h` and `fixedpoint.c`. For more information, see “Structure of the S-Function” on page A-4.

Languages

C

TLC Functions

None.

See Also

ssRegisterDataTypeFxpFSlopeFixExpBias, ssGetDataTypeIsDoubleSingleorHalf, ssGetDataTypeIsFloatingPoint ssGetDataTypeIsHalfPrecision

Version History

Introduced in R2020b

Fixed-Point Designer Examples

Create Fixed-Point Data

This example shows the basics of how to use the fixed-point numeric object `fi`.

The fixed-point numeric object is called `fi` because J.H. Wilkinson used `fi` to denote fixed-point computations in his classic texts *Rounding Errors in Algebraic Processes* (1963), and *The Algebraic Eigenvalue Problem* (1965).

Setup

This example may use display settings or preferences that are different from what you are currently using. To ensure that your current display settings and preferences are not changed by running this example, the example automatically saves and restores them. The following code captures the current states for any display settings or properties that the example changes.

```
originalFormat = get(0, 'format');
format loose
format long g
```

Capture the current state of and reset the `fi` display and logging preferences to the default values.

```
fiprefAtStartOfThisExample = get(fipref);
reset(fipref);
```

Create Fixed-Point Number with Default Properties

To assign a fixed-point data type to a number or variable with the default fixed-point properties, use the `fi` constructor. The resulting fixed-point value is called a `fi` object.

For example, create `fi` objects `a` and `b`. The first input to the `fi` constructor is the value.

```
a = fi(pi)
a =
    3.1416015625
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 13
b = fi(0.1)
b =
    0.0999984741210938
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 18
```

The default fixed-point attributes are displayed. You can specify these attributes when you construct `fi` variables.

The default `WordLength` is 16 bits. When the `FractionLength` property is not specified, it is automatically set to the fraction length that gives the best precision for the given word length while avoiding an overflow, keeping the most-significant bits of the value.

limited to 53 bits. For more information on floating-point arithmetic, refer to Chapter 1 Numerical Computing with MATLAB, by Cleve Moler.

Because most fixed-point processors have data stored in a smaller precision, and then compute with larger precisions, you may want to create a `fi` object that has more precision than double-precision floating point.

For example, initialize a 40-bit unsigned `fi` and multiply using full-precision for products. The full-precision product of 40-bit operands is 80 bits, which is greater precision than standard double-precision floating-point.

```
a = fi(0.1,0,40);
bin(a)

ans =
'110011001100110011001100110011001100110011001101'

b = a*a

b =
    0.01000000000000045

    DataTypeMode: Fixed-point: binary point scaling
    Signedness:   Unsigned
    WordLength:   80
    FractionLength: 86

bin(b)

ans =
'10100011110101110000101000111101011100001111010111000010100011110101110000101001'
```

Access Data

The data can be accessed in a number of ways which map to built-in data types and binary strings.

For example, `double(a)` returns the double-precision floating-point real-world value of `a`, quantized to the precision of `a`.

```
a = fi(pi);
double(a)

ans =
    3.1416015625
```

You can also set the real-world value in a double. For example, set the real-world value of `a` to `e`, quantized to the numeric type of `a`.

```
a.double = exp(1)

a =
    2.71826171875

    DataTypeMode: Fixed-point: binary point scaling
    Signedness:   Signed
    WordLength:   16
    FractionLength: 13
```

Use the `storedInteger` function to return the stored integer in the smallest built-in integer type available, up to 64 bits.

```
storedInteger(a)
```

```
ans = int16
      22268
```

Relationship Between Stored Integer Value and Real-World Value

In binary-point scaling, the relationship between the stored integer value and the real-world value is

$$\text{Real-world value} = (\text{Stored integer}) \cdot 2^{-\text{Fraction length}} .$$

There is also slope-bias scaling, which has the relationship

$$\text{Real-world value} = (\text{Stored integer}) \cdot \text{Slope} + \text{Bias}$$

where

$$\text{Slope} = (\text{Slope adjustment factor}) \cdot 2^{\text{Fixed exponent}} .$$

and

$$\text{Fixed exponent} = - \text{Fraction length} .$$

The math operators of `fi` work with binary-point scaling and real-valued slope-bias scaled `fi` objects.

Other Data Formats

The functions `bin`, `oct`, `dec`, and `hex` return the stored integer in binary, octal, unsigned decimal, and hexadecimal strings, respectively.

```
bin(a)
```

```
ans =
'01010110111111100'
```

```
oct(a)
```

```
ans =
'053374'
```

```
dec(a)
```

```
ans =
'22268'
```

```
hex(a)
```

```
ans =
'56fc'
```

You can use dot notation to set the stored integer from binary, octal, unsigned decimal, and hexadecimal strings.

```
fi(n)
```

```
a.bin = '0110010010001000'  
a =  
    3.1416015625  
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 16  
    FractionLength: 13  
  
    fi( $\phi$ )  
a.oct = '031707'  
a =  
    1.6180419921875  
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 16  
    FractionLength: 13  
  
    fi(e)  
a.dec = '22268'  
a =  
    2.71826171875  
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 16  
    FractionLength: 13  
  
    fi(0.1)  
a.hex = '0333'  
a =  
    0.0999755859375  
    DataTypeMode: Fixed-point: binary point scaling  
    Signedness: Signed  
    WordLength: 16  
    FractionLength: 13
```

Specify Fraction Length

When the `FractionLength` property is not specified, it is computed to give the best precision for the magnitude of the value and given word length, while avoiding overflow. You may also specify the fraction length directly as the fourth numeric argument in the `fi` constructor.

Compare the fraction length of `a`, which was explicitly set to 0, to the fraction length of `b`, which was set to best precision for the magnitude of the value.

```
a = fi(10,1,16,0)  
a =  
    10
```



```

        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 16
        FractionLength: 0

```

```
b = fi(10,1,16)
```

```
b =
    10
```

```

        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 16
        FractionLength: 11

```

The stored integer values of `a` and `b` are different, even though their real-world values are the same. This is because the real-world value of `a` is the stored integer scaled by $2^0 = 1$, while the real-world value of `b` is the stored integer scaled by $2^{-11} = 0.00048828125$.

```
storedInteger(a)
```

```
ans = int16
    10
```

```
storedInteger(b)
```

```
ans = int16
  20480
```

Specify Properties with Name-Value Pair Arguments

You can specify the numeric type properties by passing numeric arguments to the `fi` constructor, as shown above. You can also specify properties by giving the name of the property as a string followed by the value of the property.

```
a = fi(pi, 'WordLength', 20)
```

```
a =
    3.14159393310547
```

```

        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 20
        FractionLength: 17

```

Numeric Type Properties

Each `fi` object has an associated `numericType` object. The `numericType` object stores information about the `fi` object, including word length, fractionlength, and signedness.

```
T = numericType
```

```
T =
```

```

        DataTypeMode: Fixed-point: binary point scaling

```

```
Signedness: Signed
WordLength: 16
FractionLength: 15
```

The numeric type properties can be modified either when the object is created by passing in name-value pair arguments.

```
T = numerictype('WordLength',40,'FractionLength',37)
```

```
T =
```

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 40
FractionLength: 37
```

You can also assign numeric type properties by using the dot notation.

```
T.Signed = false
```

```
T =
```

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Unsigned
WordLength: 40
FractionLength: 37
```

All of the numeric type properties of a `fi` may be set at once by passing in the `numerictype` object. This allows you to, for example, create multiple `fi` objects that share the same numeric type properties.

```
a = fi(pi,'numerictype',T)
```

```
a =
```

```
3.14159265359194
```

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Unsigned
WordLength: 40
FractionLength: 37
```

```
b = fi(exp(1),'numerictype',T)
```

```
b =
```

```
2.71828182845638
```

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Unsigned
WordLength: 40
FractionLength: 37
```

You can also pass the `numerictype` object directly to the `fi` constructor.

```
a1 = fi(pi,T)
```

```
a1 =
```

```
3.14159265359194
```


Cleanup

Set any display settings or preferences that the example changed back to their original states.

```
fipref(fiprefAtStartOfThisExample);  
set(0, 'format', originalFormat);  
%#ok<*NOPTS, *NASGU>
```

See Also

fi | fipref | savefipref | numericity | storedInteger | bin | oct | dec | hex

Perform Fixed-Point Arithmetic

This example shows how to perform basic fixed-point arithmetic operations.

Save the current state of all warnings before beginning.

```
warnstate = warning;
```

Addition and Subtraction

When you add two unsigned fixed-point numbers, you may need a carry bit to correctly represent the result. For this reason, when adding two B-bit numbers with the same scaling, the resulting value has an extra bit compared to the two operands used.

```
a = fi(0.234375,0,4,6);
c = a + a
```

```
c =
    0.4688
```

```
        DataTypeMode: Fixed-point: binary point scaling
        Signedness:   Unsigned
        WordLength:   5
        FractionLength: 6
```

```
a.bin
```

```
ans =
'1111'
```

```
c.bin
```

```
ans =
'11110'
```

With signed, two's-complement numbers, a similar scenario occurs because of the sign extension required to correctly represent the result.

```
a = fi(0.078125,1,4,6);
b = fi(-0.125,1,4,6);
c = a + b
```

```
c =
   -0.0469
```

```
        DataTypeMode: Fixed-point: binary point scaling
        Signedness:   Signed
        WordLength:   5
        FractionLength: 6
```

```
a.bin
```

```
ans =
'0101'
```

```
b.bin
```

```
ans =
'1000'

c.bin

ans =
'11101'
```

If you add or subtract two numbers with different precision, the radix point first needs to be aligned to perform the operation. The result is that there is a difference of more than one bit between the result of the operation and the operands, depending on how far apart the radix points are.

```
a = fi(pi,1,16,13);
b = fi(0.1,1,12,14);
c = a + b

c =
    3.2416

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 18
    FractionLength: 14
```

Further Considerations for Addition and Subtraction

The following pattern is **not** recommended. Because scalar additions are performed at each iteration in the for loop, a bit is added to `temp` during each iteration. As a result, instead of a bit growth of `ceil(log2(Nadds))`, the bit growth is equal to `Nadds`.

```
s = rng; rng('default');
b = fi(4*rand(16,1)-2,1,32,30);
rng(s); % restore RNG state
Nadds = length(b) - 1;
temp = b(1);
for n = 1:Nadds
    temp = temp + b(n+1); % temp has 15 more bits than b
end
```

If the `sum` command is used instead, the bit growth is curbed.

```
c = sum(b) % c has 4 more bits than b

c =
    7.0059

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 36
    FractionLength: 30
```

Multiplication

In general, a full-precision product requires a word length equal to the sum of the word lengths of the operands. In this example, the word length of the product `c` is equal to the word length of `a` plus the word length of `b`. The fraction length of `c` is also equal to the fraction length of `a` plus the fraction length of `b`.

```

a = fi(pi,1,20);
b = fi(exp(1),1,16);
c = a*b

c =
    8.5397

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 36
    FractionLength: 30

```

Assignment

When you assign a fixed-point value into a predefined variable, quantization might be involved. In such cases, the right-hand side of the expression is quantized by rounding to nearest and then saturating, if necessary, before assigning to the left-hand side.

```

N = 10;

a = fi(2*rand(N,1)-1,1,16,15);
b = fi(2*rand(N,1)-1,1,16,15);
c = fi(zeros(N,1),1,16,14);

for n = 1:N
    c(n) = a(n).*b(n);
end

```

When the product $a(n) .* b(n)$ is computed with full precision, an intermediate result with wordlength 32 and fraction length 30 is generated. That result is then quantized to a word length of 16 bits and a fraction length of 14 bits. The quantized value is then assigned to the element $c(n)$.

Quantize Results Explicitly

Often, it is not desirable to round to nearest or to saturate when quantizing a result because of the extra logic and computation required. It may also be undesirable to have to assign to a left-hand side value to perform the quantization. You can use the `quantize` function such purposes. A common case is a feedback-loop. If no quantization is introduced, unbounded bit growth will occur as more input data is provided.

```

a = fi(0.1,1,16,18);
x = fi(2*rand(128,1)-1,1,16,15);
y = fi(zeros(size(x)),1,16,14);

for n = 1:length(x)
    z = y(n);
    y(n) = x(n) - quantize(a.*z,true,16,14,'Floor','Wrap');
end

```

In this example, the product $a .* z$ is computed with full precision and is then quantized to a wordlength of 16 bits and a fraction length of 14 bits. The quantization is done by rounding to floor (truncation) and allowing for wrapping if overflow occurs. Quantization still occurs at assignment because the expression $x(n) - \text{quantize}(a .* z, \dots)$ produces an intermediate result of 18 bits and y is defined to have 16 bits.

To eliminate the quantization at assignment, you can introduce an additional explicit quantization so that no round to nearest or saturation logic is used. Because the left-hand side result has the same 16 bit word length and fraction length of 14 bits as $y(n)$, no quantization is necessary.

```
a = fi(0.1,1,16,18);
x = fi(2*rand(128,1)-1,1,16,15);
y = fi(zeros(size(x)),1,16,14);
T = numerictype(true,16,14);

for n = 1:length(x)
    z = y(n);
    y(n) = quantize(x(n),T,'Floor','Wrap') - ...
           quantize(a.*z,T,'Floor','Wrap');
end
```

Non-Full-Precision Sums

Full-precision sums are not always desirable and can result in complicated and inefficient generated C code. For example, the intermediate result $x(n) - \text{quantize}(\dots)$ in the above example has an 18-bit word length. Instead, it may be desirable to keep all results of addition and subtraction to 16 bits. You can use the `accumpos` and `accumneg` functions to keep the results of addition and subtraction to 16 bits.

```
a = fi(0.1,1,16,18);
x = fi(2*rand(128,1)-1,1,16,15);
y = fi(zeros(size(x)),1,16,14);

T = numerictype(true, 16, 14);

for n = 1:length(x)
    z = y(n);
    y(n) = quantize(x(n),T); % defaults: 'Floor','Wrap'
    y(n) = accumneg(y(n),quantize(a.*z, T)); % defaults: 'Floor','Wrap'
end
```

Model Accumulators

You can use the `accumpos` and `accumneg` functions to model accumulators. The behavior of `accumpos` and `accumneg` corresponds to the `+=` and `-=` operators in C, respectively. A common example is an FIR filter in which the coefficients and input data are represented with 16 bits. The multiplication is performed in full precision, yielding 32 bits, and an accumulator with 8 guard bits. In total, 40 bits are used to enable up to 256 accumulations without the possibility of overflow.

```
b = fi(1/256*[1:128,128:-1:1],1,16); % Filter coefficients
x = fi(2*rand(300,1)-1,1,16,15); % Input data
z = fi(zeros(256,1),1,16,15); % Used to store the states
y = fi(zeros(size(x)),1,40,31); % Initialize Output data

for n = 1:length(x)
    acc = fi(0,1,40,31); % Reset accumulator
    z(1) = x(n); % Load input sample
    for k = 1:length(b)
        acc = accumpos(acc,b(k).*z(k)); % Multiply and accumulate
    end
    z(2:end) = z(1:end-1); % Update states
    y(n) = acc; % Assign output
end
```


Matrix Arithmetic

To simplify syntax and shorten simulation time, you can use matrix arithmetic. For the FIR filter example, you can replace the inner loop with an inner product.

```
z = fi(zeros(256,1),1,16,15); % Used to store the states
y = fi(zeros(size(x)),1,40,31);

for n = 1:length(x)
    z(1) = x(n);
    y(n) = b*z;
    z(2:end) = z(1:end-1);
end
```

The inner product $b*z$ is performed with full precision. Because this is a matrix operation, the bit growth is due to both the multiplication involved and the addition of the resulting products. Therefore, the bit growth depends on the length of the operands. In this example, b and z have length 256, resulting in an 8-bit growth due to the additions. The inner product results in $32 + 8 - 40$ bits, with a fraction length of 31 bits. No quantization occurs in the assignment $y(n) = b*z$ because y was initialized to this format.

If you had to perform an inner product for more than 256 coefficients, the bit growth would be more than 8 bits beyond the 32 needed for the product. If you only had a 40-bit accumulator, you could model the behavior by either introducing a quantizer, as in $y(n) = \text{quantize}(Q, b*z)$, or you could use the `accumpos` function.

Model a Counter

You can use `accumpos` to model a simple counter which wraps after reaching its maximum value. For example, you can model a 3-bit counter as follows.

```
c = fi(0,0,3,0);
Ncounts = 20; % Number of times to count
for n = 1:Ncounts
    c = accumpos(c,1);
end
```

Because the 3-bit counter wraps back to 0 after reaching 7, the final value of the counter is $\text{mod}(20, 8) = 4$.

Arithmetic with Other Built-In Data Types

In the C language, the result of an operation between an integer data type and a double data type promotes to a double. However, in MATLAB®, the result of an operation between a built-in integer data type and a double data type is an integer. In this respect, the `fi` object behaves like the built-in integer data types in MATLAB. The result of an operation between a `fi` and a `double` is a `fi`.

`fi * double`

When doing multiplication between `fi` and `double`, the `double` is cast to a `fi` with the same word length and signedness of the `fi`, and best-precision fraction length. The result of the operation is a `fi`.

```
a = fi(pi);
b = 0.5 * a

b =
    1.5708
```

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 32
FractionLength: 28
```

fi + double or fi - double

When doing addition or subtraction between `fi` and `double`, the `double` is cast to a `fi` with the same `numericType` as the `fi`. The result of the operation is a `fi`.

```
a = fi(pi);
b = a + 1
```

```
b =
    4.1416
```

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 17
FractionLength: 13
```

fi * int8

When doing arithmetic between `fi` and one of the built-in integer data types, `[u]int[8,16,32,64]`, the word length and signedness of the integer are preserved. The result of the operation is a `fi`.

```
a = fi(pi);
b = int8(2) * a
```

```
b =
    6.2832
```

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 24
FractionLength: 13
```

Restore warning states.

```
warning(warnstate);
%#ok< *NASGU, *NOPTS>
```

See Also

`fi` | “Fixed-Point Arithmetic”

Set Fixed-Point Math Attributes

This example shows how to set fixed point math attributes in MATLAB® code.

Use the `fimath` object to control fixed-point math attributes for assignment, addition, subtraction, and multiplication. Use the `setfimath` function to attach a `fimath` object to a `fi` object. Use the `removefimath` function to remove a `fimath` object from a `fi` object.

You can use the MATLAB Coder™ software to generate C code from these examples.

Set and Remove Fixed Point Math Attributes

The `user_written_sum` function shows an example of how to insulate fixed-point operations from global and local `fimath` settings by using the `setfimath` and `removefimath` functions. You can also return from functions with no `fimath` attached to output variables. This gives you local control over fixed-point math settings without interfering with the settings in other functions.

```
function y = user_written_sum(u)
    % Setup
    F = fimath('RoundingMethod','Floor',...
             'OverflowAction','Wrap',...
             'SumMode','KeepLSB',...
             'SumWordLength',32);
    u = setfimath(u,F);
    y = fi(0,true,32,get(u,'FractionLength'),F);
    % Algorithm
    for i=1:length(u)
        y(:) = y + u(i);
    end
    % Cleanup
    y = removefimath(y);
end
```

The `fimath` controls the arithmetic inside the function, but the returned value has no attached `fimath`. This is due to the use of `setfimath` and `removefimath` inside the `user_written_sum` function.

```
u = fi(1:10,true,16,11);
y = user_written_sum(u)
```

```
y =
    55
```

```
        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 32
        FractionLength: 11
```

Generate C Code

If you have a MATLAB Coder license, you can run these commands to generate C code.

```
u = fi(1:10,true,16,11);
codegen user_written_sum -args {u} -config:lib -launchreport
```

The `fimath`, `setfimath`, and `removefimath` functions control the fixed-point math, but the underlying data contained in the variables does not change and so the generated C code does not produce any data copies.

```
int user_written_sum(const short u[10])
{
    int i;
    int y;
    y = 0;
    /* Algorithm */
    for (i = 0; i < 10; i++) {
        y += u[i];
    }
    /* Cleanup */
    return y;
}
```

Mismatched fimath

When you operate on `fi` objects, their `fimath` properties must be equal or you get an error.

```
A = fi(pi, 'ProductMode', 'KeepLSB');
B = fi(2, 'ProductMode', 'SpecifyPrecision');
try
    C = A*B
catch me
    disp(me.message)
end
```

The `embedded.fimath` of both operands must be equal.

To avoid this error, you can remove `fimath` from one of the variables in the expression. In this example, the `fimath` is removed from `B` in the context of the expression without modifying `B` itself. The product is computed using the `fimath` attached to `A`.

```
C = A * removefimath(B)
```

```
C =
    6.2832
```

```
        DataTypeMode: Fixed-point: binary point scaling
          Signedness: Signed
         WordLength: 32
    FractionLength: 26

    RoundingMethod: Nearest
    OverflowAction: Saturate
       ProductMode: KeepLSB
    ProductWordLength: 32
           SumMode: FullPrecision
```

Change fimath on Temporary Variables

If you have variables with no attached `fimath`, but you want to control a particular operation, then you can attach a `fimath` in the context of the expression without modifying the variables.

For example, compute the product using the `fimath` defined by `F`.

```

F = fimath('ProductMode','KeepLSB',...
          'OverflowAction','Wrap',...
          'RoundingMethod','Floor');
A = fi(pi);
B = fi(2);
C = A * setfimath(B,F)

C =
    6.2832

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 32
    FractionLength: 26

    RoundingMethod: Floor
    OverflowAction: Wrap
    ProductMode: KeepLSB
    ProductWordLength: 32
    SumMode: FullPrecision

```

The variable **B** is not changed.

```

B

B =
    2

    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 13

```

Remove fimath Conflict in a Loop

You can compute products and sums to match the accumulator of a DSP with floor rounding and wrap overflow, and use nearest rounding and saturate overflow on the output. To avoid mismatched `fimath` errors, you can remove the `fimath` on the output variable when it is used in a computation with the other variables.

In the `setfimath_removefimath_in_a_loop` function, the products are 32 bits and the accumulator is 40 bits, keeping the least-significant bits with floor rounding and wrap overflow like C's native integer rules. The output uses nearest rounding and saturate overflow.

```

function [y,z] = setfimath_removefimath_in_a_loop(b,a,x,zi)
% Setup
F_floor = fimath('RoundingMethod','Floor',...
                'OverflowAction','Wrap',...
                'ProductMode','KeepLSB',...
                'ProductWordLength',32,...
                'SumMode','KeepLSB',...
                'SumWordLength',40);
F_nearest = fimath('RoundingMethod','Nearest',...
                  'OverflowAction','Wrap');
% Set fimaths that are local to this function
b = setfimath(b,F_floor);
a = setfimath(a,F_floor);
x = setfimath(x,F_floor);

```

```

z = setfimath(zi,F_floor);
% Create y with nearest rounding
y = setfimath(zeros(size(x),'like',zi),F_nearest);
% Algorithm
for j=1:length(x)
    % Nearest assignment into y
    y(j) = b(1)*x(j) + z(1);
    % Remove y's fimath conflict with other fimaths
    z(1) = (b(2)*x(j) + z(2)) - a(2) * removefimath(y(j));
    z(2) = b(3)*x(j) - a(3) * removefimath(y(j));
end
% Cleanup: Remove fimath from outputs
y = removefimath(y);
z = removefimath(z);
end

```

Generate C Code

If you have a MATLAB Coder license, you can run these commands to generate C code using the specified hardware characteristics.

```

N = 256;
t = 1:N;
xstep = [ones(N/2,1);-ones(N/2,1)];
num = [0.0299545822080925 0.0599091644161849 0.0299545822080925];
den = [1 -1.4542435862515900 0.5740619150839550];

b = fi(num,true,16);
a = fi(den,true,16);
x = fi(xstep,true,16,15);
zi = fi(zeros(2,1),true,16,14);

B = coder.Constant(b);
A = coder.Constant(a);

config_obj = coder.config('lib');
config_obj.GenerateReport = true;
config_obj.LaunchReport = true;
config_obj.TargetLang = 'C';
config_obj.DataTypeReplacement = 'CoderTypedefs';
config_obj.GenerateComments = true;
config_obj.GenCodeOnly = true;
config_obj.HardwareImplementation.ProdBitPerChar=8;
config_obj.HardwareImplementation.ProdBitPerShort=16;
config_obj.HardwareImplementation.ProdBitPerInt=32;
config_obj.HardwareImplementation.ProdBitPerLong=40;

codegen -config config_obj setfimath_removefimath_in_a_loop -args {B,A,x,zi}

```

The `fimath`, `setfimath` and `removefimath` functions control the fixed-point math, but the underlying data contained in the variables does not change and so the generated C code does not produce any data copies.

```

void setfimath_removefimath_in_a_loop(const int16_T x[256], const int16_T zi[2],
                                     int16_T y[256], int16_T z[2])
{
    int32_T j;
    int16_T i;
    int16_T il;

```

```

/* Set fimaths that are local to this function */
/* Create y with nearest rounding */
/* Algorithm */
i = zi[0];
i1 = zi[1];
for (j = 0; j < 256; j++) {
    int64_T i3;
    int32_T y_tmp;
    int16_T i2;
    int16_T i4;
    /* Nearest assignment into y */
    i2 = x[j];
    y_tmp = 15705 * i2;
    i3 = y_tmp + ((int64_T)i << 20);
    i4 = (int16_T)((i3 >> 20) + ((i3 & 524288L) != 0L));
    y[j] = i4;
    /* Remove y's fimath conflict with other fimaths */
    i = (int16_T)((31410 * i2 + ((int64_T)i1 << 20)) -
                ((int64_T)(-23826 * i4) << 6) >>
                20);
    i1 = (int16_T)((y_tmp - ((int64_T)(9405 * i4) << 6) >> 20);
}
z[1] = i1;
z[0] = i;
/* Cleanup: Remove fimath from outputs */
}

```

Polymorphic Code

You can use the `setfimath` and `removefimath` functions to write MATLAB code that can be used for both floating-point and fixed-point types.

```

function y = user_written_function(u)
    % Setup
    F = fimath('RoundingMethod','Floor',...
              'OverflowAction','Wrap',...
              'SumMode','KeepLSB',...
              'SumWordLength',32);
    u = setfimath(u,F);
    % Algorithm
    y = u + u;
    % Cleanup
    y = removefimath(y);
end

```

Fixed-Point Inputs

When the function is called with fixed-point inputs, then `fimath F` is used for the arithmetic and the output has no attached `fimath`.

```

u = fi(pi/8,true,16,15,'RoundingMethod','Convergent');
y = user_written_function(u)

y =
    0.7854

```

```

DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed

```

```
WordLength: 32
FractionLength: 15
```

If you have a MATLAB Coder license, you can run these commands to generate C code.

```
u = fi(pi/8,true,16,15,'RoundingMethod','Convergent');
codegen user_written_function -args {u} -config:lib -launchreport
```

The `fimath`, `setfimath` and `removefimath` functions control the fixed-point math, but the underlying data contained in the variables does not change and so the generated C code does not produce any data copies.

```
int user_written_function(short u)
{
    /* Algorithm */
    return u + u;
    /* Cleanup */
}
```

Floating-Point Double Inputs

The `user_written_function` example works with floating-point types because the `setfimath` and `removefimath` functions are pass-through for floating-point types.

```
u = double(pi/8);
codegen user_written_function -args {0} -config:lib -launchreport
```

When compiled with floating-point input, you get the following generated C code.

```
double user_written_function(double u)
{
    /* Algorithm */
    return u + u;
    /* Cleanup */
}
```

More Polymorphic Code

The `user_written_sum_polymorphic` function is written so that the output is created to be the same type as the input. Both floating-point and fixed-point inputs can be used with the `user_written_sum_polymorphic` function.

```
function y = user_written_sum_polymorphic(u)
    % Setup
    F = fimath('RoundingMethod','Floor',...
        'OverflowAction','Wrap',...
        'SumMode','KeepLSB',...
        'SumWordLength',32);
    u = setfimath(u,F);
    if isfi(u)
        y = fi(0,true,32,get(u,'FractionLength'),F);
    else
        y = zeros(1,1,class(u));
    end
    % Algorithm
    for i=1:length(u)
        y(:) = y + u(i);
    end
```



```

    % Cleanup
    y = removefimath(y);
end

```

Generate Fixed-Point C Code

If you have a MATLAB Coder license, you can run these commands to generate C code.

```

u = fi(1:10,true,16,11);
codegen user_written_sum_polymorphic -args {u} -config:lib -launchreport

```

The `fimath`, `setfimath` and `removefimath` functions control the fixed-point math, but the underlying data contained in the variables does not change and so the generated C code does not produce any data copies.

```

int user_written_sum_polymorphic(const short u[10])
{
    int i;
    int y;
    y = 0;
    /* Algorithm */
    for (i = 0; i < 10; i++) {
        y += u[i];
    }
    /* Cleanup */
    return y;
}

```

Generate Floating-Point C Code

If you have a MATLAB Coder license, you can run these commands to generate C code.

```

u = 1:10;
codegen user_written_sum_polymorphic -args {u} -config:lib -launchreport

```

```

double user_written_sum_polymorphic(const double u[10])
{
    double y;
    int i;
    y = 0.0;
    /* Algorithm */
    for (i = 0; i < 10; i++) {
        y += u[i];
    }
    /* Cleanup */
    return y;
}

```

setfimath on Integer Types

Following the established pattern of treating built-in integers like `fi` objects, `setfimath` converts integer input to the equivalent `fi` with attached `fimath`.

```

function y = user_written_u_plus_u(u)
    % Setup
    F = fimath('RoundingMethod','Floor',...
        'OverflowAction','Wrap',...
        'SumMode','KeepLSB',...

```

```
        'SumWordLength',32);  
u = setfimath(u,F);  
% Algorithm  
y = u + u;  
% Cleanup  
y = removefimath(y);  
end
```

If you have a MATLAB Coder license, you can run these commands to generate C code.

```
u = int8(5);  
codegen user_written_u_plus_u -args {u} -config:lib -launchreport
```

The output type was specified by the `fimath` to be 32-bit.

```
int user_written_u_plus_u(signed char u)  
{  
    /* Algorithm */  
    return u + u;  
    /* Cleanup */  
}
```

Disable editor warnings.

```
 %#ok<NASGU>
```

See Also

`fi` | `fimath` | `setfimath` | `removefimath`

Related Examples

- “`fimath` Properties Usage for Fixed-Point Arithmetic” on page 3-10

View Fixed-Point Number Circles

This example shows how to illustrate the definitions of unsigned and signed two's complement integer and fixed-point numbers.

Unsigned Integer Number Circle

Unsigned integers are represented in the binary number system in the following way. Let

$$b = [b(n) \ b(n-1) \ \dots \ b(2) \ b(1)]$$

be the binary digits of an n-bit unsigned integer, where each $b(i)$ is either one or zero. Then the value of b is

$$u = b(n)*2^{(n-1)} + b(n-1)*2^{(n-2)} + \dots + b(2)*2^{(1)} + b(1)*2^{(0)}$$

For example, define a 3-bit unsigned integer quantizer and enumerate its range.

```
q = quantizer('ufixed',[3 0]);
[a,b] = range(q);
u = (a:eps(q):b)'
```

```
u = 8x1
```

```
0
1
2
3
4
5
6
7
```

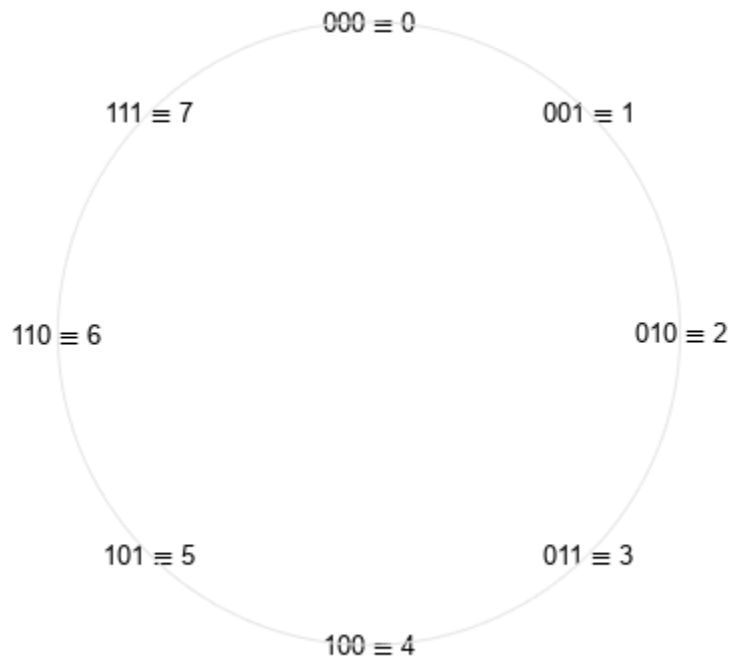
Display those values in binary.

```
b = num2bin(q,u)
```

```
b = 8x3 char array
'000'
'001'
'010'
'011'
'100'
'101'
'110'
'111'
```

Array them around a clock face with their corresponding binary and decimal values.

```
numbercircle(q);
```



Unsigned Fixed-Point Number Circle

Unsigned fixed-point values are unsigned integers that are scaled by a power of two. The negative exponent of the power of two is called the fraction length.

If the unsigned integer u is defined as before and the fraction length is f , then the value of the unsigned fixed-point number is

$$uf = u \cdot 2^{-f}$$

For example, define a 3-bit unsigned fixed-point quantizer with a fraction length of 1 and enumerate its range.

```
q = quantizer('ufixed',[3 1]);
[a,b] = range(q);
uf = (a:eps(q):b)'
```

```
uf = 8×1
    0
    0.5000
    1.0000
    1.5000
    2.0000
    2.5000
    3.0000
    3.5000
```

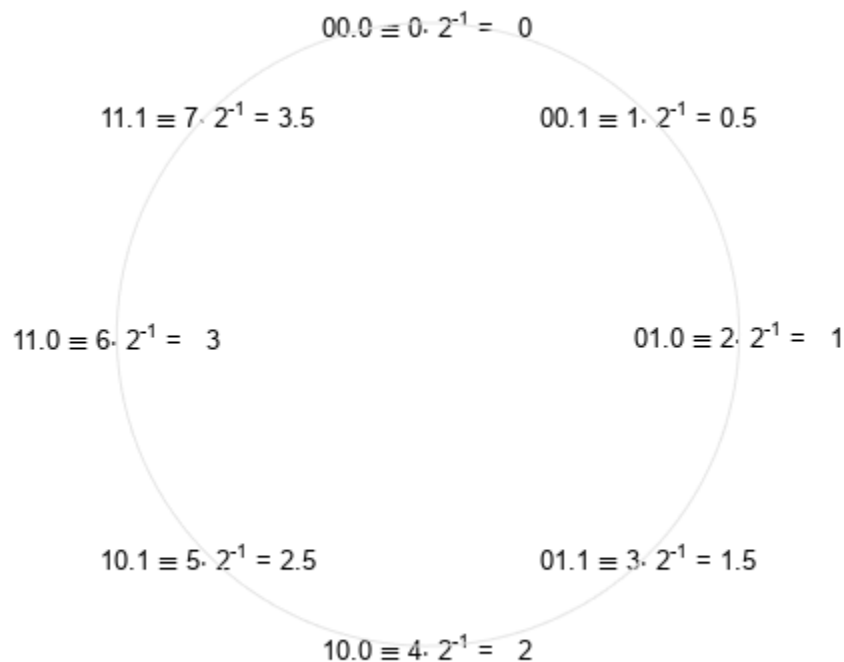
Display those values in binary.

```
b = num2bin(q,uf)
```

```
b = 8x3 char array
'000'
'001'
'010'
'011'
'100'
'101'
'110'
'111'
```

Array them around a clock face with their corresponding binary and decimal values.

```
numbercircle(q);
```



Unsigned Fractional Fixed-Point Number Circle

Unsigned fractional fixed-point numbers are fixed-point numbers whose fraction length f is equal to the wordlength n , which produces a scaling such that the range of numbers is between 0 and $1 - 2^{-f}$, inclusive. This is the most common form of fixed-point numbers because it has the nice property that all of the numbers are less than one and the product of two numbers less than one is a number less than one. Therefore, multiplication does not overflow.

The definition of unsigned fractional fixed-point is the same as unsigned fixed-point, with the restriction that $f=n$, where n is the wordlength in bits.

$$uf = u \cdot 2^{-f}$$

For example, define a 3-bit unsigned fractional fixed-point quantizer, which implies a fraction length of 3.

```
q = quantizer('ufixed',[3 3]);  
[a,b] = range(q);  
uf = (a:eps(q):b)'
```

```
uf = 8x1
```

```
      0  
0.1250  
0.2500  
0.3750  
0.5000  
0.6250  
0.7500  
0.8750
```

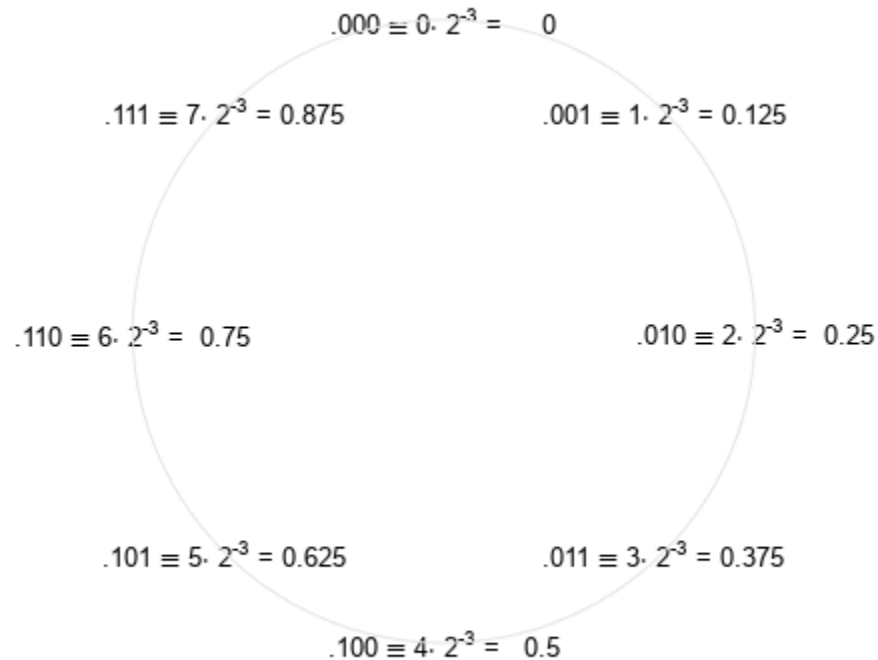
Display those values in binary.

```
b = num2bin(q,uf)
```

```
b = 8x3 char array  
'000'  
'001'  
'010'  
'011'  
'100'  
'101'  
'110'  
'111'
```

Array them around a clock face with their corresponding binary and decimal values.

```
numbercircle(q);
```



Signed Two's-Complement Integer Number Circle

Signed integers are represented in two's-complement in the binary number system in the following way. Let

$$b = [b(n) \ b(n-1) \ \dots \ b(2) \ b(1)]$$

be the binary digits of an n -bit signed integer, where each $b(i)$ is either one or zero. Then the value of b is

$$s = -b(n) \cdot 2^{(n-1)} + b(n-1) \cdot 2^{(n-2)} + \dots + b(2) \cdot 2^{(1)} + b(1) \cdot 2^{(0)}$$

Note that the difference between this and the unsigned number is the negative weight on the most-significant-bit (MSB).

For example, define a 3-bit signed integer quantizer and enumerate its range.

```
q = quantizer('fixed',[3 0]);
[a,b] = range(q);
s = (a:eps(q):b)'
```

```
s = 8×1
```

```
-4
-3
-2
-1
```

```
0  
1  
2  
3
```

Display those values in binary.

```
b = num2bin(q,s)
```

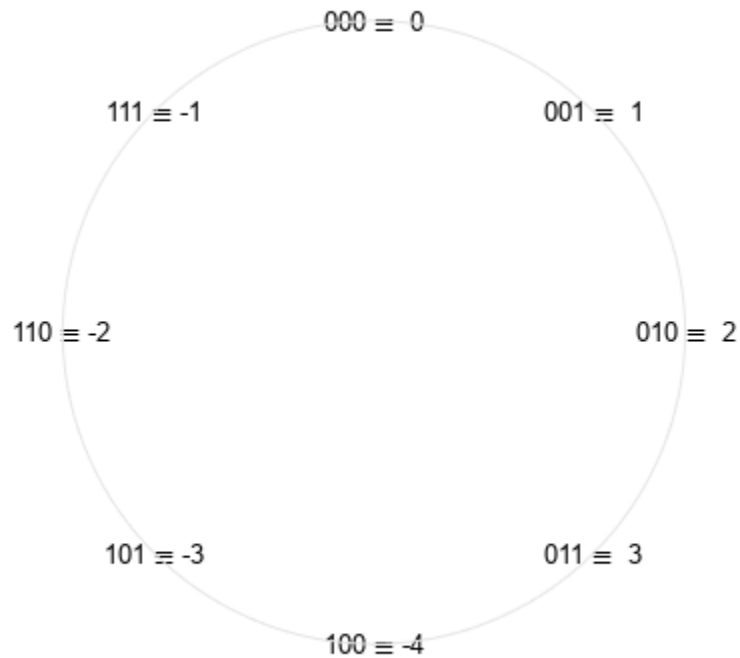
```
b = 8x3 char array
```

```
'100'  
'101'  
'110'  
'111'  
'000'  
'001'  
'010'  
'011'
```

Note that the most-significant-bit of negative numbers is 1 and the most-significant-bit of positive numbers is 0.

Array them around a clock face with their corresponding binary and decimal values.

```
numbercircle(q);
```



The reason for this ungainly looking definition of negative numbers is that addition of all numbers, both positive and negative, is carried out as if they were all positive and then the $n+1$ carry bit is discarded. The result will be correct if there is no overflow.

Signed Fixed-Point Number Circle

Signed fixed-point values are signed integers that are scaled by a power of two. The negative exponent of the power of two is called the fractionlength.

If the signed integer s is defined as before and the fraction length is f , then the value of the signed fixed-point number is

$$sf = s \cdot 2^{-f}$$

For example, define a 3-bit signed fixed-point quantizer with a fraction length of 1 and enumerate its range.

```
q = quantizer('fixed',[3 1]);
[a,b] = range(q);
sf = (a:eps(q):b)'
```

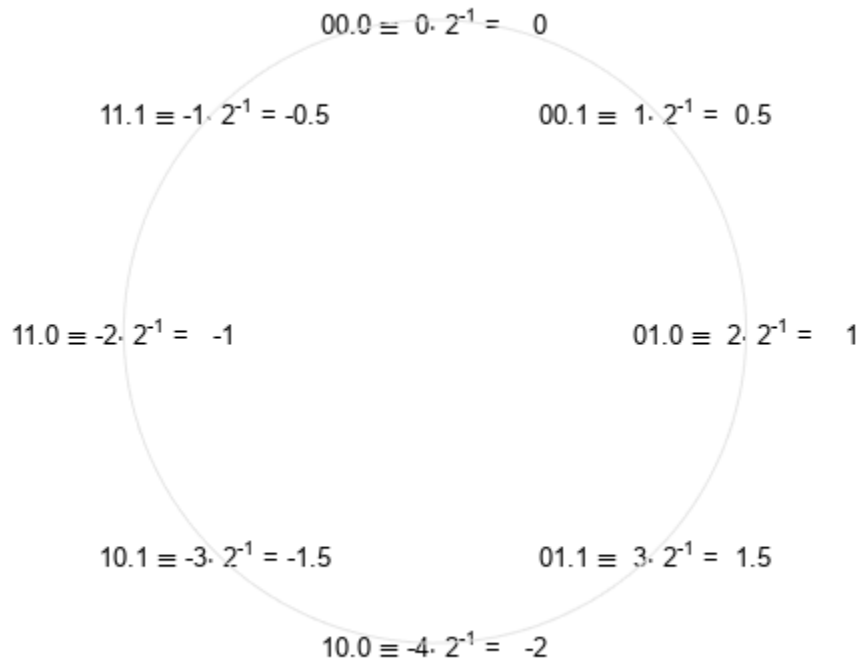
```
sf = 8x1
    -2.0000
    -1.5000
    -1.0000
    -0.5000
         0
     0.5000
     1.0000
     1.5000
```

Display those values in binary.

```
b = num2bin(q,sf)
b = 8x3 char array
    '100'
    '101'
    '110'
    '111'
    '000'
    '001'
    '010'
    '011'
```

Array them around a clock face with their corresponding binary and decimal values.

```
numbercircle(q);
```



Signed Fractional Fixed-Point Number Circle

Signed fractional fixed-point numbers are fixed-point numbers whose fraction length f is one less than the wordlength n , which produces a scaling such that the range of numbers is between -1 and $1 - 2^{-f}$, inclusive. This is the most common form of fixed-point numbers because it has the nice property that the product of two numbers less than one is a number less than one, and so multiplication does not overflow. The only exception is the case when we are multiplying -1 by -1 , because $+1$ is not an element of this number system. Some processors have a special multiplication instruction for this situation, and some add an extra bit in the product to guard against this overflow.

The definition of signed fractional fixed-point is the same as signed fixed-point, with the restriction that $f=n-1$, where n is the word length in bits.

$$sf = s \cdot 2^{-f}$$

For example, define a 3-bit signed fractional fixed-point quantizer, which implies a fraction length of 2.

```
q = quantizer('fixed',[3 2]);
[a,b] = range(q);
sf = (a:eps(q):b)'
```

```
sf = 8×1
    -1.0000
    -0.7500
    -0.5000
```

```

-0.2500
  0
 0.2500
 0.5000
 0.7500

```

Display those values in binary.

```
b = num2bin(q,sf)
```

```
b = 8x3 char array
```

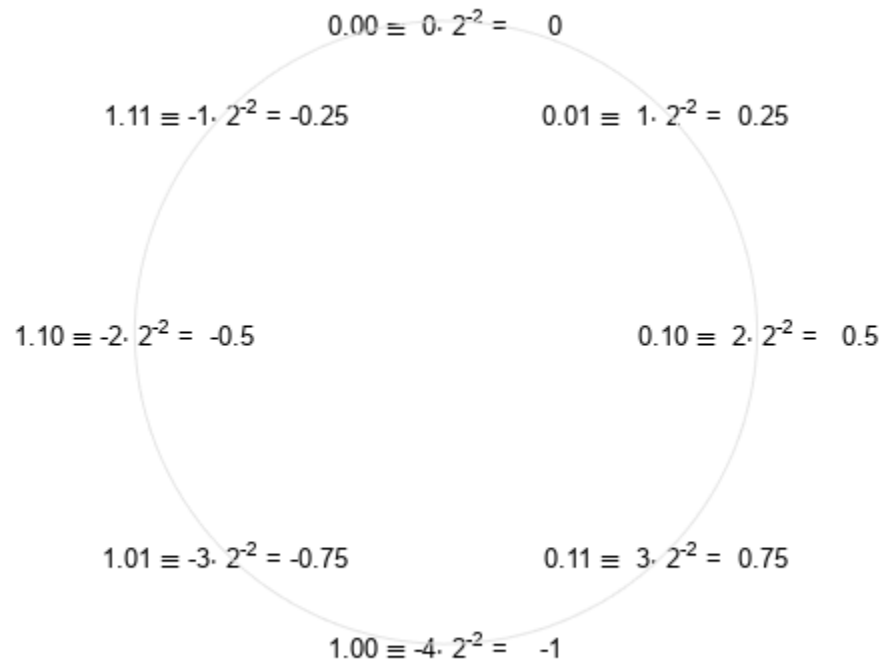
```

'100'
'101'
'110'
'111'
'000'
'001'
'010'
'011'

```

Array them around a clock face with their corresponding binary and decimal values.

```
numbercircle(q);
```



`%#ok<*NOPTS, *NASGU>`

See Also

“Data Types and Scaling in Digital Hardware” on page 35-2 | “Arithmetic Operations” on page 1-9

Perform Binary-Point Scaling

This example shows how to perform binary point scaling using a `fi` object.

Construct `fi` Object

Use the `fi` constructor, `a = fi(v,s,w,f)`, to return a `fi` object with value `v`, signedness `s`, word length `w`, and fraction length `f`. If `s` is true (signed), the leading or most significant bit (MSB) in the resulting `fi` object is always the sign bit. The fraction length `f` gives the scaling, $2^{(-f)}$. The fraction length or the scaling determines the position of the binary point in the `fi` object.

For example, create a signed 8-bit `fi` object with a value of 0.5 and a scaling of $2^{(-7)}$.

```
a = fi(0.5,true,8,7)
```

```
a =
    0.5000
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 8
    FractionLength: 7
```

Fraction Length Positive and Less than Word Length

When the fraction length `f` is positive and less than the word length, the binary point lies `f` places to the left of the least significant bit (LSB) and within the word.

For example, in a signed 3-bit `fi` with fraction length of 1 and value -0.5, the binary point lies 1 place to the left of the LSB. In this case, each bit is set to 1 and the binary equivalent of the `fi` with its binary point is 11.1.

The real world value of -0.5 is obtained by multiplying each bit by its scaling factor, starting with the LSB and working up to the signed MSB.

$$(1*2^{-1}) + (1*2^0) + (-1*2^1) = -0.5$$

`storedInteger(a)` returns the stored signed, unscaled integer value -1.

$$(1*2^0) + (1*2^1) + (-1*2^2) = -1$$

```
a = fi(-0.5,true,3,1)
```

```
a =
   -0.5000
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 3
    FractionLength: 1
```

```
bin(a)
```

```
ans =
'111'
```

```
storedInteger(a)
```

```
ans = int8
     -1
```

Fraction Length Positive and Greater than Word Length

When the fraction length f is positive and greater than the word length, the binary point lies f places to the left of the LSB and outside the word.

For example the binary equivalent of a signed 3-bit word with fraction length of 4 and value of -0.0625 is `._111`. Here, `_` in the `._111` denotes an unused bit that is not a part of the 3-bit word. The first 1 after the `_` is the MSB or the sign bit.

The real world value of -0.0625 is computed as follows (LSB to MSB).

$$(1*2^{-4}) + (1*2^{-3}) + (-1*2^{-2}) = -0.0625$$

```
b = fi(-0.0625,true,3,4)
```

```
b =
     -0.0625
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 3
      FractionLength: 4
```

```
bin(b)
```

```
ans =
'111'
```

```
storedInteger(b)
```

```
ans = int8
     -1
```

Fraction Length is Negative Integer and Less than Word Length

When the fraction length f is negative, the binary point lies f places to the right of LSB and is outside the physical word.

For instance, in `c = fi(-4,true,3,-2)` the binary point lies 2 places to the right of the LSB `111__`. Here, the two right most spaces are unused bits that are not part of the 3-bit word. The right most 1 is the LSB and the leading 1 is the sign bit.

The real world value of -4 is obtained by multiplying each bit by its scaling factor $2^{(-f)}$, for instance $2^{(-(-2))} = 2^{(2)}$ for the LSB, and then adding the products together.

$$(1*2^2) + (1*2^3) + (-1*2^4) = -4$$

```
c = fi(-4,true,3,-2)
```

```
c =
     -4
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 3
      FractionLength: -2
```

```

bin(c)

ans =
'111'

storedInteger(c)

ans = int8
    -1

```

Fraction Length Set Automatically to the Best Precision Possible and is Negative

Create a signed 3-bit `fi` where the fraction length is set automatically depending on the value that the `fi` is supposed to contain. The resulting `fi` has a value of 6, with a wordlength of 3 bits and a fraction length of -1. Here the binary point is 1 place to the right of the LSB: 011_.. The _ is again an unused bit and the first 1 before the _ is the LSB. The leading 1 is the sign bit.

The real world value of 6 is obtained as follows:

$$(1*2^1) + (1*2^2) + (-0*2^3) = 6$$

```
d = fi(5,true,3)
```

```
d =
    6
```

```

        DataTypeMode: Fixed-point: binary point scaling
        Signedness: Signed
        WordLength: 3
        FractionLength: -1

```

```

bin(d)

ans =
'011'

storedInteger(d)

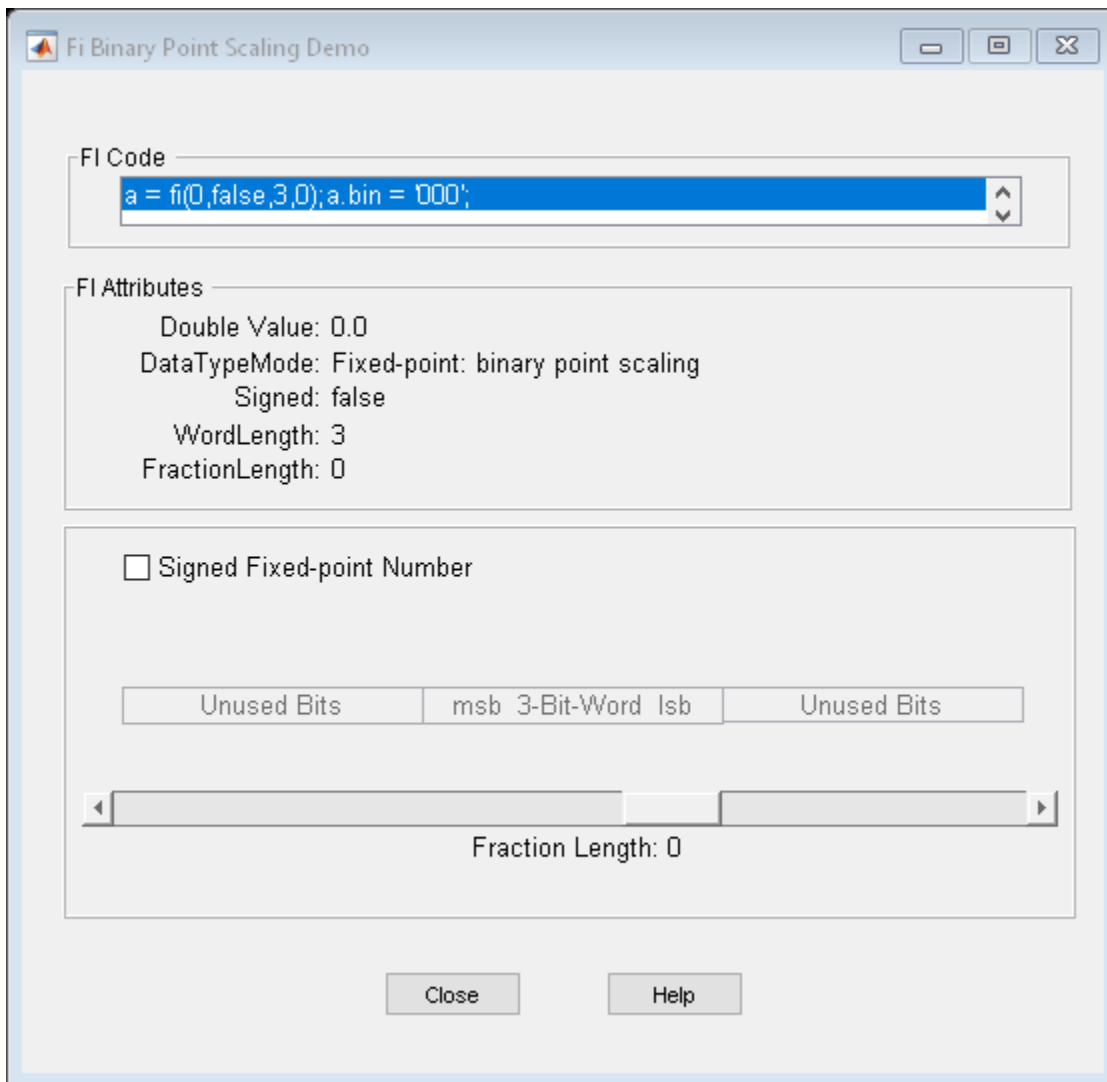
ans = int8
    3

```

Interactive `fi` Binary Point Scaling Example

To run an interactive binary-point scaling example, enter at the MATLAB® Command Window:

```
fibinscaling
```



This interactive example allows you to change the fraction length of a 3-bit fixed-point number by moving the binary point using a slider. The fraction length can be varied from -3 to 5. You can change the value of the 3 bits to '0' or '1' for either signed or unsigned numbers.

`%#ok<*NOPTS,*NASGU>`

See Also

`fi` | `bin` | `storedInteger`

Compute Quantization Error

This example shows how to compute and compare the statistics of the signal quantization error when using various rounding methods. Quantization occurs when a data type cannot represent a value exactly. In these cases, the value must be rounded to the nearest value that can be represented by the data type.

First, a random signal is created that spans the range of the quantizer object. Next, the signal is quantized, respectively, with rounding methods 'fix', 'floor', 'ceil', 'nearest', and 'convergent', and the statistics of the signal are estimated.

The theoretical probability density function of the quantization error is computed with the `errpdf` function, the theoretical mean of the quantization error is computed with the `errmean` function, and the theoretical variance of the quantization error is computed with the `errvar` function.

Create Uniformly Distributed Random Signal

Create a uniformly distributed random signal that spans the domain -1 to 1 of the fixed-point quantizer object `q`.

```
q = quantizer([8 7]);
r = realmax(q);
u = r*(2*rand(50000,1) - 1);
xi = linspace(-2*eps(q),2*eps(q),256);
```

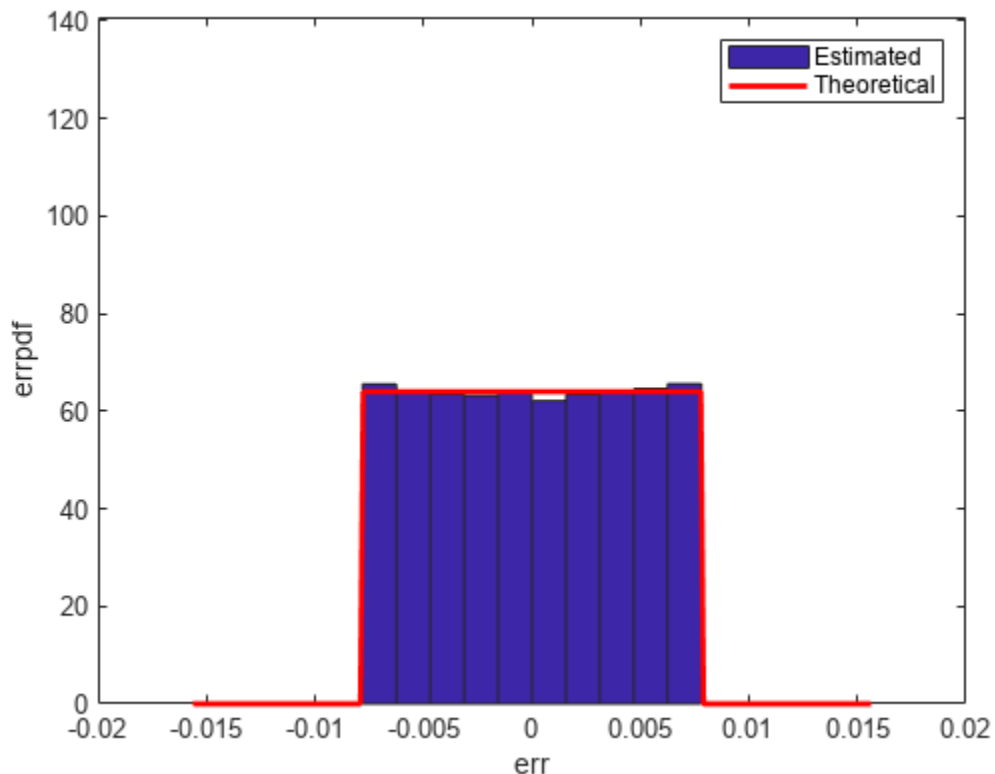
Fix: Round Towards Zero

With 'fix' rounding, the probability density function is twice as wide as the others. For this reason, the variance is four times that of the others.

```
q = quantizer('fix',[8 7]);
err = quantize(q,u) - u;
f_t = errpdf(q,xi);
mu_t = errmean(q);
v_t = errvar(q);

qerrordemoplot(q,f_t,xi,mu_t,v_t,err)

Estimated error variance (dB) = -46.8586
Theoretical error variance (dB) = -46.9154
Estimated mean = 7.788e-06
Theoretical mean = 0
```



The theoretical variance is $\text{eps}(q)^2/3$ and the theoretical mean is 0.

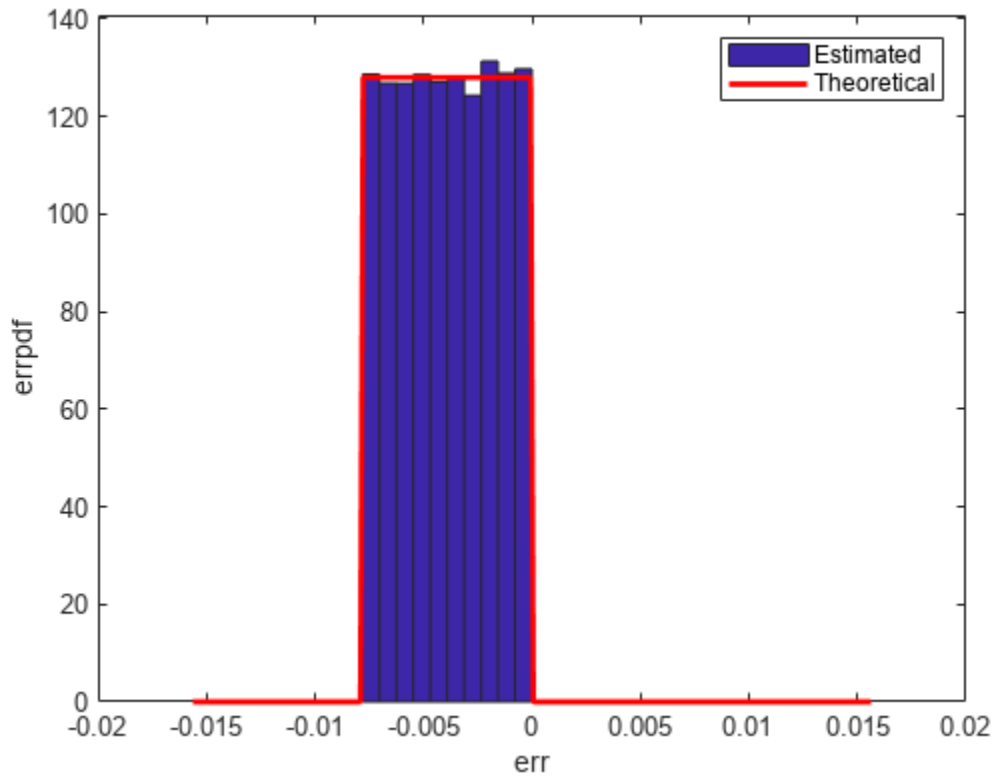
Floor: Round Towards Negative Infinity

'floor' rounding is often called truncation when used with integers and fixed-point numbers that are represented using two's complement notation. It is the most common rounding mode of DSP processors because it requires no hardware to implement. 'floor' does not produce quantized values that are as close to the true values as 'round' will, but it has the same variance. Using 'floor', small signals that vary in sign will be detected, whereas in 'round' they will be lost.

```
q = quantizer('floor',[8 7]);
err = quantize(q,u) - u;
f_t = errpdf(q,xi);
mu_t = errmean(q);
v_t = errvar(q);

qerrordemoplot(q,f_t,xi,mu_t,v_t,err)

Estimated error variance (dB) = -52.9148
Theoretical error variance (dB) = -52.936
Estimated mean = -0.0038956
Theoretical mean = -0.0039062
```



The theoretical variance is $\text{eps}(q)^2/12$ and the theoretical mean is $-\text{eps}(q)/2$.

Ceil: Round Towards Positive Infinity

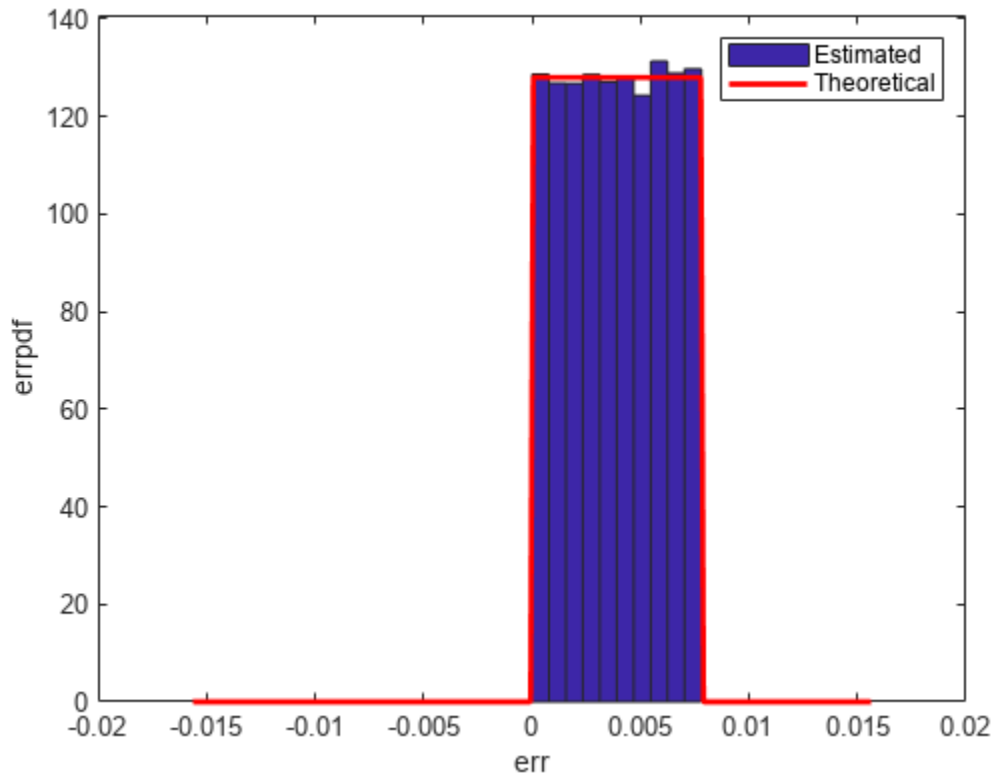
```

q = quantizer('ceil',[8 7]);
err = quantize(q,u) - u;
f_t = errpdf(q,xi);
mu_t = errmean(q);
v_t = errvar(q);

qerrordemoplot(q,f_t,xi,mu_t,v_t,err)

Estimated error variance (dB) = -52.9148
Theoretical error variance (dB) = -52.936
Estimated mean = 0.0039169
Theoretical mean = 0.0039062

```



The theoretical variance is $\text{eps}(q)^2/12$ and the theoretical mean is $\text{eps}(q)/2$.

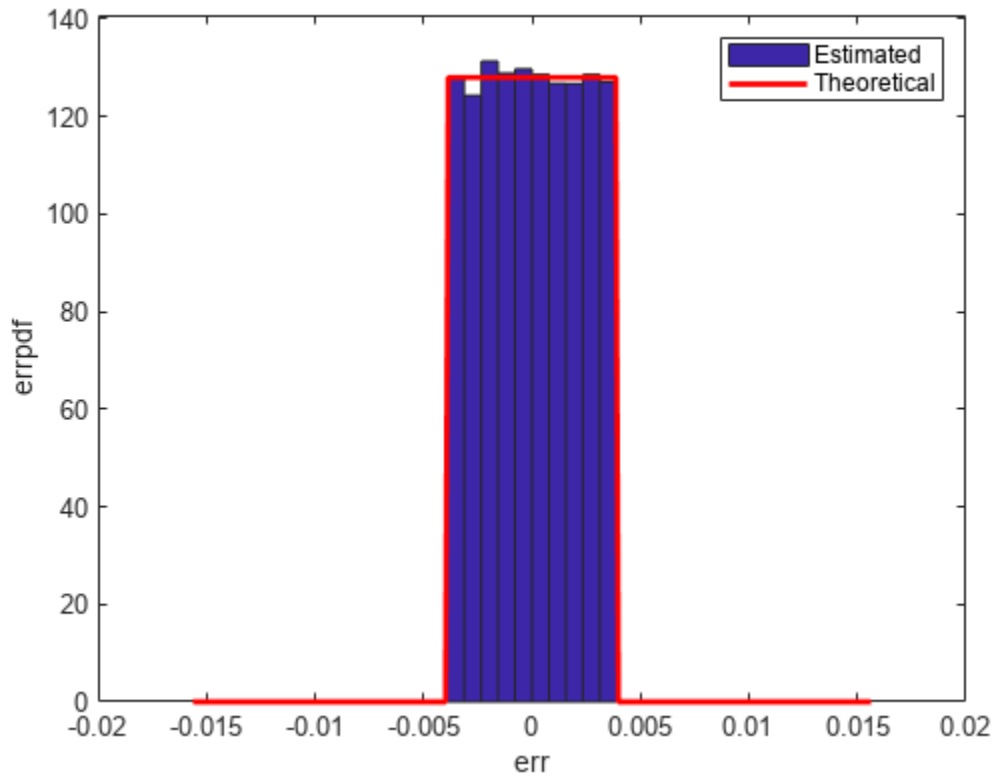
Round: Round to Nearest; In a Tie Round to Largest Magnitude

'round' is more accurate than 'floor', but all values smaller than $\text{eps}(q)$ get rounded to zero and are lost.

```
q = quantizer('nearest',[8 7]);
err = quantize(q,u) - u;
f_t = errpdf(q,xi);
mu_t = errmean(q);
v_t = errvar(q);

qerrordemoplot(q,f_t,xi,mu_t,v_t,err)

Estimated error variance (dB) = -52.9579
Theoretical error variance (dB) = -52.936
Estimated mean = -2.212e-06
Theoretical mean = 0
```



The theoretical variance is $\text{eps}(q)^2/12$ and the theoretical mean is 0.

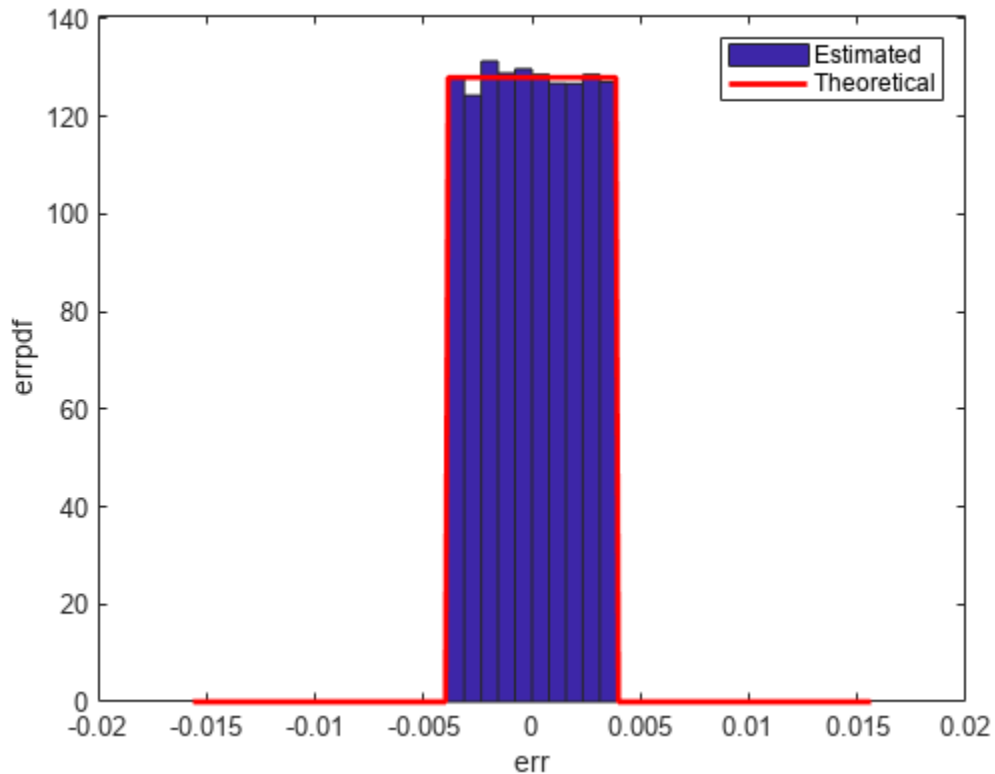
Convergent: Round to Nearest; In a Tie Round to Even

'convergent' rounding eliminates the bias introduced by ordinary 'round' caused by always rounding the tie in the same direction.

```
q = quantizer('convergent',[8 7]);
err = quantize(q,u) - u;
f_t = errpdf(q,xi);
mu_t = errmean(q);
v_t = errvar(q);

qerrordemoplot(q,f_t,xi,mu_t,v_t,err)

Estimated error variance (dB) = -52.9579
Theoretical error variance (dB) = -52.936
Estimated mean = -2.212e-06
Theoretical mean = 0
```



The theoretical variance is $\text{eps}(q)^2/12$ and the theoretical mean is 0.

Compare Nearest and Convergent Rounding

The error probability density function for convergent rounding is difficult to distinguish from that of round-to-nearest by looking at the plot.

The error probability density function of convergent is

$f(\text{err}) = 1/\text{eps}(q)$, for $-\text{eps}(q)/2 \leq \text{err} \leq \text{eps}(q)/2$, and 0 otherwise

while the error probability density function of round is

$f(\text{err}) = 1/\text{eps}(q)$, for $-\text{eps}(q)/2 < \text{err} \leq \text{eps}(q)/2$, and 0 otherwise

The error probability density function of convergent is symmetric, while round is slightly biased towards the positive.

The only difference is the direction of rounding in a tie.

```
x = (-3.5:3.5)';
[x convergent(x) nearest(x)]
```

```
ans = 8x3
```

```
-3.5000    -4.0000    -3.0000
-2.5000    -2.0000    -2.0000
-1.5000    -2.0000    -1.0000
```

| | | |
|---------|--------|--------|
| -0.5000 | 0 | 0 |
| 0.5000 | 0 | 1.0000 |
| 1.5000 | 2.0000 | 2.0000 |
| 2.5000 | 2.0000 | 3.0000 |
| 3.5000 | 4.0000 | 4.0000 |

See Also

quantizer | quantize | "Rounding" on page 36-2

Detect Limit Cycles in Fixed-Point State-Space Systems

This example shows how to analyze a fixed-point state-space system to detect limit cycles.

The example focuses on detecting large scale limit cycles due to overflow with zero inputs and highlights the conditions that are sufficient to prevent such oscillations.

Select a State-Space Representation of the System

See that the system is stable by observing that the eigenvalues of the state-transition matrix A have magnitudes less than 1.

```
A = [0 1; -.5 1]; B = [0; 1]; C = [1 0]; D = 0;
eig(A)
```

```
ans = 2×1 complex

    0.5000 + 0.5000i
    0.5000 - 0.5000i
```

Filter Implementation

The function `fisisostatespacefilter` implements a single-input single-output state-space filter.

```
type('fisisostatespacefilter.m')
```

```
function [y,z] = fisisostatespacefilter(A,B,C,D,x,z)
%FISISOSTATESPACEFILTER Single-input, single-output statespace filter
% [Y,Zf] = FISISOSTATESPACEFILTER(A,B,C,D,X,Zi) filters data X with
% initial conditions Zi with the state-space filter defined by matrices
% A, B, C, D. Output Y and final conditions Zf are returned.
```

```
% Copyright 2004-2011 The MathWorks, Inc.
```

```
y = x;
z(:,2:length(x)+1) = 0;
for k=1:length(x)
    y(k) = C*z(:,k) + D*x(k);
    z(:,k+1) = A*z(:,k) + B*x(k);
end
```

Floating-Point Filter

Create a floating-point filter and observe the trajectory of the states.

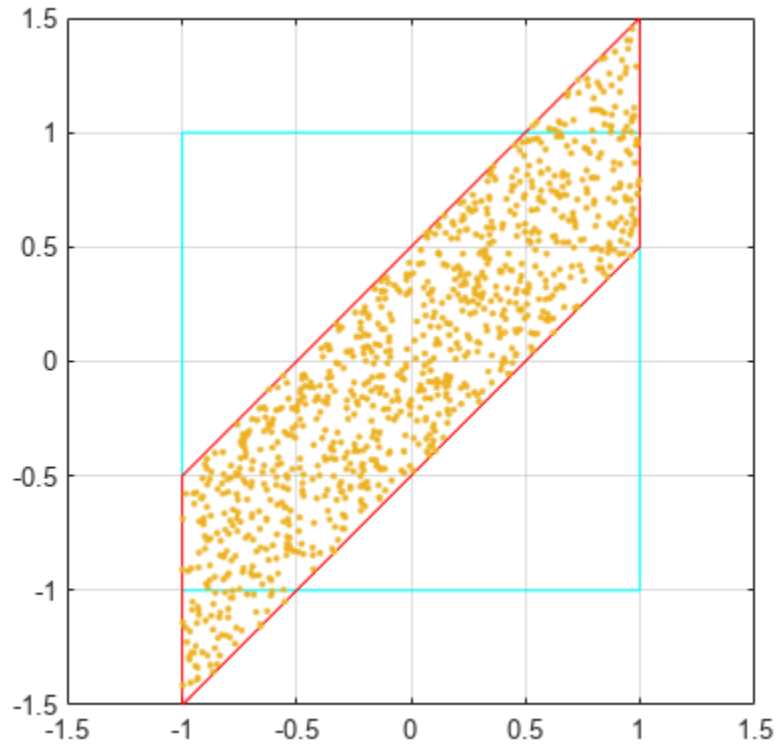
First, choose random states within the unit square and observe where they are projected after one step of being multiplied by the state-transition matrix A.

```
rng('default');
clf
x1 = [-1 1 1 -1 -1];
y1 = [-1 -1 1 1 -1];
plot(x1,y1,'c')
axis([-1.5 1.5 -1.5 1.5]); axis square; grid;
hold on
```


Plot the projection of the square.

```
p = A*[x1;y1];
plot(p(1,:),p(2,:), 'r')

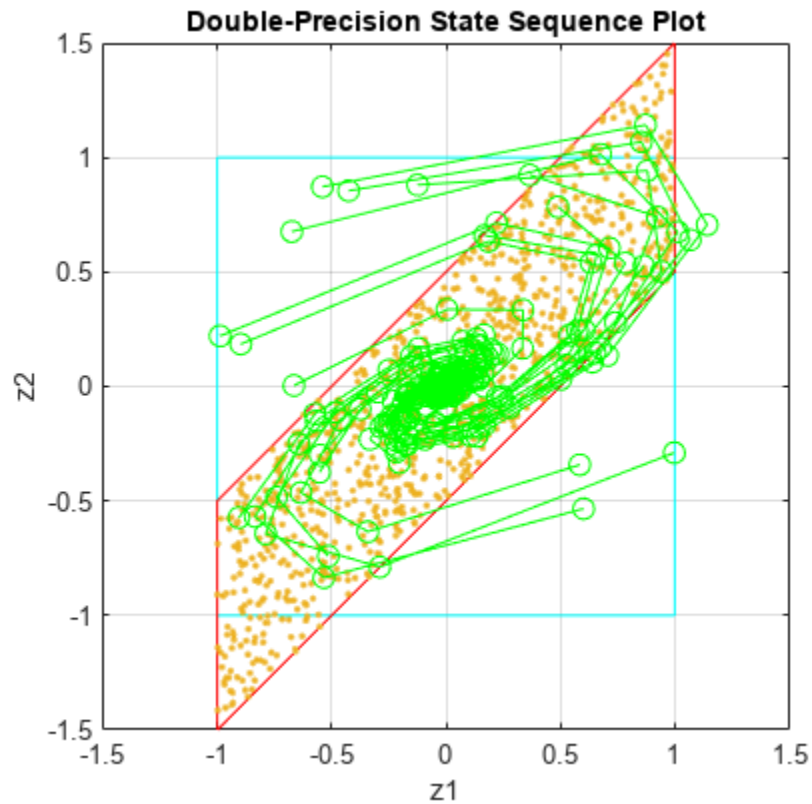
r = 2*rand(2,1000)-1;
pr = A*r;
plot(pr(1,:),pr(2,:), '.')
```



Follow Random Initial States Through Time

Drive the filter with a random initial state, normalized to be inside the unit square, with the input all zero, and run the filter.

```
x = zeros(10,1);
zi = [0;0];
q = quantizer([16 15]);
for k=1:20
    y = x;
    zi(:) = randquant(q,size(A,1),1);
    [y,zf] = fisisostatespacefilter(A,B,C,D,x,zi);
    plot(zf(1,:), zf(2,:), 'go-', 'markersize',8);
end
title('Double-Precision State Sequence Plot');
xlabel('z1'); ylabel('z2')
```



Some of the states wander outside the unit square and eventually wind down to the zero state at the origin, $z = [0; 0]$.

State Trajectory

Because the eigenvalues are less than one in magnitude, the system is stable and all initial states wind down to the origin with zero input. However, the eigenvalues don't tell the whole story about the trajectory of the states, as in this example, where the states were projected outward first before they start to contract.

The singular values of A give us a better indication of the overall state trajectory.

```
svd(A)
```

```
ans = 2x1
```

```
1.4604
0.3424
```

The largest singular value is about 1.46, which indicates that states aligned with the corresponding singular vector will be projected away from the origin.

Create Fixed-Point Filter

Create a fixed-point filter and check for limit cycles.

The MATLAB® code for the filter remains the same. It becomes a fixed-point filter when it is driven with fixed-point inputs. For the sake of illustrating overflow oscillation, choose product and sum data types that will overflow.

```
rng('default');
F = fimath('OverflowAction','Wrap',...
          'ProductMode','SpecifyPrecision',...
          'ProductWordLength',16,'ProductFractionLength',15,...
          'SumMode','SpecifyPrecision',...
          'SumWordLength',16,'SumFractionLength',15);
```

```
A = fi(A,'fimath',F)
```

```
A =
```

```
    0    1.0000
-0.5000 1.0000
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 14
```

```
    RoundingMethod: Nearest
    OverflowAction: Wrap
    ProductMode: SpecifyPrecision
    ProductWordLength: 16
    ProductFractionLength: 15
    SumMode: SpecifyPrecision
    SumWordLength: 16
    SumFractionLength: 15
    CastBeforeSum: true
```

```
B = fi(B,'fimath',F)
```

```
B =
```

```
    0
    1
```

```
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Signed
    WordLength: 16
    FractionLength: 14
```

```
    RoundingMethod: Nearest
    OverflowAction: Wrap
    ProductMode: SpecifyPrecision
    ProductWordLength: 16
    ProductFractionLength: 15
    SumMode: SpecifyPrecision
    SumWordLength: 16
    SumFractionLength: 15
    CastBeforeSum: true
```

```
C = fi(C,'fimath',F)
```

```
C =
```

```
    1    0
```

```
    DataTypeMode: Fixed-point: binary point scaling
```

```
Signedness: Signed
WordLength: 16
FractionLength: 14

RoundingMethod: Nearest
OverflowAction: Wrap
ProductMode: SpecifyPrecision
ProductWordLength: 16
ProductFractionLength: 15
SumMode: SpecifyPrecision
SumWordLength: 16
SumFractionLength: 15
CastBeforeSum: true

D = fi(D, 'fimath', F)

D =
    0

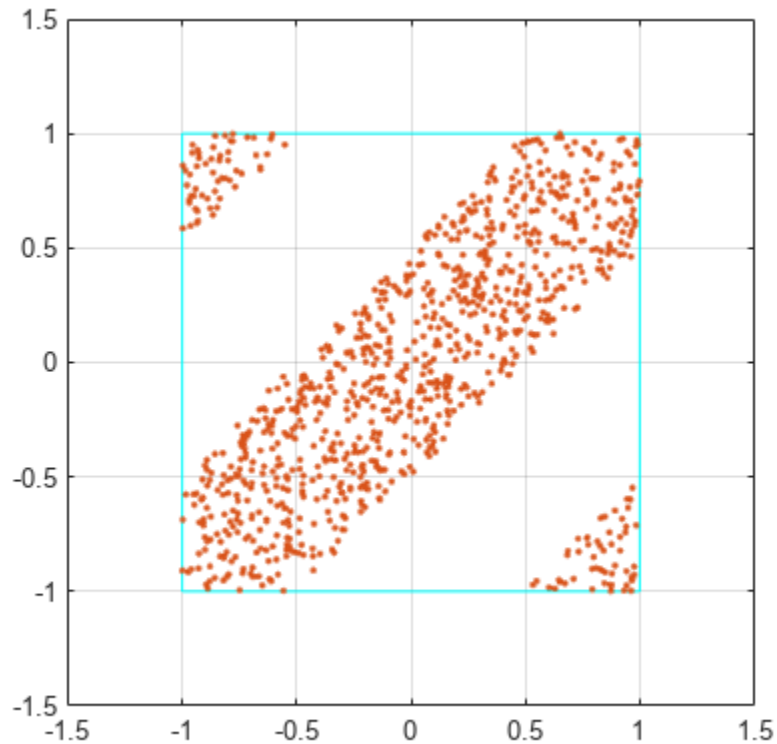
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 16
FractionLength: 15

RoundingMethod: Nearest
OverflowAction: Wrap
ProductMode: SpecifyPrecision
ProductWordLength: 16
ProductFractionLength: 15
SumMode: SpecifyPrecision
SumWordLength: 16
SumFractionLength: 15
CastBeforeSum: true
```

Plot Projection of Square in Fixed-Point

Choose random states within the unit square and observe where they are projected after one step of being multiplied by the state-transition matrix A. This time the matrix A is fixed-point.

```
clf
r = 2*rand(2,1000)-1;
pr = A*r;
plot([-1 1 1 -1 -1],[-1 -1 1 1 -1], 'c')
axis([-1.5 1.5 -1.5 1.5]); axis square; grid;
hold on
plot(pr(1,:),pr(2:,:), '.')
```

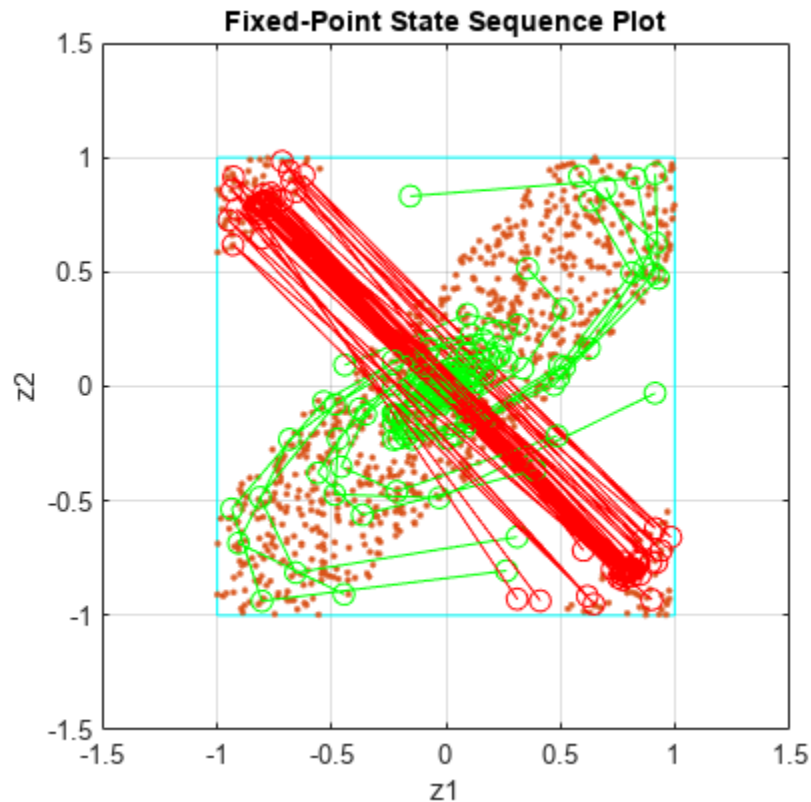


The triangles that projected out of the square in floating-point are now wrapped back into the interior of the square.

Execute Fixed-Point Filter

Drive the filter with fixed-point data types.

```
x = fi(zeros(10,1),1,16,15,'fimath',F);
zi = fi([0;0],1,16,15,'fimath',F);
q = assignmentquantizer(zi);
e = double(eps(zi));
rng('default');
for k=1:20
    y = x;
    zi(:) = randquant(q,size(A,1),1);
    [y,zf] = fisisostatespacefilter(A,B,C,D,x,zi);
    if abs(double(zf(end)))>0.5, c='ro-'; else, c='go-'; end
    plot(zf(1,:), zf(2,:),c,'markersize',8);
end
title('Fixed-Point State Sequence Plot');
xlabel('z1'); ylabel('z2')
```



Trying this for other randomly chosen initial states illustrates that once a state enters one of the triangular regions, then it is projected into the other triangular region, and back and forth, and never escapes.

Sufficient Conditions for Preventing Overflow Limit Cycles

There are two sufficient conditions to prevent overflow limit cycles in a system:

- The system is stable: $\text{abs}(\text{eig}(A)) < 1$
- The matrix A is normal: $A' * A = A * A'$

For the current representation, the second condition does not hold.

Apply Similarity Transform to Create Normal A

Apply a similarity transformation to the original system to create a normal state-transition matrix A_2 .

```
T = [-2 0; -1 1];
Tinv = [-.5 0; -.5 1];
A2 = Tinv*A*T; B2 = Tinv*B; C2 = C*T; D2 = D;
```

Similarity transformations preserve eigenvalues. The system transfer function of the transformed system remains same as before. However, the transformed state transformation matrix A_2 is normal.

Check for Limit Cycles on Transformed System

Plot the projection of the square of the normal-form system.

```

clf
r = 2*rand(2,1000)-1;
pr = A2*r;
plot([-1 1 1 -1 -1],[-1 -1 1 1 -1],'c')
axis([-1.5 1.5 -1.5 1.5]); axis square; grid;
hold on
plot(pr(1,:),pr(2:,:),'.')

```

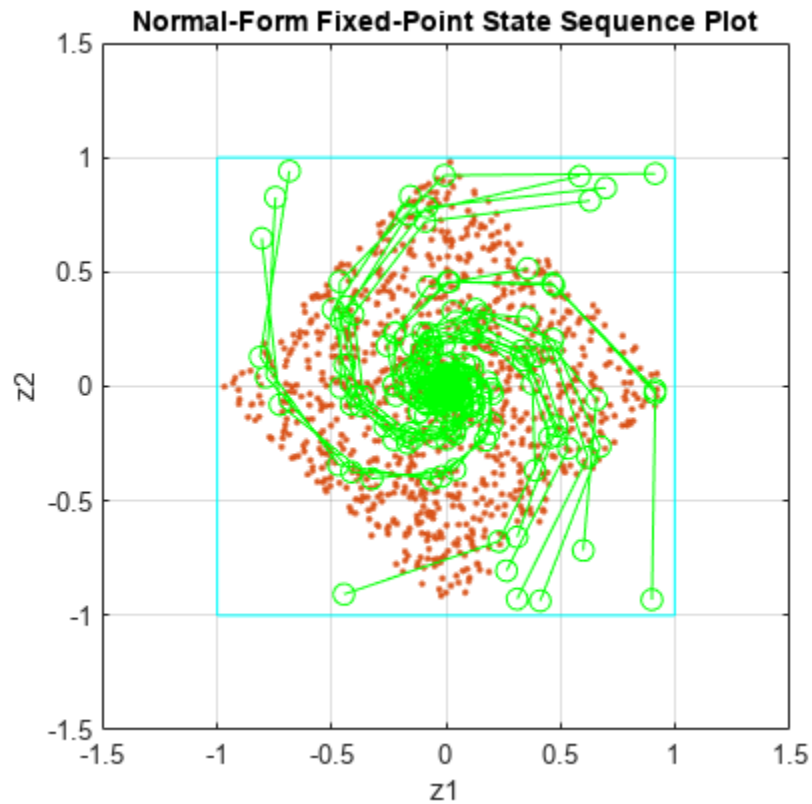
Now the projection of random initial states inside the unit square all contract uniformly. This is the result of the state transition matrix A_2 being normal. The states are also rotated by 90 degrees counterclockwise.

Plot the state sequences again for the same initial states as before. Observe that the outputs now spiral towards the origin.

```

x = fi(zeros(10,1),1,16,15,'fimath',F);
zi = fi([0;0],1,16,15,'fimath',F);
q = assignmentquantizer(zi);
e = double(eps(zi));
rng('default');
for k=1:20
    y = x;
    zi(:) = randquant(q,size(A,1),1);
    [y,zf] = fisisostatespacefilter(A2,B2,C2,D2,x,zi);
    if abs(double(zf(end)))>0.5, c='ro-'; else, c='go-'; end
    plot(zf(1,:), zf(2,:),c,'markersize',8);
end
title('Normal-Form Fixed-Point State Sequence Plot');
xlabel('z1'); ylabel('z2')

```



Trying this for other randomly chosen initial states illustrates that there is no region from which the filter is unable to recover.

```
 %#ok<*NASGU,*NOPTS>
```

References

[1] Richard A. Roberts and Clifford T. Mullis, "Digital Signal Processing", Addison-Wesley, Reading, Massachusetts, 1987, ISBN 0-201-16350-0, Section 9.3.

[2] S. K. Mitra, "Digital Signal Processing: A Computer Based Approach", McGraw-Hill, New York, 1998, ISBN 0-07-042953-7.

See Also

quantizer | fi | fimath

Develop Fixed-Point Algorithms

This example shows how to develop and verify a simple fixed-point algorithm. This example follows these steps for algorithm development:

- 1) Implement a second-order filter algorithm and simulate in double-precision floating-point.
- 2) Instrument the code to visualize the dynamic range of the output and state.
- 3) Convert the algorithm to fixed point by changing the data type of the variables. The algorithm itself does not change.
- 4) Compare and plot the fixed-point and floating-point results.

Define Double-Precision Floating-Point Variables

Develop the algorithm in double-precision floating-point. The algorithm used in this example is a second-order lowpass filter that removes the high frequencies in the input signal.

Define numerator coefficients, a , and denominator coefficients, b .

```
b = [ 0.25 0.5      0.25    ];
a = [ 1     0.09375 0.28125 ];
```

Generate a random input that has both high and low frequencies.

```
s = rng; rng(0, 'v5uniform');
x = randn(1000,1);
rng(s); % restore |rng| state
```

Pre-allocate the output, y , and state, z , for speed.

```
y = zeros(size(x));
z = [0;0];
```

Implement Data Type Independent Algorithm

This is a second-order filter that implements the standard difference equation:

$$|y(n) = b(1)*x(n) + b(2)*x(n-1) + b(3)*x(n-2) - a(2)*y(n-1) - a(3)*y(n-2)|$$

```
for k=1:length(x)
    y(k) = b(1)*x(k) + z(1);
    z(1) = (b(2)*x(k) + z(2)) - a(2)*y(k);
    z(2) = b(3)*x(k)          - a(3)*y(k);
end
```

Save the floating-point result.

```
ydouble = y;
```

Instrument Floating-Point Code to Visualize Dynamic Range

To convert to fixed point, we need to know the range of the variables. Depending on the complexity of an algorithm, this task can be simple or quite challenging. In this example, the range of the input value is known, so selecting an appropriate fixed-point data type is simple. Concentrate on the output

(y) and states (z) since their range is unknown. To view the dynamic range of the output and states, modify the code slightly to instrument it.

Create a `NumericTypeScope` object and view the dynamic range of states (z)

```
z = [0;0]; % Reset states

hscope1 = NumericTypeScope;

for k=1:length(x)
    y(k) = b(1)*x(k) + z(1);
    z(1) = (b(2)*x(k) + z(2)) - a(2)*y(k);
    z(2) = b(3)*x(k) - a(3)*y(k);

    % Process the data and update the visual
    step(hscope1,z);
end

% clear the information stored in the object hscope1
reset(hscope1);

% Create a |NumericTypeScope| object and view the dynamic range of the
% output (|y|)
hscope2 = NumericTypeScope;
step(hscope2,y);
```

First, analyze the information displayed for variable z (state). From the histogram, observe that the dynamic range lies within $[2^{(1)} \ 2^{(-12)}]$.

By default, the scope uses a word length of 16 bits with zero tolerable overflows. This results in a data type of `numericType(true,16,14)` since at least 2 integer bits are needed to avoid overflows. You can get more information on the statistical data from the Input Data and Resulting Type panels. The Input Data panel shows that the data has both positive and negative values and hence a signed quantity which is reflected in the suggested `numericType`. Also, the maximum data value is 1.51 which can be represented by the suggested type.

Next, look at variable y (output). From the histogram plot we see that the dynamic range lies within $[2^{(1)} \ 2^{(-13)}]$.

By default, the scope uses a word length of 16 bits with zero tolerable overflows. This results in a data type of `numericType(true,16,14)` since at least 2 integer bits are needed to avoid overflows. With this suggested type there are no overflows or underflows.

Define Fixed-Point Variables

Convert variables to use fixed-point data types and run the algorithm again.

Enable logging to see the overflows and underflows introduced by the selected data types.

```
FIPREF_STATE = get(fipref);
reset(fipref)
fp = fipref;
default_loggingmode = fp.LoggingMode;
fp.LoggingMode = 'On';
```

Capture the present state of and reset the global `fimath` to the factory settings.

```
globalFimathAtStart = fimath;
resetglobalfimath;
```

Define the fixed-point types for the variables in the below format: `fi(Data, Signed, WordLength, FractionLength)`

```
b = fi(b, 1, 8, 6);
a = fi(a, 1, 8, 6);

x = fi(x, 1, 16, 13);
y = fi(zeros(size(x)), 1, 16, 13);
z = fi([0;0], 1, 16, 14);
```

Implement the Same Data Type Independent Algorithm

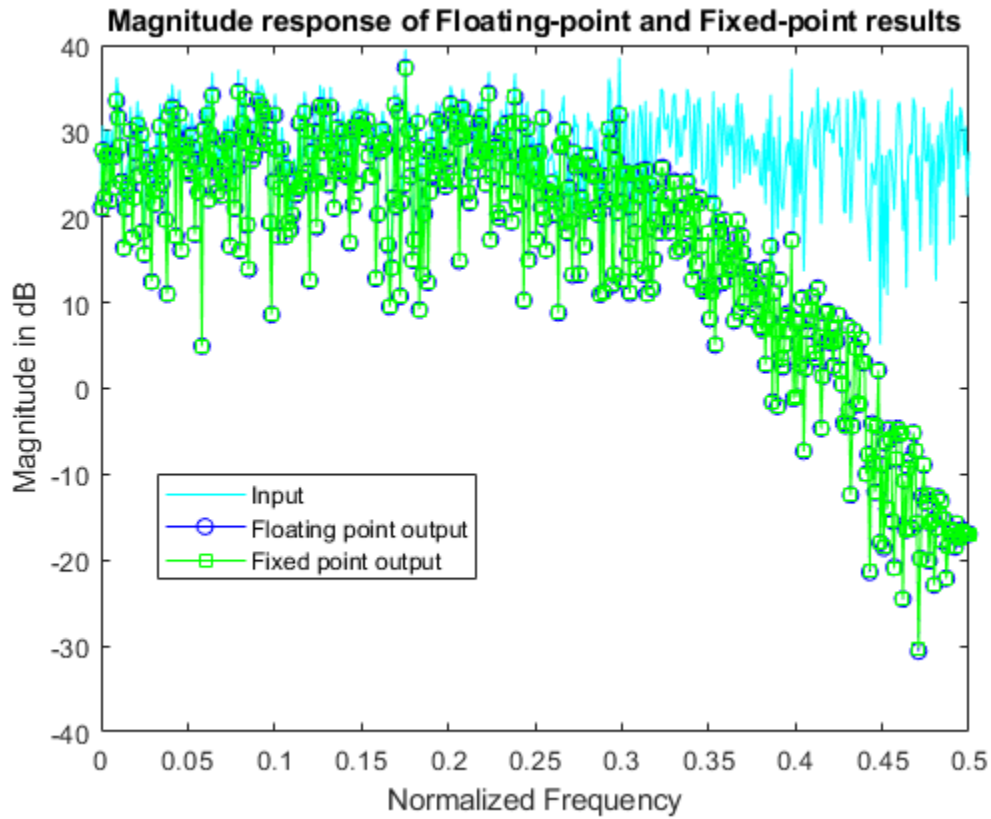
```
for k=1:length(x)
    y(k) = b(1)*x(k) + z(1);
    z(1) = (b(2)*x(k) + z(2)) - a(2)*y(k);
    z(2) = b(3)*x(k) - a(3)*y(k);
end
% Reset the logging mode.
fp.LoggingMode = default_loggingmode;
```

In this example, we have redefined the fixed-point variables with the same names as the floating-point so that we could inline the algorithm code for clarity. However, it is a better practice to enclose the algorithm code in a MATLAB® file function that could be called with either floating-point or fixed-point variables.

Compare and Plot the Floating-Point and Fixed-Point Results

Plot the magnitude response of the floating-point and fixed-point results and the response of the filter to see if the filter behaves as expected when it is converted to fixed-point.

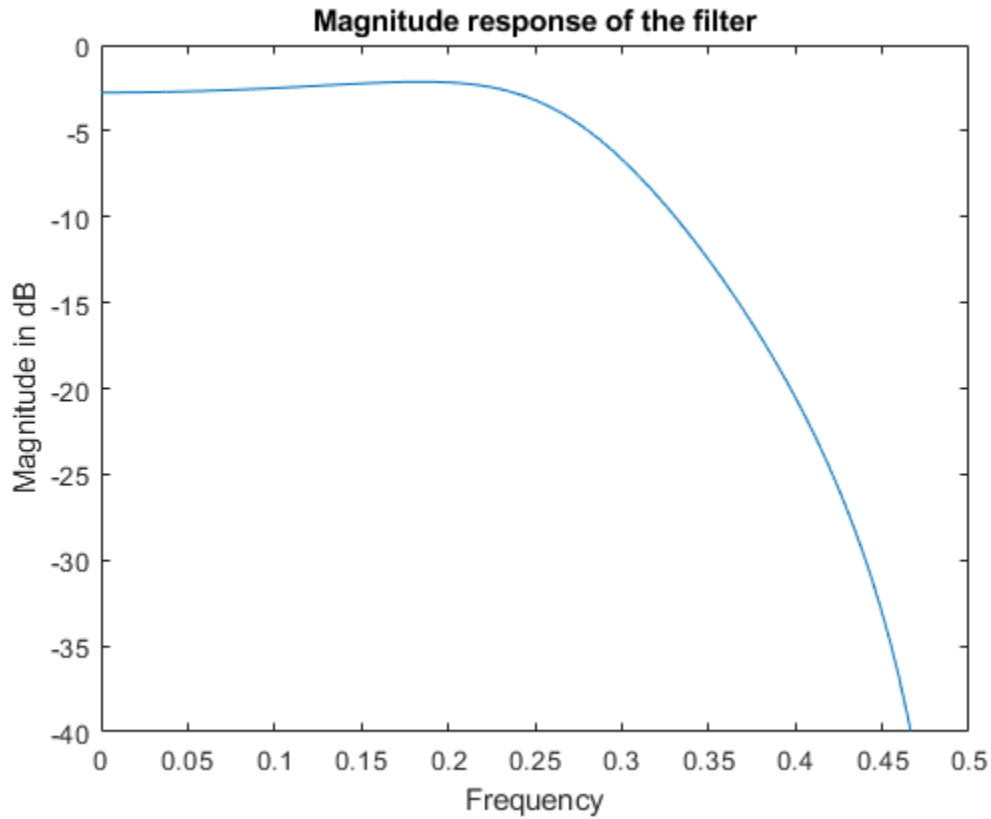
```
n = length(x);
f = linspace(0,0.5,n/2);
x_response = 20*log10(abs(fft(double(x))));
ydouble_response = 20*log10(abs(fft(ydouble)));
y_response = 20*log10(abs(fft(double(y))));
plot(f,x_response(1:n/2),'c-',...
     f,ydouble_response(1:n/2),'bo-',...
     f,y_response(1:n/2),'gs-');
ylabel('Magnitude in dB');
xlabel('Normalized Frequency');
legend('Input','Floating point output','Fixed point output','Location','Best');
title('Magnitude response of Floating-point and Fixed-point results');
```



```

h = fft(double(b),n)./fft(double(a),n);
h = h(1:end/2);
clf
hax = axes;
plot(hax,f,20*log10(abs(h)));
set(hax,'YLim',[-40 0]);
title('Magnitude response of the filter');
ylabel('Magnitude in dB')
xlabel('Frequency');

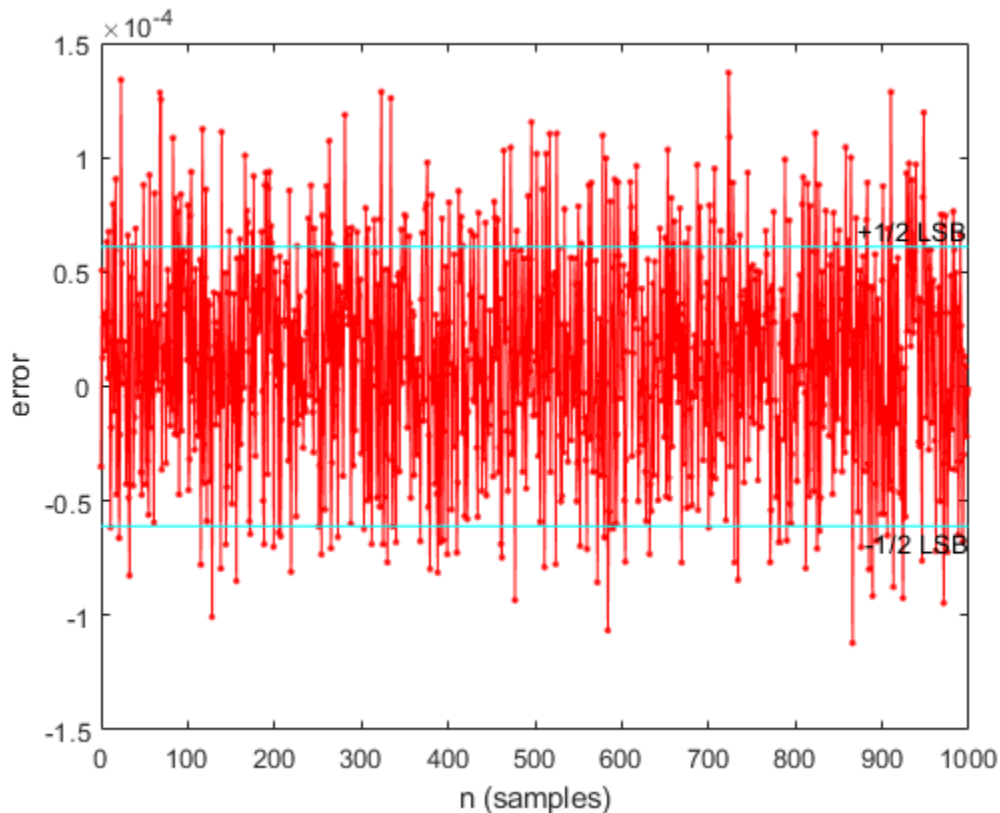
```



Notice that the high frequencies in the input signal are attenuated by the low-pass filter which is the expected behavior.

Plot the Error

```
clf
n = (0:length(y)-1)';
e = double(lsb(y));
plot(n,double(y)-ydouble,'-r', ...
      [n(1) n(end)],[e/2 e/2],'c', ...
      [n(1) n(end)],[-e/2 -e/2],'c')
text(n(end),e/2,'+1/2 LSB','HorizontalAlignment','right','VerticalAlignment','bottom')
text(n(end),-e/2,'-1/2 LSB','HorizontalAlignment','right','VerticalAlignment','top')
xlabel('n (samples)'); ylabel('error')
```



Implement the Algorithm in Simulink®

If you have Simulink® and Fixed-Point Designer™, you can run this model, which is the equivalent of the algorithm above. The output, `y_sim` is a fixed-point variable equal to the variable `y` calculated above in MATLAB code.

As in the MATLAB code, the fixed-point parameters in the blocks can be modified to match an actual system; these have been set to match the MATLAB code in the example above. Double-click on the blocks to see the settings.

```
% Set up the From Workspace variable
x_sim.time = n;
x_sim.signals.values = x;
x_sim.signals.dimensions = 1;

% Run the simulation
out_sim = sim('fitdf2filter_demo', 'SaveOutput', 'on', ...
             'SrcWorkspace', 'current');

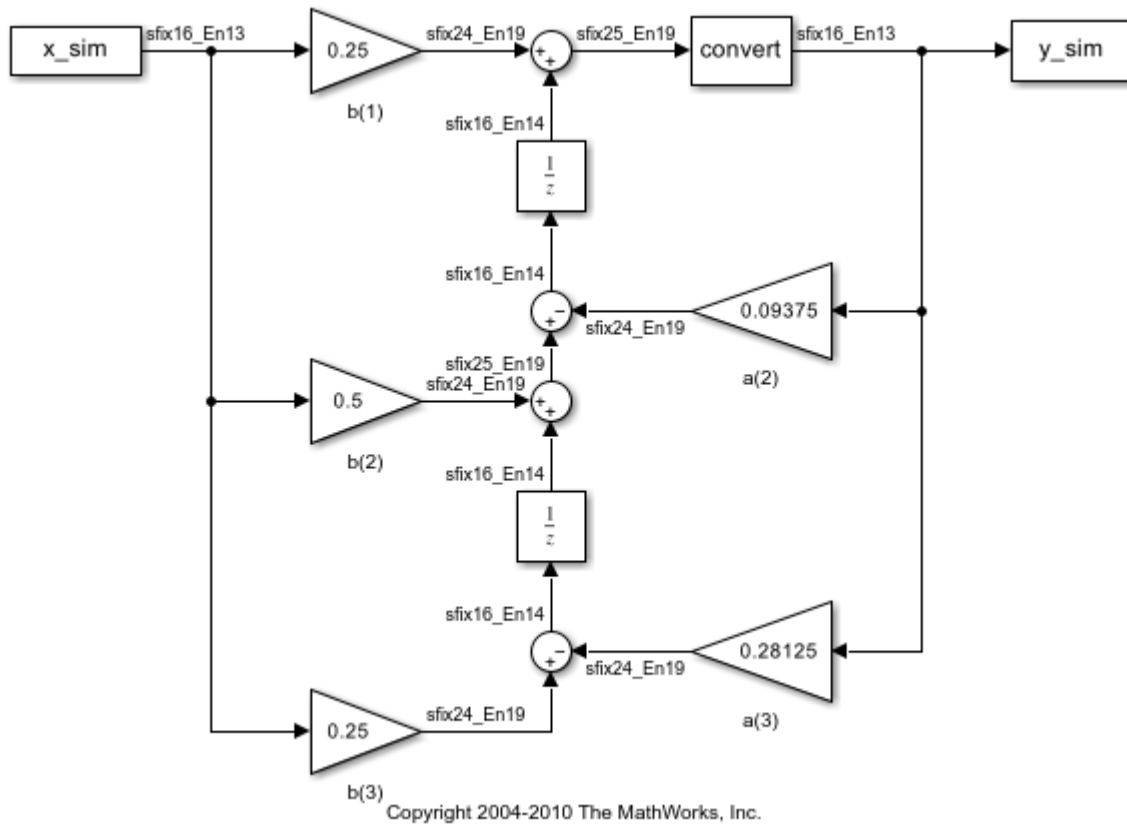
% Open the model
open_system('fitdf2filter_demo')

% Verify that the Simulink results are the same as the MATLAB file
isequal(y, out_sim.get('y_sim'))
```

ans =

logical

1



Assumptions Made for this Example

In order to simplify the example, we have taken the default math parameters: round-to-nearest, saturate on overflow, full precision products and sums. You can modify all of these parameters to match an actual system.

The settings were chosen as a starting point in algorithm development. Save a copy of this MATLAB file, start playing with the parameters, and see what effects they have on the output. How does the algorithm behave with a different input? See the help for `fi`, `fimath`, and `numericType` for information on how to set other parameters, such as rounding mode, and overflow mode.

```
% Reset the global fimath
globalfimath(globalFimathAtStart);
fipref(FIPREF_STATE);
```

Perform QR Factorization Using CORDIC

This example shows how to write MATLAB® code that works for both floating-point and fixed-point data types. The algorithm used in this example is the QR factorization implemented via CORDIC (Coordinate Rotation Digital Computer).

A good way to write an algorithm intended for a fixed-point target is to write it in MATLAB using builtin floating-point types so you can verify that the algorithm works. When you refine the algorithm to work with fixed-point types, then the best thing to do is to write it so that the same code continues working with floating-point. That way, when you are debugging, then you can switch the inputs back and forth between floating-point and fixed-point types to determine if a difference in behavior is because of fixed-point effects such as overflow and quantization versus an algorithmic difference. Even if the algorithm is not well suited for a floating-point target (as is the case of using CORDIC in the following example), it is still advantageous to have your MATLAB code work with floating-point for debugging purposes.

In contrast, you may have a completely different strategy if your target is floating point. For example, the QR algorithm is often done in floating-point with Householder transformations and row or column pivoting. But in fixed-point it is often more efficient to use CORDIC to apply Givens rotations with no pivoting.

This example addresses the first case, where your target is fixed-point, and you want an algorithm that is independent of data type because it is easier to develop and debug.

In this example you will learn various coding methods that can be applied across systems. The significant design patterns used in this example are the following:

- **Data Type Independence:** the algorithm is written in such a way that the MATLAB code is independent of data type, and will work equally well for fixed-point, double-precision floating-point, and single-precision floating-point.
- **Overflow Prevention:** method to guarantee not to overflow. This demonstrates how to prevent overflows in fixed-point.
- **Solving Systems of Equations:** method to use computational efficiency. Narrow your code scope by isolating what you need to define.

The main part in this example is an implementation of the QR factorization in fixed-point arithmetic using CORDIC for the Givens rotations. The algorithm is written in such a way that the MATLAB code is independent of data type, and will work equally well for fixed-point, double-precision floating-point, and single-precision floating-point.

The QR factorization of M-by-N matrix A produces an M-by-N upper triangular matrix R and an M-by-M orthogonal matrix Q such that $A = Q \cdot R$. A matrix is upper triangular if it has all zeros below the diagonal. An M-by-M matrix Q is orthogonal if $Q' \cdot Q = \text{eye}(M)$, the identity matrix.

The QR factorization is widely used in least-squares problems, such as the recursive least squares (RLS) algorithm used in adaptive filters.

The CORDIC algorithm is attractive for computing the QR algorithm in fixed-point because you can apply orthogonal Givens rotations with CORDIC using only shift and add operations.

Setup

So this example does not change your preferences or settings, we store the original state here, and restore them at the end.

```
originalFormat = get(0, 'format'); format short
originalFipref = get(fipref);      reset(fipref);
originalGlobalFimath = fimath;    resetglobalfimath;
```

Defining the CORDIC QR Algorithm

The CORDIC QR algorithm is given in the following MATLAB function, where A is an M-by-N real matrix, and niter is the number of CORDIC iterations. Output Q is an M-by-M orthogonal matrix, and R is an M-by-N upper-triangular matrix such that $Q \cdot R = A$.

```
function [Q,R] = cordicqr(A,niter)
    Kn = inverse_cordic_growth_constant(niter);
    [m,n] = size(A);
    R = A;
    Q = coder.nullcopy(repmat(A(:,1),1,m)); % Declare type and size of Q
    Q(:) = eye(m); % Initialize Q
    for j=1:n
        for i=j+1:m
            [R(j,j:end),R(i,j:end),Q(:,j),Q(:,i)] = ...
                cordicgivens(R(j,j:end),R(i,j:end),Q(:,j),Q(:,i),niter,Kn);
        end
    end
end
```

This function was written to be independent of data type. It works equally well with builtin floating-point types (double and single) and with the fixed-point `fi` object.

One of the trickiest aspects of writing data-type independent code is to specify data type and size for a new variable. In order to preserve data types without having to explicitly specify them, the output R was set to be the same as input A, like this:

```
R = A;
```

In addition to being data-type independent, this function was written in such a way that MATLAB Coder™ will be able to generate efficient C code from it. In MATLAB, you most often declare and initialize a variable in one step, like this:

```
Q = eye(m)
```

However, `Q=eye(m)` would always produce Q as a double-precision floating point variable. If A is fixed-point, then we want Q to be fixed-point; if A is single, then we want Q to be single; etc.

Hence, you need to declare the type and size of Q in one step, and then initialize it in a second step. This gives MATLAB Coder the information it needs to create an efficient C program with the correct types and sizes. In the finished code you initialize output Q to be an M-by-M identity matrix and the same data type as A, like this:

```
Q = coder.nullcopy(repmat(A(:,1),1,m)); % Declare type and size of Q
Q(:) = eye(m); % Initialize Q
```

The `coder.nullcopy` function declares the size and type of Q without initializing it. The expansion of the first column of A with `repmat` won't appear in code generated by MATLAB; it is only used to

specify the size. The `repmat` function was used instead of `A(:,1:m)` because `A` may have more rows than columns, which will be the case in a least-squares problem. You have to be sure to always assign values to every element of an array when you declare it with `coder.nullcopy`, because if you don't then you will have uninitialized memory.

You will notice this pattern of assignment again and again. This is another key enabler of data-type independent code.

The heart of this function is applying orthogonal Givens rotations in-place to the rows of `R` to zero out sub-diagonal elements, thus forming an upper-triangular matrix. The same rotations are applied in-place to the columns of the identity matrix, thus forming orthogonal `Q`. The Givens rotations are applied using the `cordicgivens` function, as defined in the next section. The rows of `R` and columns of `Q` are used as both input and output to the `cordicgivens` function so that the computation is done in-place, overwriting `R` and `Q`.

```
[R(j,j:end),R(i,j:end),Q(:,j),Q(:,i)] = ...
    cordicgivens(R(j,j:end),R(i,j:end),Q(:,j),Q(:,i),niter,Kn);
```

Defining the CORDIC Givens Rotation

The `cordicgivens` function applies a Givens rotation by performing CORDIC iterations to rows $x=R(j,j:end)$, $y=R(i,j:end)$ around the angle defined by $x(1)=R(j,j)$ and $y(1)=R(i,j)$ where $i>j$, thus zeroing out $R(i,j)$. The same rotation is applied to columns $u = Q(:,j)$ and $v = Q(:,i)$, thus forming the orthogonal matrix `Q`.

```
function [x,y,u,v] = cordicgivens(x,y,u,v,niter,Kn)
    if x(1)<0
        % Compensation for 3rd and 4th quadrants
        x(:) = -x; u(:) = -u;
        y(:) = -y; v(:) = -v;
    end
    for i=0:niter-1
        x0 = x;
        u0 = u;
        if y(1)<0
            % Counter-clockwise rotation
            % x and y form R, u and v form Q
            x(:) = x - bitsra(y, i); u(:) = u - bitsra(v, i);
            y(:) = y + bitsra(x0,i); v(:) = v + bitsra(u0,i);
        else
            % Clockwise rotation
            % x and y form R, u and v form Q
            x(:) = x + bitsra(y, i); u(:) = u + bitsra(v, i);
            y(:) = y - bitsra(x0,i); v(:) = v - bitsra(u0,i);
        end
    end
    % Set y(1) to exactly zero so R will be upper triangular without round off
    % showing up in the lower triangle.
    y(1) = 0;
    % Normalize the CORDIC gain
    x(:) = Kn * x; u(:) = Kn * u;
    y(:) = Kn * y; v(:) = Kn * v;
end
```

The advantage of using CORDIC in fixed-point over the standard Givens rotation is that CORDIC does not use square root or divide operations. Only bit-shifts, addition, and subtraction are needed in the

main loop, and one scalar-vector multiply at the end to normalize the CORDIC gain. Also, CORDIC rotations work well in pipelined architectures.

The bit shifts in each iteration are performed with the bit shift right arithmetic (`bitsra`) function instead of `bitshift`, multiplication by 0.5, or division by 2, because `bitsra`

- generates more efficient embedded code,
- works equally well with positive and negative numbers,
- works equally well with floating-point, fixed-point and integer types, and
- keeps this code independent of data type.

It is worthwhile to note that there is a difference between sub-scripted assignment (`subsasgn`) into a variable `a(:) = b` versus overwriting a variable `a = b`. Sub-scripted assignment into a variable like this

```
x(:) = x + bitsra(y, i);
```

always preserves the type of the left-hand-side argument `x`. This is the recommended programming style in fixed-point. For example fixed-point types often grow their word length in a sum, which is governed by the `SumMode` property of the `fimath` object, so that the right-hand-side `x + bitsra(y, i)` can have a different data type than `x`.

If, instead, you overwrite the left-hand-side like this

```
x = x + bitsra(y, i);
```

then the left-hand-side `x` takes on the type of the right-hand-side sum. This programming style leads to changing the data type of `x` in fixed-point code, and is discouraged.

Defining the Inverse CORDIC Growth Constant

This function returns the inverse of the CORDIC growth factor after `niter` iterations. It is needed because CORDIC rotations grow the values by a factor of approximately 1.6468, depending on the number of iterations, so the gain is normalized in the last step of `cordicgivens` by a multiplication by the inverse $K_n = 1/1.6468 = 0.60725$.

```
function Kn = inverse_cordic_growth_constant(niter)
    Kn = 1/prod(sqrt(1+2.^(-2*(0:double(niter)-1))));
end
```

Exploring CORDIC Growth as a Function of Number of Iterations

The function for CORDIC growth is defined as

```
growth = prod(sqrt(1+2.^(-2*(0:double(niter)-1))))
```

and the inverse is

```
inverse_growth = 1 ./ growth
```

Growth is a function of the number of iterations `niter`, and quickly converges to approximately 1.6468, and the inverse converges to approximately 0.60725. You can see in the following table that the difference from one iteration to the next ceases to change after 27 iterations. This is because the calculation hit the limit of precision in double floating-point at 27 iterations.

| niter | growth | diff(growth) | 1./growth | diff(1./growth) |
|-------|--------------------|--------------|--------------------|-----------------|
| 0 | 1.0000000000000000 | 0 | 1.0000000000000000 | 0 |

```

1  1.414213562373095  0.414213562373095  0.707106781186547  -0.292893218813453
2  1.581138830084190  0.166925267711095  0.632455532033676  -0.074651249152872
3  1.629800601300662  0.048661771216473  0.613571991077896  -0.018883540955780
4  1.642484065752237  0.012683464451575  0.608833912517752  -0.004738078560144
5  1.645688915757255  0.003204850005018  0.607648256256168  -0.001185656261584
6  1.646492278712479  0.000803362955224  0.607351770141296  -0.000296486114872
7  1.646693254273644  0.000200975561165  0.607277644093526  -0.000074126047770
8  1.646743506596901  0.000050252323257  0.607259112298893  -0.000018531794633
9  1.646756070204878  0.000012563607978  0.607254479332562  -0.000004632966330
10 1.646759211139822  0.000003140934944  0.607253321089875  -0.000001158242687
11 1.646759996375617  0.000000785235795  0.607253031529134  -0.000000289560741
12 1.646760192684695  0.000000196309077  0.607252959138945  -0.000000072390190
13 1.646760241761972  0.000000049077277  0.607252941041397  -0.000000018097548
14 1.646760254031292  0.000000012269320  0.607252936517010  -0.000000004524387
15 1.646760257098622  0.000000003067330  0.607252935385914  -0.000000001131097
16 1.646760257865455  0.000000000766833  0.607252935103139  -0.000000000282774
17 1.646760258057163  0.000000000191708  0.607252935032446  -0.000000000070694
18 1.646760258105090  0.000000000047927  0.607252935014772  -0.000000000017673
19 1.646760258117072  0.000000000011982  0.607252935010354  -0.000000000004418
20 1.646760258120067  0.000000000002995  0.607252935009249  -0.000000000001105
21 1.646760258120816  0.000000000000749  0.607252935008973  -0.000000000000276
22 1.646760258121003  0.000000000000187  0.607252935008904  -0.000000000000069
23 1.646760258121050  0.000000000000047  0.607252935008887  -0.000000000000017
24 1.646760258121062  0.000000000000012  0.607252935008883  -0.000000000000004
25 1.646760258121065  0.000000000000003  0.607252935008882  -0.000000000000001
26 1.646760258121065  0.000000000000001  0.607252935008881  -0.000000000000000
27 1.646760258121065  0 0.607252935008881 0
28 1.646760258121065  0 0.607252935008881 0
29 1.646760258121065  0 0.607252935008881 0
30 1.646760258121065  0 0.607252935008881 0
31 1.646760258121065  0 0.607252935008881 0
32 1.646760258121065  0 0.607252935008881 0

```

Comparing CORDIC to the Standard Givens Rotation

The `cordicgivens` function is numerically equivalent to the following standard Givens rotation algorithm from Golub & Van Loan, *Matrix Computations*. In the `cordicqr` function, if you replace the call to `cordicgivens` with a call to `givensrotation`, then you will have the standard Givens QR algorithm.

```

function [x,y,u,v] = givensrotation(x,y,u,v)
    a = x(1); b = y(1);
    if b==0
        % No rotation necessary.  c = 1; s = 0;
        return;
    else
        if abs(b) > abs(a)
            t = -a/b; s = 1/sqrt(1+t^2); c = s*t;
        else
            t = -b/a; c = 1/sqrt(1+t^2); s = c*t;
        end
    end
    x0 = x;          u0 = u;
    % x and y form R,  u and v form Q
    x(:) = c*x0 - s*y;  u(:) = c*u0 - s*v;
    y(:) = s*x0 + c*y;  v(:) = s*u0 + c*v;
end

```

The givens rotation function uses division and square root, which are expensive in fixed-point, but good for floating-point algorithms.

Example of CORDIC Rotations

Here is a 3-by-3 example that follows the CORDIC rotations through each step of the algorithm. The algorithm uses orthogonal rotations to zero out the subdiagonal elements of R using the diagonal elements as pivots. The same rotations are applied to the identity matrix, thus producing orthogonal Q such that $Q^*R = A$.

Let A be a random 3-by-3 matrix, and initialize $R = A$, and $Q = \text{eye}(3)$.

$$R = A = \begin{bmatrix} -0.8201 & 0.3573 & -0.0100 \\ -0.7766 & -0.0096 & -0.7048 \\ -0.7274 & -0.6206 & -0.8901 \end{bmatrix}$$

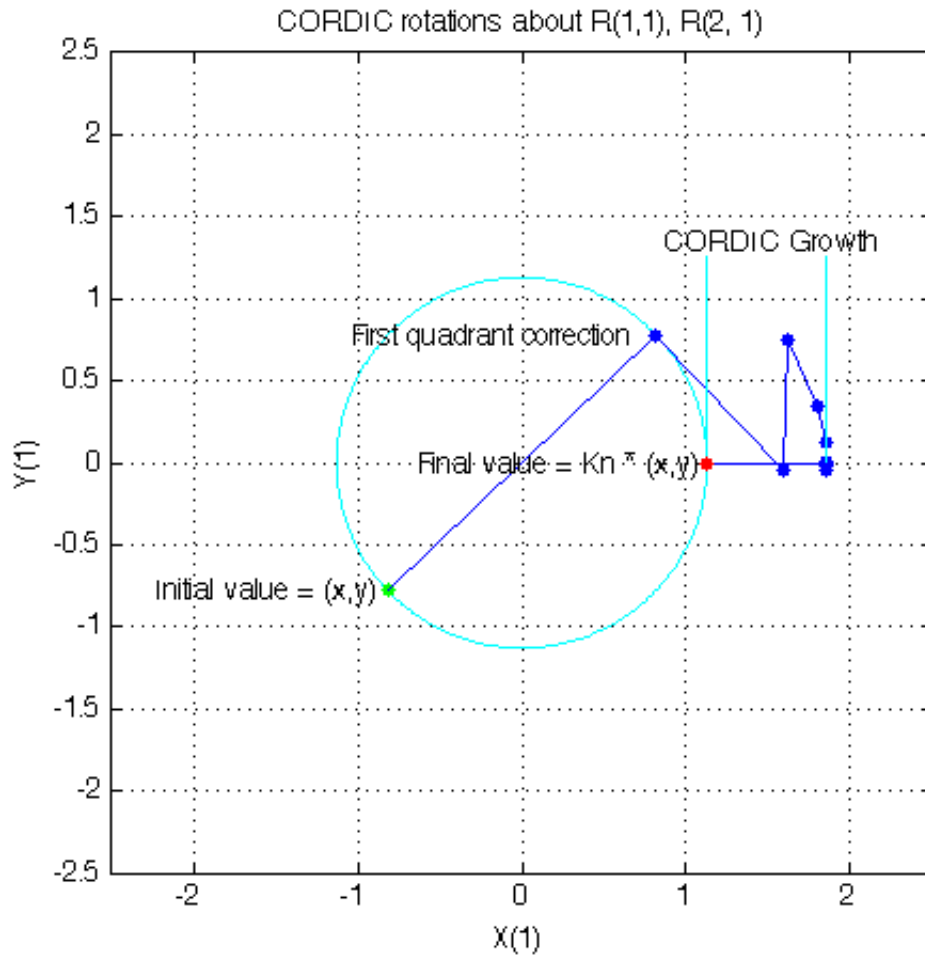
$$Q = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

The first rotation is about the first and second row of R and the first and second column of Q. Element $R(1,1)$ is the pivot and $R(2,1)$ rotates to 0.

$$\begin{array}{l} \text{R before the first rotation} \\ x \begin{bmatrix} -0.8201 & 0.3573 & -0.0100 \end{bmatrix} \\ y \begin{bmatrix} -0.7766 & -0.0096 & -0.7048 \\ -0.7274 & -0.6206 & -0.8901 \end{bmatrix} \end{array} \rightarrow \begin{array}{l} \text{R after the first rotation} \\ x \begin{bmatrix} 1.1294 & -0.2528 & 0.4918 \end{bmatrix} \\ y \begin{bmatrix} 0 & 0.2527 & 0.5049 \\ -0.7274 & -0.6206 & -0.8901 \end{bmatrix} \end{array}$$

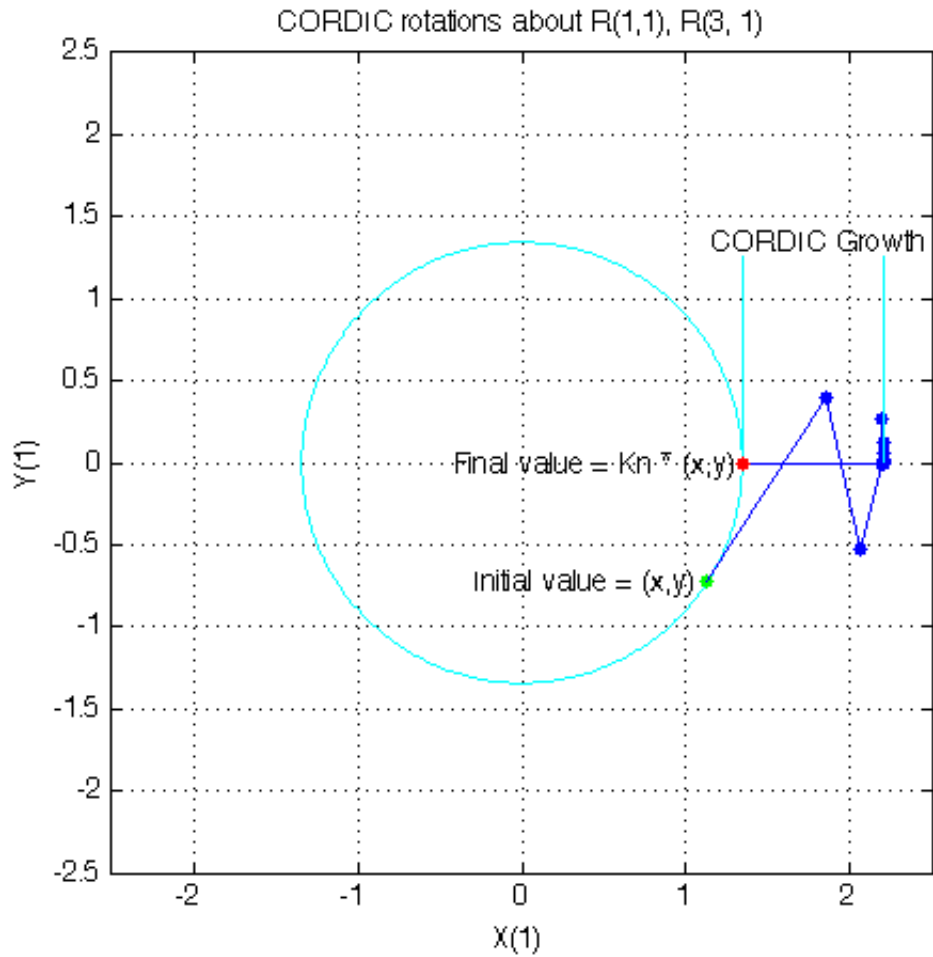
$$\begin{array}{l} \text{Q before the first rotation} \\ u \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \\ v \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \end{array} \rightarrow \begin{array}{l} \text{Q after the first rotation} \\ u \begin{bmatrix} -0.7261 \\ -0.6876 \\ 0 \end{bmatrix} \\ v \begin{bmatrix} 0.6876 \\ -0.7261 \\ 0 \end{bmatrix} \end{array}$$

In the following plot, you can see the growth in x in each of the CORDIC iterations. The growth is factored out at the last step by multiplying it by $K_n = 0.60725$. You can see that $y(1)$ iterates to 0. Initially, the point $[x(1), y(1)]$ is in the third quadrant, and is reflected into the first quadrant before the start of the CORDIC iterations.



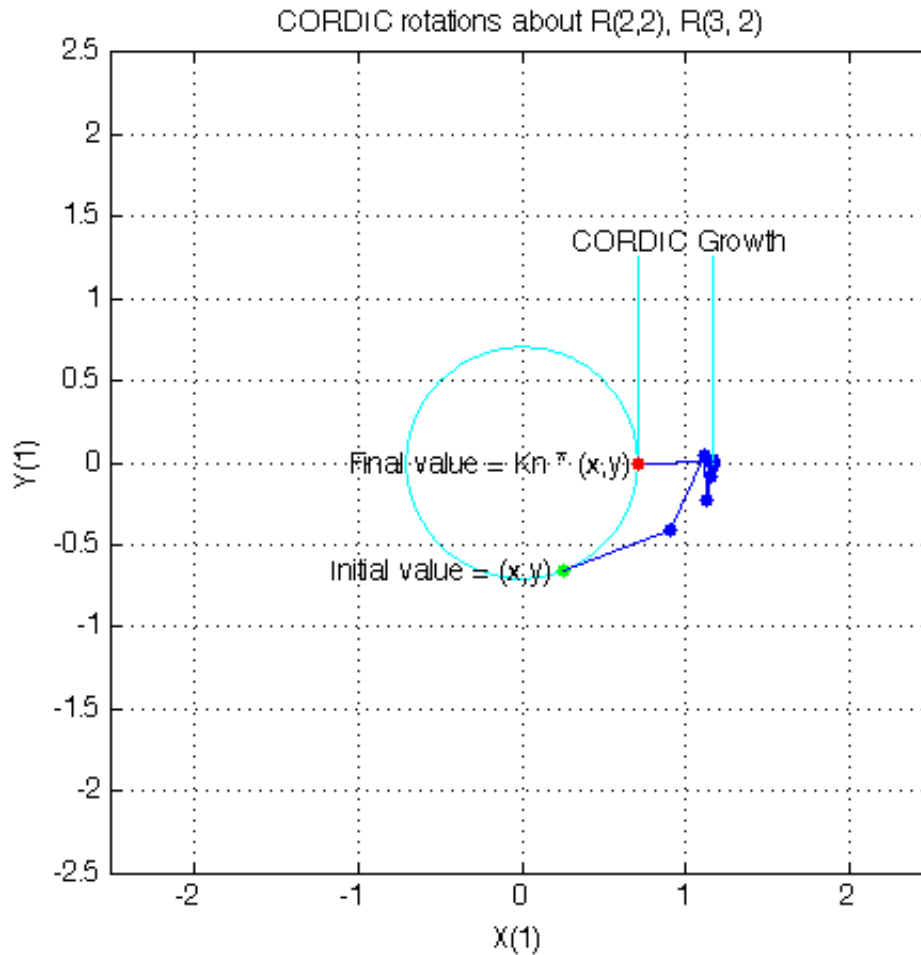
The second rotation is about the first and third row of R and the first and third column of Q. Element R(1, 1) is the pivot and R(3, 1) rotates to 0.

| | | | | | |
|------------------------------|-------------------------------|-----|-----------------------------|----------------------------|-----------|
| R before the second rotation | | -> | R after the second rotation | | |
| x | [1.1294 -0.2528 0.4918] | | x | [1.3434 0.1235 0.8954] | |
| | 0 0.2527 0.5049 | | | 0 0.2527 0.5049 | |
| y | [-0.7274] -0.6206 -0.8901 | -> | y | [0 -0.6586 -0.4820] | |
| Q before the second rotation | | | Q after the second rotation | | |
| u | | | u | | |
| | v | | | v | |
| [-0.7261] | 0.6876 | [0] | [-0.6105] | 0.6876 | [-0.3932] |
| [-0.6876] | -0.7261 | [0] | [-0.5781] | -0.7261 | [-0.3723] |
| [0] | 0 | [1] | [-0.5415] | 0 | [0.8407] |



The third rotation is about the second and third row of R and the second and third column of Q. Element R(2, 2) is the pivot and R(3, 2) rotates to 0.

| | | | | | | | | | |
|---|-----------------------------|-----------|-----------|----|---|----------------------------|-----------|-----------|--|
| | R before the third rotation | | | | | R after the third rotation | | | |
| | 1.3434 | 0.1235 | 0.8954 | | | 1.3434 | 0.1235 | 0.8954 | |
| x | 0 | [0.2527 | 0.5049] | -> | x | 0 | [0.7054 | 0.6308] | |
| y | 0 | [-0.6586 | -0.4820] | -> | y | 0 | [0 | 0.2987] | |
| | Q before the third rotation | | | | | Q after the third rotation | | | |
| | | u | v | | | | u | v | |
| | -0.6105 | [0.6876] | [-0.3932] | | | -0.6105 | [0.6134] | [0.5011] | |
| | -0.5781 | [-0.7261] | [-0.3723] | -> | | -0.5781 | [0.0875] | [-0.8113] | |
| | -0.5415 | [0 | 0.8407] | | | -0.5415 | [-0.7849] | [0.3011] | |



This completes the QR factorization. R is upper triangular, and Q is orthogonal.

$$R = \begin{bmatrix} 1.3434 & 0.1235 & 0.8954 \\ 0 & 0.7054 & 0.6308 \\ 0 & 0 & 0.2987 \end{bmatrix}$$

$$Q = \begin{bmatrix} -0.6105 & 0.6134 & 0.5011 \\ -0.5781 & 0.0875 & -0.8113 \\ -0.5415 & -0.7849 & 0.3011 \end{bmatrix}$$

You can verify that Q is within roundoff error of being orthogonal by multiplying and seeing that it is close to the identity matrix.

$$Q*Q' = \begin{bmatrix} 1.0000 & 0.0000 & 0.0000 \\ 0.0000 & 1.0000 & 0 \\ 0.0000 & 0 & 1.0000 \end{bmatrix}$$

$$Q'*Q = \begin{bmatrix} 1.0000 & 0.0000 & -0.0000 \\ 0.0000 & 1.0000 & -0.0000 \\ -0.0000 & -0.0000 & 1.0000 \end{bmatrix}$$

You can see the error difference by subtracting the identity matrix.

```
Q*Q' - eye(size(Q)) =
           0  2.7756e-16  3.0531e-16
  2.7756e-16  4.4409e-16  0
  3.0531e-16  0  6.6613e-16
```

You can verify that Q^*R is close to A by subtracting to see the error difference.

```
Q*R - A = -3.7802e-11 -7.2325e-13 -2.7756e-17
          -3.0512e-10  1.1708e-12 -4.4409e-16
          3.6836e-10 -4.3487e-13 -7.7716e-16
```

Determining the Optimal Output Type of Q for Fixed Word Length

Since Q is orthogonal, you know that all of its values are between -1 and $+1$. In floating-point, there is no decision about the type of Q : it should be the same floating-point type as A . However, in fixed-point, you can do better than making Q have the identical fixed-point type as A . For example, if A has word length 16 and fraction length 8, and if we make Q also have word length 16 and fraction length 8, then you force Q to be less accurate than it could be and waste the upper half of the fixed-point range.

The best type for Q is to make it have full range of its possible outputs, plus accommodate the 1.6468 CORDIC growth factor in intermediate calculations. Therefore, assuming that the word length of Q is the same as the word length of input A , then the best fraction length for Q is 2 bits less than the word length (one bit for 1.6468 and one bit for the sign).

Hence, our initialization of Q in `cordicqr` can be improved like this.

```
if isfi(A) && (isfixed(A) || isscaleddouble(A))
    Q = fi(one*eye(m), get(A, 'NumericType'), ...
          'FractionLength', get(A, 'WordLength')-2);
else
    Q = coder.nullcopy(repmat(A(:,1),1,m));
    Q(:) = eye(m);
end
```

A slight disadvantage is that this section of code is dependent on data type. However, you gain a major advantage by picking the optimal type for Q , and the main algorithm is still independent of data type. You can do this kind of input parsing in the beginning of a function and leave the main algorithm data-type independent.

Preventing Overflow in Fixed Point R

This section describes how to determine a fixed-point output type for R in order to prevent overflow. In order to pick an output type, you need to know how much the magnitude of the values of R will grow.

Given real matrix A and its QR factorization computed by Givens rotations without pivoting, an upper-bound on the magnitude of the elements of R is the square-root of the number of rows of A times the magnitude of the largest element in A . Furthermore, this growth will never be greater during an intermediate computation. In other words, let $[m,n]=\text{size}(A)$, and $[Q,R]=\text{givensqr}(A)$. Then

```
max(abs(R(:))) <= sqrt(m) * max(abs(A(:))).
```

This is true because each element of R is formed from orthogonal rotations from its corresponding column in A , so the largest that any element $R(i,j)$ can get is if all of the elements of

its corresponding column $A(:, j)$ were rotated to a single value. In other words, the largest possible value will be bounded by the 2-norm of $A(:, j)$. Since the 2-norm of $A(:, j)$ is equal to the square-root of the sum of the squares of the m elements, and each element is less-than-or-equal-to the largest element of A , then

$$\text{norm}(A(:, j)) \leq \sqrt{m} * \max(\text{abs}(A(:))).$$

That is, for all j

$$\begin{aligned} \text{norm}(A(:, j)) &= \sqrt{A(1, j)^2 + A(2, j)^2 + \dots + A(m, j)^2} \\ &\leq \sqrt{m * \max(\text{abs}(A(:)))^2} \\ &= \sqrt{m} * \max(\text{abs}(A(:))). \end{aligned}$$

and so for all i, j

$$\text{abs}(R(i, j)) \leq \text{norm}(A(:, j)) \leq \sqrt{m} * \max(\text{abs}(A(:))).$$

Hence, it is also true for the largest element of R

$$\max(\text{abs}(R(:))) \leq \sqrt{m} * \max(\text{abs}(A(:))).$$

This becomes useful in fixed-point where the elements of A are often very close to the maximum value attainable by the data type, so we can set a tight upper bound without knowing the values of A . This is important because we want to set an output type for R with a minimum number of bits, only knowing the upper bound of the data type of A . You can use `fi` method `upperbound` to get this value.

Therefore, for all i, j

$$\text{abs}(R(i, j)) \leq \sqrt{m} * \text{upperbound}(A)$$

Note that $\sqrt{m} * \text{upperbound}(A)$ is also an upper bound for the elements of A :

$$\text{abs}(A(i, j)) \leq \text{upperbound}(A) \leq \sqrt{m} * \text{upperbound}(A)$$

Therefore, when picking fixed-point data types, $\sqrt{m} * \text{upperbound}(A)$ is an upper bound that will work for both A and R .

Attaining the maximum is easy and common. The maximum will occur when all elements get rotated into a single element, like the following matrix with orthogonal columns:

$$A = \begin{bmatrix} 7 & -7 & 7 & 7 \\ 7 & 7 & -7 & 7 \\ 7 & -7 & -7 & -7 \\ 7 & 7 & 7 & -7 \end{bmatrix};$$

Its maximum value is 7 and its number of rows is $m=4$, so we expect that the maximum value in R will be bounded by $\max(\text{abs}(A(:))) * \sqrt{m} = 7 * \sqrt{4} = 14$. Since A in this example is orthogonal, each column gets rotated to the max value on the diagonal.

```
niter = 52;
[Q,R] = cordicqr(A,niter)
```

Q =

$$\begin{bmatrix} 0.5000 & -0.5000 & 0.5000 & 0.5000 \\ 0.5000 & 0.5000 & -0.5000 & 0.5000 \\ 0.5000 & -0.5000 & -0.5000 & -0.5000 \end{bmatrix}$$

```
0.5000  0.5000  0.5000 -0.5000
```

R =

```
14.0000  0.0000 -0.0000 -0.0000
      0  14.0000 -0.0000  0.0000
      0      0  14.0000  0.0000
      0      0      0  14.0000
```

Another simple example of attaining maximum growth is a matrix that has all identical elements, like a matrix of all ones. A matrix of ones will get rotated into $1*\sqrt{m}$ in the first row and zeros elsewhere. For example, this 9-by-5 matrix will have all $1*\sqrt{9}=3$ in the first row of R.

```
m = 9; n = 5;
A = ones(m,n)
niter = 52;
[Q,R] = cordicqr(A,niter)
```

A =

```
1  1  1  1  1
1  1  1  1  1
1  1  1  1  1
1  1  1  1  1
1  1  1  1  1
1  1  1  1  1
1  1  1  1  1
1  1  1  1  1
1  1  1  1  1
```

Q =

Columns 1 through 7

```
0.3333  0.5567 -0.6784  0.3035 -0.1237  0.0503  0.0158
0.3333  0.0296  0.2498 -0.1702 -0.6336  0.1229 -0.3012
0.3333  0.2401  0.0562 -0.3918  0.4927  0.2048 -0.5395
0.3333  0.0003  0.0952 -0.1857  0.2148  0.4923  0.7080
0.3333  0.1138  0.0664 -0.2263  0.1293 -0.8348  0.2510
0.3333 -0.3973 -0.0143  0.3271  0.4132 -0.0354 -0.2165
0.3333  0.1808  0.3538 -0.1012 -0.2195      0  0.0824
0.3333 -0.6500 -0.4688 -0.2380 -0.2400      0      0
0.3333 -0.0740  0.3400  0.6825 -0.0331      0      0
```

Columns 8 through 9

```
0.0056 -0.0921
-0.5069 -0.1799
0.0359  0.3122
-0.2351 -0.0175
-0.2001  0.0610
-0.0939 -0.6294
0.7646 -0.2849
0.2300  0.2820
```

```

0      0.5485

R =

3.0000    3.0000    3.0000    3.0000    3.0000
0         0.0000    0.0000    0.0000    0.0000
0         0         0.0000    0.0000    0.0000
0         0         0         0.0000    0.0000
0         0         0         0         0.0000
0         0         0         0         0
0         0         0         0         0
0         0         0         0         0
0         0         0         0         0

```

As in the `cordicqr` function, the Givens QR algorithm is often written by overwriting `A` in-place with `R`, so being able to cast `A` into `R`'s data type at the beginning of the algorithm is convenient.

In addition, if you compute the Givens rotations with CORDIC, there is a growth-factor that converges quickly to approximately 1.6468. This growth factor gets normalized out after each Givens rotation, but you need to accommodate it in the intermediate calculations. Therefore, the number of additional bits that are required including the Givens and CORDIC growth are $\log_2(1.6468 * \text{sqrt}(m))$. The additional bits of head-room can be added either by increasing the word length, or decreasing the fraction length.

A benefit of increasing the word length is that it allows for the maximum possible precision for a given word length. A disadvantage is that the optimal word length may not correspond to a native type on your processor (e.g. increasing from 16 to 18 bits), or you may have to increase to the next larger native word size which could be quite large (e.g. increasing from 16 to 32 bits, when you only needed 18).

A benefit of decreasing fraction length is that you can do the computation in-place in the native word size of `A`. A disadvantage is that you lose precision.

Another option is to pre-scale the input by right-shifting. This is equivalent to decreasing the fraction length, with the additional disadvantage of changing the scaling of your problem. However, this may be an attractive option to you if you prefer to only work in fractional arithmetic or integer arithmetic.

Example of Fixed Point Growth in R

If you have a fixed-point input matrix `A`, you can define fixed-point output `R` with the growth defined in the previous section.

Start with a random matrix `X`.

```

X = [0.0513    -0.2097    0.9492    0.2614
     0.8261     0.6252    0.3071   -0.9415
     1.5270     0.1832    0.1352   -0.1623
     0.4669    -1.0298    0.5152   -0.1461];

```

Create a fixed-point `A` from `X`.

```
A = sfi(X)
```

```
A =
```

```

0.0513   -0.2097   0.9492   0.2614
0.8261   0.6252   0.3071  -0.9415
1.5270   0.1832   0.1352  -0.1623
0.4669  -1.0298   0.5152  -0.1461

```

```

        DataTypeMode: Fixed-point: binary point scaling
          Signedness: Signed
          WordLength: 16
        FractionLength: 14

```

```
m = size(A,1)
```

```
m =
```

```
4
```

The growth factor is 1.6468 times the square-root of the number of rows of A. The bit growth is the next integer above the base-2 logarithm of the growth.

```
bit_growth = ceil(log2(cordic_growth_constant * sqrt(m)))
```

```
bit_growth =
```

```
2
```

Initialize R with the same values as A, and a word length increased by the bit growth.

```
R = sfi(A, get(A, 'WordLength')+bit_growth, get(A, 'FractionLength'))
```

```
R =
```

```

0.0513   -0.2097   0.9492   0.2614
0.8261   0.6252   0.3071  -0.9415
1.5270   0.1832   0.1352  -0.1623
0.4669  -1.0298   0.5152  -0.1461

```

```

        DataTypeMode: Fixed-point: binary point scaling
          Signedness: Signed
          WordLength: 18
        FractionLength: 14

```

Use R as input and overwrite it.

```
niter = get(R, 'WordLength') - 1
[Q,R] = cordicqr(R, niter)
```

```
niter =
```

```
17
```

```
Q =
```

```

0.0284   -0.1753    0.9110    0.3723
0.4594    0.4470    0.3507   -0.6828
0.8490    0.0320   -0.2169    0.4808
0.2596   -0.8766   -0.0112   -0.4050

```

```

DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 18
FractionLength: 16

```

R =

```

1.7989    0.1694    0.4166   -0.6008
0         1.2251   -0.4764   -0.3438
0         0         0.9375   -0.0555
0         0         0         0.7214

```

```

DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 18
FractionLength: 14

```

Verify that $Q*Q'$ is near the identity matrix.

```
double(Q)*double(Q')
```

ans =

```

1.0000   -0.0001    0.0000    0.0000
-0.0001    1.0001    0.0000   -0.0000
0.0000    0.0000    1.0000   -0.0000
0.0000   -0.0000   -0.0000    1.0000

```

Verify that $Q*R - A$ is small relative to the precision of A.

```
err = double(Q)*double(R) - double(A)
```

err =

```

1.0e-03 *
-0.1048   -0.2355    0.1829   -0.2146
0.3472    0.2949    0.0260   -0.2570
0.2776   -0.1740   -0.1007    0.0966
0.0138   -0.1558    0.0417   -0.0362

```

Increasing Precision in R

The previous section showed you how to prevent overflow in R while maintaining the precision of A. If you leave the fraction length of R the same as A, then R cannot have more precision than A, and your precision requirements may be such that the precision of R must be greater.

An extreme example of this is to define a matrix with an integer fixed-point type (i.e. fraction length is zero). Let matrix X have elements that are the full range for signed 8 bit integers, between -128 and +127.

```
X = [-128 -128 -128 127
     -128 127 127 -128
      127 127 127 127
      127 127 -128 -128];
```

Define fixed-point A to be equivalent to an 8-bit integer.

```
A = sfi(X,8,0)
```

```
A =
```

```
-128 -128 -128 127
-128 127 127 -128
 127 127 127 127
 127 127 -128 -128
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 8
      FractionLength: 0
```

```
m = size(A,1)
```

```
m =
```

```
4
```

The necessary growth is 1.6468 times the square-root of the number of rows of A.

```
bit_growth = ceil(log2(cordic_growth_constant*sqrt(m)))
```

```
bit_growth =
```

```
2
```

Initialize R with the same values as A, and allow for bit growth like you did in the previous section.

```
R = sfi(A, get(A, 'WordLength')+bit_growth, get(A, 'FractionLength'))
```

```
R =
```

```
-128 -128 -128 127
-128 127 127 -128
 127 127 127 127
 127 127 -128 -128
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 10
      FractionLength: 0
```

Compute the QR factorization, overwriting R.

```
niter = get(R, 'WordLength') - 1;
[Q,R] = cordicqr(R, niter)
```

Q =

```
-0.5039  -0.2930  -0.4062  -0.6914
-0.5039   0.8750   0.0039   0.0078
 0.5000   0.2930   0.3984  -0.7148
 0.4922   0.2930  -0.8203   0.0039
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 10
      FractionLength: 8
```

R =

```
257  126   -1   -1
  0  225  151 -148
  0   0  211  104
  0   0   0 -180
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 10
      FractionLength: 0
```

Notice that R is returned with integer values because you left the fraction length of R at 0, the same as the fraction length of A.

The scaling of the least-significant bit (LSB) of A is 1, and you can see that the error is proportional to the LSB.

```
err = double(Q)*double(R)-double(A)
```

err =

```
-1.5039  -1.4102  -1.4531  -0.9336
-1.5039   6.3828   6.4531  -1.9961
 1.5000   1.9180   0.8086  -0.7500
-0.5078   0.9336  -1.3398  -1.8672
```

You can increase the precision in the QR factorization by increasing the fraction length. In this example, you needed 10 bits for the integer part (8 bits to start with, plus 2 bits growth), so when you increase the fraction length you still need to keep the 10 bits in the integer part. For example, you can increase the word length to 32 and set the fraction length to 22, which leaves 10 bits in the integer part.

```
R = sfi(A, 32, 22)
```

R =


```
-128 -128 -128 127
-128 127 127 -128
127 127 127 127
127 127 -128 -128
```

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 32
FractionLength: 22
```

```
niter = get(R, 'WordLength') - 1;
[Q,R] = cordicqr(R, niter)
```

Q =

```
-0.5020 -0.2913 -0.4088 -0.7043
-0.5020 0.8649 0.0000 0.0000
0.4980 0.2890 0.4056 -0.7099
0.4980 0.2890 -0.8176 0.0000
```

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 32
FractionLength: 30
```

R =

```
255.0020 127.0029 0.0039 0.0039
0 220.5476 146.8413 -147.9930
0 0 208.4793 104.2429
0 0 0 -179.6037
```

```
DataTypeMode: Fixed-point: binary point scaling
Signedness: Signed
WordLength: 32
FractionLength: 22
```

Now you can see fractional parts in R, and $Q \cdot R - A$ is small.

```
err = double(Q)*double(R)-double(A)
```

err =

```
1.0e-05 *
-0.1234 -0.0014 -0.0845 0.0267
-0.1234 0.2574 0.1260 -0.1094
0.0720 0.0289 -0.0400 -0.0684
0.0957 0.0818 -0.1034 0.0095
```

The number of bits you choose for fraction length will depend on the precision requirements for your particular algorithm.

Picking Default Number of Iterations

The number of iterations is dependent on the desired precision, but limited by the word length of A. With each iteration, the values are right-shifted one bit. After the last bit gets shifted off and the value becomes 0, then there is no additional value in continuing to rotate. Hence, the most precision will be attained by choosing `niter` to be one less than the word length.

For floating-point, the number of iterations is bounded by the size of the mantissa. In double, 52 iterations is the most you can do to continue adding to something with the same exponent. In single, it is 23. See the reference page for `eps` for more information about floating-point accuracy.

Thus, we can make our code more usable by not requiring the number of iterations to be input, and assuming that we want the most precision possible by changing `cordicqr` to use this default for `niter`.

```
function [Q,R] = cordicqr(A,varargin)
    if nargin>=2 && ~isempty(varargin{1})
        niter = varargin{1};
    elseif isa(A,'double') || isfi(A) && isdouble(A)
        niter = 52;
    elseif isa(A,'single') || isfi(A) && issingle(A)
        niter = single(23);
    elseif isfi(A)
        niter = int32(get(A,'WordLength') - 1);
    else
        assert(0,'First input must be double, single, or fi.');
```

A disadvantage of doing this is that this makes a section of our code dependent on data type. However, an advantage is that the function is much more convenient to use because you don't have to specify `niter` if you don't want to, and the main algorithm is still data-type independent. Similar to picking an optimal output type for Q, you can do this kind of input parsing in the beginning of a function and leave the main algorithm data-type independent.

Here is an example from a previous section, without needing to specify an optimal `niter`.

```
A = [7   -7   7   7
      7   7  -7   7
      7  -7  -7  -7
      7   7   7  -7];
```

```
[Q,R] = cordicqr(A)
```

Q =

```
0.5000   -0.5000   0.5000   0.5000
0.5000    0.5000  -0.5000   0.5000
0.5000   -0.5000  -0.5000  -0.5000
0.5000    0.5000   0.5000  -0.5000
```

R =

```
14.0000    0.0000  -0.0000  -0.0000
      0   14.0000  -0.0000   0.0000
      0      0   14.0000   0.0000
```

```
0      0      0  14.0000
```

Example: QR Factorization Not Unique

When you compare the results from `cordicqr` and the `qr` function in MATLAB, you will notice that the QR factorization is not unique. It is only important that Q is orthogonal, R is upper triangular, and $Q^*R - A$ is small.

Here is a simple example that shows the difference.

```
m = 3;
A = ones(m)
```

A =

```
1      1      1
1      1      1
1      1      1
```

The built-in QR function in MATLAB uses a different algorithm and produces:

```
[Q0,R0] = qr(A)
```

Q0 =

```
-0.5774  -0.5774  -0.5774
-0.5774   0.7887  -0.2113
-0.5774  -0.2113   0.7887
```

R0 =

```
-1.7321  -1.7321  -1.7321
      0      0      0
      0      0      0
```

And the `cordicqr` function produces:

```
[Q,R] = cordicqr(A)
```

Q =

```
0.5774   0.7495   0.3240
0.5774  -0.6553   0.4871
0.5774  -0.0942  -0.8110
```

R =

```
1.7321   1.7321   1.7321
      0   0.0000   0.0000
      0      0  -0.0000
```

Notice that the elements of Q from function `cordicqr` are different from Q_0 from built-in QR. However, both results satisfy the requirement that Q is orthogonal:

$Q_0^*Q_0'$

ans =

```

1.0000    0.0000    0
0.0000    1.0000    0
0         0         1.0000

```

Q^*Q'

ans =

```

1.0000    0.0000    0.0000
0.0000    1.0000   -0.0000
0.0000   -0.0000    1.0000

```

And they both satisfy the requirement that $Q^*R - A$ is small:

$Q_0^*R_0 - A$

ans =

```

1.0e-15 *
-0.1110   -0.1110   -0.1110
-0.1110   -0.1110   -0.1110
-0.1110   -0.1110   -0.1110

```

$Q^*R - A$

ans =

```

1.0e-15 *
-0.2220    0.2220    0.2220
0.4441         0         0
0.2220    0.2220    0.2220

```

Solving Systems of Equations Without Forming Q

Given matrices A and B , you can use the QR factorization to solve for X in the following equation:

$$A^*X = B.$$

If A has more rows than columns, then X will be the least-squares solution. If X and B have more than one column, then several solutions can be computed at the same time. If $A = Q^*R$ is the QR factorization of A , then the solution can be computed by back-solving

$$R^*X = C$$

where $C = Q' * B$. Instead of forming Q and multiplying to get $C = Q' * B$, it is more efficient to compute C directly. You can compute C directly by applying the rotations to the rows of B instead of to the columns of an identity matrix. The new algorithm is formed by the small modification of initializing $C = B$, and operating along the rows of C instead of the columns of Q .

```
function [R,C] = cordicrc(A,B,niter)
    Kn = inverse_cordic_growth_constant(niter);
    [m,n] = size(A);
    R = A;
    C = B;
    for j=1:n
        for i=j+1:m
            [R(j,j:end),R(i,j:end),C(j,:),C(i,:)] = ...
                cordicgivens(R(j,j:end),R(i,j:end),C(j,:),C(i,:),niter,Kn);
        end
    end
end
```

You can verify the algorithm with this example. Let A be a random 3-by-3 matrix, and B be a random 3-by-2 matrix.

```
A = [-0.8201    0.3573   -0.0100
     -0.7766   -0.0096   -0.7048
     -0.7274   -0.6206   -0.8901];
```

```
B = [-0.9286    0.3575
      0.6983    0.5155
      0.8680    0.4863];
```

Compute the QR factorization of A .

```
[Q,R] = cordicqr(A)
```

$Q =$

```
-0.6105    0.6133    0.5012
-0.5781    0.0876   -0.8113
-0.5415   -0.7850    0.3011
```

$R =$

```
1.3434    0.1235    0.8955
  0        0.7054    0.6309
  0         0         0.2988
```

Compute $C = Q' * B$ directly.

```
[R,C] = cordicrc(A,B)
```

$R =$

```
1.3434    0.1235    0.8955
  0        0.7054    0.6309
  0         0         0.2988
```

```
C =  
-0.3068 -0.7795  
-1.1897 -0.1173  
-0.7706 -0.0926
```

Subtract, and you will see that the error difference is on the order of roundoff.

```
Q'*B - C
```

```
ans =  
  
1.0e-15 *  
-0.0555 0.3331  
 0 0  
0.1110 0.2914
```

Now try the example in fixed-point. Declare A and B to be fixed-point types.

```
A = sfi(A)
```

```
A =  
-0.8201 0.3573 -0.0100  
-0.7766 -0.0096 -0.7048  
-0.7274 -0.6206 -0.8901  
  
DataTypeMode: Fixed-point: binary point scaling  
Signedness: Signed  
WordLength: 16  
FractionLength: 15
```

```
B = sfi(B)
```

```
B =  
-0.9286 0.3575  
0.6983 0.5155  
0.8680 0.4863  
  
DataTypeMode: Fixed-point: binary point scaling  
Signedness: Signed  
WordLength: 16  
FractionLength: 15
```

The necessary growth is 1.6468 times the square-root of the number of rows of A.

```
bit_growth = ceil(log2(cordic_growth_constant*sqrt(m)))
```

```
bit_growth =
```

2

Initialize R with the same values as A, and allow for bit growth.

```
R = sfi(A, get(A, 'WordLength')+bit_growth, get(A, 'FractionLength'))
```

R =

```
-0.8201    0.3573   -0.0100
-0.7766   -0.0096   -0.7048
-0.7274   -0.6206   -0.8901
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 18
      FractionLength: 15
```

The growth in C is the same as R, so initialize C and allow for bit growth the same way.

```
C = sfi(B, get(B, 'WordLength')+bit_growth, get(B, 'FractionLength'))
```

C =

```
-0.9286    0.3575
 0.6983    0.5155
 0.8680    0.4863
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 18
      FractionLength: 15
```

Compute $C = Q^*B$ directly, overwriting R and C.

```
[R,C] = cordicrc(R,C)
```

R =

```
1.3435    0.1233    0.8954
 0         0.7055    0.6308
 0         0         0.2988
```

```
      DataTypeMode: Fixed-point: binary point scaling
      Signedness: Signed
      WordLength: 18
      FractionLength: 15
```

C =

```
-0.3068   -0.7796
-1.1898   -0.1175
-0.7706   -0.0926
```

```
      DataTypeMode: Fixed-point: binary point scaling
```

```
Signedness: Signed
WordLength: 18
FractionLength: 15
```

An interesting use of this algorithm is that if you initialize B to be the identity matrix, then output argument C is Q'. You may want to use this feature to have more control over the data type of Q. For example,

```
A = [-0.8201    0.3573   -0.0100
      -0.7766   -0.0096   -0.7048
      -0.7274   -0.6206   -0.8901];
B = eye(size(A,1))
```

```
B =
```

```
    1    0    0
    0    1    0
    0    0    1
```

```
[R,C] = cordicrc(A,B)
```

```
R =
```

```
    1.3434    0.1235    0.8955
         0    0.7054    0.6309
         0         0    0.2988
```

```
C =
```

```
   -0.6105   -0.5781   -0.5415
    0.6133    0.0876   -0.7850
    0.5012   -0.8113    0.3011
```

Then C is orthogonal

```
C'*C
```

```
ans =
```

```
    1.0000    0.0000    0.0000
    0.0000    1.0000   -0.0000
    0.0000   -0.0000    1.0000
```

```
and R = C*A
```

```
R - C*A
```

```
ans =
```

```
    1.0e-15 *
```



```

0.6661   -0.0139   -0.1110
0.5551   -0.2220    0.6661
-0.2220   -0.1110    0.2776

```

Links to the Documentation

Fixed-Point Designer™

- `bitsra` Bit shift right arithmetic
- `fi` Construct fixed-point numeric object
- `fimath` Construct `fimath` object
- `fipref` Construct `fipref` object
- `get` Property values of object
- `globalfimath` Configure global `fimath` and return handle object
- `isfi` Determine whether variable is `fi` object
- `sfi` Construct signed fixed-point numeric object
- `upperbound` Upper bound of range of `fi` object
- `fiaccel` Accelerate fixed-point code

MATLAB

- `bitshift` Shift bits specified number of places
- `ceil` Round toward positive infinity
- `double` Convert to double precision floating point
- `eps` Floating-point relative accuracy
- `eye` Identity matrix
- `log2` Base 2 logarithm and dissect floating-point numbers into exponent and mantissa
- `prod` Product of array elements
- `qr` Orthogonal-triangular factorization
- `repmat` Replicate and tile array
- `single` Convert to single precision floating point
- `size` Array dimensions
- `sqrt` Square root
- `subsasgn` Subscripted assignment

Functions Used in this Example

These are the MATLAB functions used in this example.

CORDICQR computes the QR factorization using CORDIC.

- `[Q,R] = cordicqr(A)` chooses the number of CORDIC iterations based on the type of `A`.
- `[Q,R] = cordicqr(A,niter)` uses `niter` number of CORDIC iterations.

CORDICRC computes `R` from the QR factorization of `A`, and also returns `C = Q'*B` without computing `Q`.

- `[R,C] = cordicrc(A,B)` chooses the number of CORDIC iterations based on the type of A.
- `[R,C] = cordicrc(A,B,niter)` uses `niter` number of CORDIC iterations.

CORDIC_GROWTH_CONSTANT returns the CORDIC growth constant.

- `cordic_growth = cordic_growth_constant(niter)` returns the CORDIC growth constant as a function of the number of CORDIC iterations, `niter`.

GIVENSQR computes the QR factorization using standard Givens rotations.

- `[Q,R] = givensqr(A)`, where A is M-by-N, produces an M-by-N upper triangular matrix R and an M-by-M orthogonal matrix Q so that $A = Q \cdot R$.

CORDICQR_MAKEPLOTS makes the plots in this example by executing the following from the MATLAB command line.

```
load A_3_by_3_for_cordicqr_demo.mat
niter=32;
[Q,R] = cordicqr_makeplots(A,niter)
```

References

- 1 Ray Andraka, "A survey of CORDIC algorithms for FPGA based computers," 1998, ACM 0-89791-978-5/98/01.
- 2 Anthony J Cox and Nicholas J Higham, "Stability of Householder QR factorization for weighted least squares problems," in Numerical Analysis, 1997, Proceedings of the 17th Dundee Conference, Griffiths DF, Higham DJ, Watson GA (eds). Addison-Wesley, Longman: Harlow, Essex, U.K., 1998; 57-73.
- 3 Gene H. Golub and Charles F. Van Loan, *Matrix Computations*, 3rd ed, Johns Hopkins University Press, 1996, section 5.2.3 Givens QR Methods.
- 4 Daniel V. Rabinkin, William Song, M. Michael Vai, and Huy T. Nguyen, "Adaptive array beamforming with fixed-point arithmetic matrix inversion using Givens rotations," Proceedings of Society of Photo-Optical Instrumentation Engineers (SPIE) -- Volume 4474 Advanced Signal Processing Algorithms, Architectures, and Implementations XI, Franklin T. Luk, Editor, November 2001, pp. 294--305.
- 5 Jack E. Volder, "The CORDIC Trigonometric Computing Technique," Institute of Radio Engineers (IRE) Transactions on Electronic Computers, September, 1959, pp. 330-334.
- 6 Musheng Wei and Qiaohua Liu, "On growth factors of the modified Gram-Schmidt algorithm," Numerical Linear Algebra with Applications, Vol. 15, issue 7, September 2008, pp. 621-636.

Cleanup

```
fipref(originalFipref);
globalfimath(originalGlobalFimath);
close all
set(0, 'format', originalFormat);
%#ok<*MNEFF,*NASGU,*NOPTS,*ASGLU>
```

Compute Square Root Using CORDIC

This example shows how to compute square root using a CORDIC kernel algorithm in MATLAB®. CORDIC-based algorithms are critical to many embedded applications, including motor controls, navigation, signal processing, and wireless communications.

CORDIC is an acronym for COordinate Rotation DIGital Computer. The Givens rotation-based CORDIC algorithm (see [1,2]) is one of the most hardware efficient algorithms because it only requires iterative shift-add operations. The CORDIC algorithm eliminates the need for explicit multipliers, and is suitable for calculating a variety of functions, such as sine, cosine, arcsine, arccosine, arctangent, vector magnitude, divide, square root, hyperbolic and logarithmic functions.

The fixed-point CORDIC algorithm requires the following operations:

- 1 table lookup **per iteration**
- 2 shifts **per iteration**
- 3 additions **per iteration**

Note that for hyperbolic CORDIC-based algorithms, such as square root, certain iterations ($i = 4, 13, 40, 121, \dots, k, 3k+1, \dots$) are repeated to achieve result convergence.

CORDIC Kernel Algorithms Using Hyperbolic Computation Modes

You can use a CORDIC computing mode algorithm to calculate hyperbolic functions, such as hyperbolic trigonometric, square root, log, exp, etc.

CORDIC Equations in Hyperbolic Vectoring Mode

The hyperbolic vectoring mode is used for computing square root.

For the vectoring mode, the CORDIC equations are as follows:

$$x_{i+1} = x_i + y_i * d_i * 2^{-i}$$

$$y_{i+1} = y_i + x_i * d_i * 2^{-i}$$

$$z_{i+1} = z_i - d_i * \operatorname{atanh}(2^{-i})$$

where

$d_i = +1$ if $y_i < 0$, and -1 otherwise.

This mode provides the following result as N approaches $+\infty$:

- $x_N \approx A_N \sqrt{x_0^2 - y_0^2}$
- $y_N \approx 0$
- $z_N \approx z_0 + \operatorname{atanh}(y_0/x_0)$

where

$$A_N = \prod_{i=1}^N \sqrt{1 - 2^{-2i}}$$

Typically N is chosen to be a large-enough constant value. Thus, A_N may be pre-computed.

Note also that for **square root** we will use only the x_N result.

Implement a CORDIC Hyperbolic Vectoring Algorithm in MATLAB

A MATLAB code implementation example of the CORDIC hyperbolic vectoring kernel algorithm follows (for the case of scalar x , y , and z). This same code can be used for both fixed-point and floating-point data types.

CORDIC Hyperbolic Vectoring Kernel

```
k = 4; % Used for the repeated (3*k + 1) iteration steps
for idx = 1:n
    xtmp = bitsra(x, idx); % multiply by 2^(-idx)
    ytmp = bitsra(y, idx); % multiply by 2^(-idx)
    if y < 0
        x(:) = x + ytmp;
        y(:) = y + xtmp;
        z(:) = z - atanhLookupTable(idx);
    else
        x(:) = x - ytmp;
        y(:) = y - xtmp;
        z(:) = z + atanhLookupTable(idx);
    end
    if idx==k
        xtmp = bitsra(x, idx); % multiply by 2^(-idx)
        ytmp = bitsra(y, idx); % multiply by 2^(-idx)
        if y < 0
            x(:) = x + ytmp;
            y(:) = y + xtmp;
            z(:) = z - atanhLookupTable(idx);
        else
            x(:) = x - ytmp;
            y(:) = y - xtmp;
            z(:) = z + atanhLookupTable(idx);
        end
        k = 3*k + 1;
    end
end % idx loop
```

Compute Square Root Using the CORDIC Hyperbolic Vectoring Kernel

The judicious choice of initial values allows the CORDIC kernel hyperbolic vectoring mode algorithm to compute square root.

First, the following initialization steps are performed:

- x_0 is set to $v + 0.25$.
- y_0 is set to $v - 0.25$.

After N iterations, these initial values lead to the following output as N approaches $+\infty$:

$$x_N \approx A_N \sqrt{(v + 0.25)^2 - (v - 0.25)^2}$$

This may be further simplified as follows:

$$x_N \approx A_N \sqrt{v}$$

where A_N is the CORDIC gain as defined above.

Note: for square root, z and `atanhLookupTable` have no impact on the result. Hence, z and `atanhLookupTable` are not used.

MATLAB Implementation of a CORDIC Square Root Kernel

A MATLAB code implementation example of the CORDIC Square Root Kernel algorithm follows (for the case of scalar x and y). This same code can be used for both fixed-point and floating-point data types.

CORDIC Square Root Kernel

```
k = 4; % Used for the repeated (3*k + 1) iteration steps
```

```
for idx = 1:n
    xtmp = bitsra(x, idx); % multiply by 2^(-idx)
    ytmp = bitsra(y, idx); % multiply by 2^(-idx)
    if y < 0
        x(:) = x + ytmp;
        y(:) = y + xtmp;
    else
        x(:) = x - ytmp;
        y(:) = y - xtmp;
    end

    if idx==k
        xtmp = bitsra(x, idx); % multiply by 2^(-idx)
        ytmp = bitsra(y, idx); % multiply by 2^(-idx)
        if y < 0
            x(:) = x + ytmp;
            y(:) = y + xtmp;
        else
            x(:) = x - ytmp;
            y(:) = y - xtmp;
        end
        k = 3*k + 1;
    end
end % idx loop
```

This code is identical to the CORDIC hyperbolic vectoring kernel implementation above, except that z and `atanhLookupTable` are not used. This is a cost savings of 1 table lookup and 1 addition per iteration.

Example

Use the `cordicsqrt` function to compute the approximate square root of `v_fix` using ten CORDIC kernel iterations:

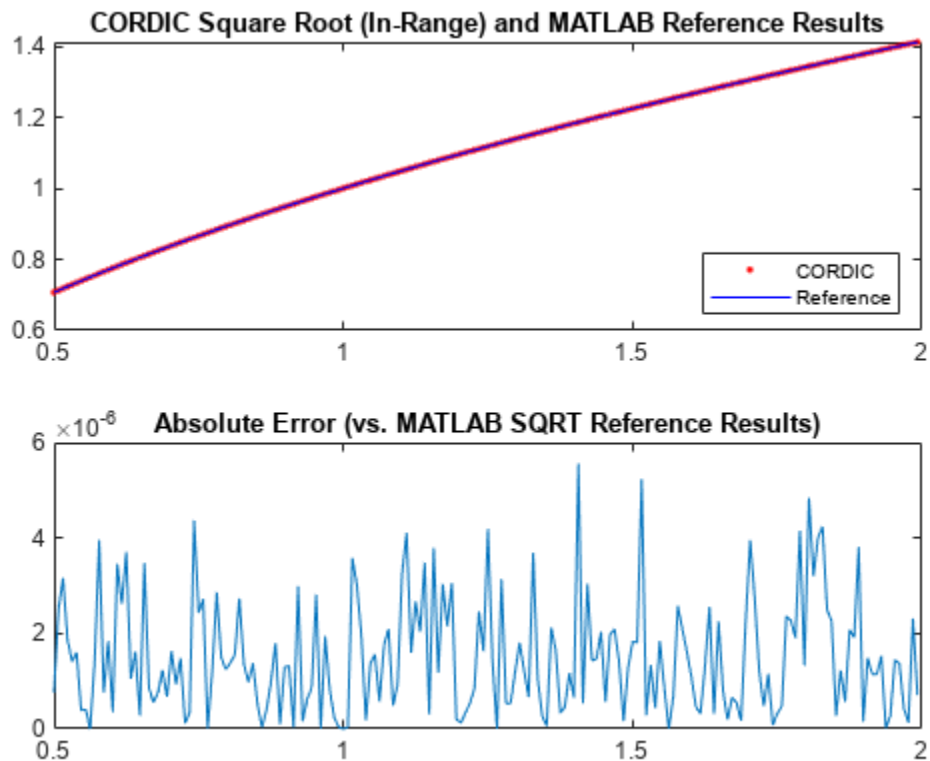
```
step = 2^-7;
v_fix = fi(0.5:step:(2-step), 1, 20); % fixed-point inputs in range [.5, 2)
niter = 10; % number of CORDIC iterations
x_sqr = cordicsqrt(v_fix, niter);
```

Get the real-world value (RWV) of the CORDIC outputs for comparison and plot the error between the MATLAB reference and CORDIC square root values.

```

x_cdc = double(x_sqr); % CORDIC results (scaled by An_hp)
v_ref = double(v_fix); % Reference floating-point input values
x_ref = sqrt(v_ref); % MATLAB reference floating-point results
figure;
subplot(211);
plot(v_ref, x_cdc, 'r.', v_ref, x_ref, 'b-');
legend('CORDIC', 'Reference', 'Location', 'SouthEast');
title('CORDIC Square Root (In-Range) and MATLAB Reference Results');
subplot(212);
absErr = abs(x_ref - x_cdc);
plot(v_ref, absErr);
title('Absolute Error (vs. MATLAB SQRT Reference Results)');

```



Overcoming Algorithm Input Range Limitations

Many square root algorithms normalize the input value, v , to within the range of $[0.5, 2)$. This pre-processing is typically done using a fixed word length normalization, and can be used to support small as well as large input value ranges.

The CORDIC-based square root algorithm implementation is particularly sensitive to inputs outside of this range. The `cordicsqrt` function overcomes this algorithm range limitation through a normalization approach based on the following mathematical relationships:

$$v = u * 2^n, \text{ for some } 0.5 < u < 2 \text{ and some even integer } n.$$

Thus:

$$\sqrt{v} = \sqrt{u} * 2^{n/2}$$

In the `cordicsqrt` function, the values for u and n , described above, are found during normalization of the input v . n is the number of leading zero most significant bits (MSBs) in the binary representation of the input v . These values are found through a series of bitwise logic and shifts. Note that because n must be even, if the number of leading zero MSBs is odd, one additional bit shift is made to make n even. The resulting value after these shifts is the value $0.5 < u < 2$.

u becomes the input to the CORDIC-based square root kernel, where an approximation to \sqrt{u} is calculated. The result is then scaled by $2^{n/2}$ so that it is back in the correct output range. This is achieved through a simple bit shift by $n/2$ bits. The (left or right) shift direction depends on the sign of n .

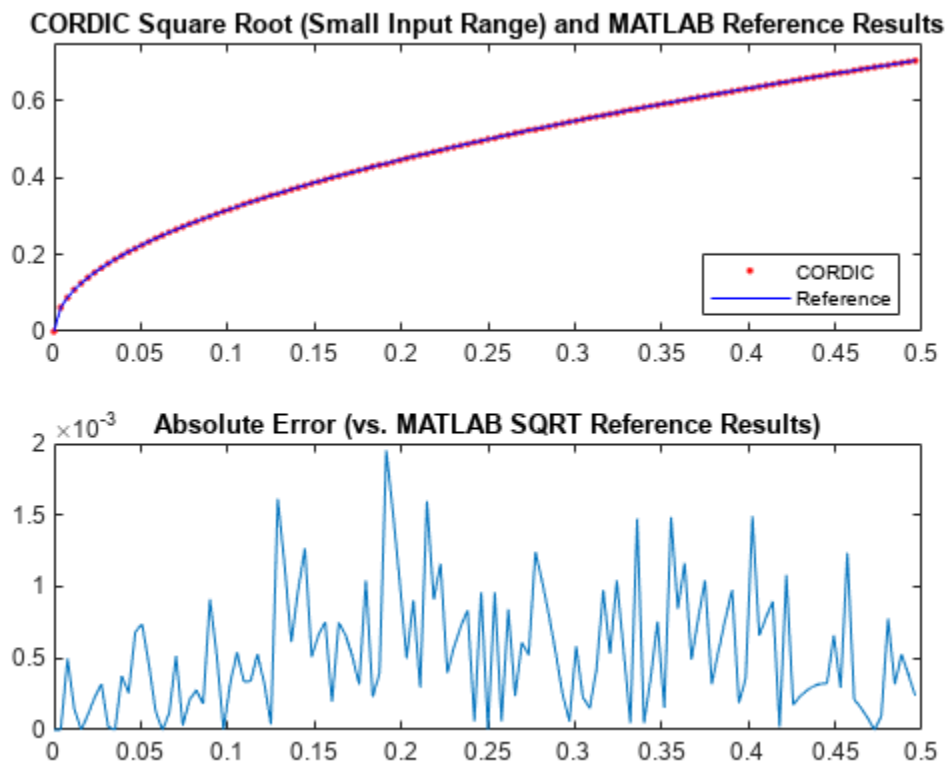
Example

Compute the square root of 10-bit fixed-point input data with a small non-negative range using CORDIC. Compare the CORDIC-based algorithm results to the floating-point MATLAB reference results over the same input range.

```
step      = 2^-8;
u_ref     = 0:step:(0.5-step); % Input array (small range of values)
u_in_arb  = fi(u_ref,0,10); % 10-bit unsigned fixed-point input data values
u_len     = numel(u_ref);
sqrt_ref  = sqrt(double(u_in_arb)); % MATLAB sqrt reference results
niter     = 10;
results   = zeros(u_len, 2);
results(:,2) = sqrt_ref(:);
```

Compute the equivalent real-world value (RWV) result for plotting. Plot the RWV of CORDIC and MATLAB reference results.

```
x_out = cordicsqrt(u_in_arb, niter);
results(:,1) = double(x_out);
figure;
subplot(211);
plot(u_ref, results(:,1), 'r.', u_ref, results(:,2), 'b-');
legend('CORDIC', 'Reference', 'Location', 'SouthEast');
title('CORDIC Square Root (Small Input Range) and MATLAB Reference Results');
axis([0 0.5 0 0.75]);
subplot(212);
absErr = abs(results(:,2) - results(:,1));
plot(u_ref, absErr);
title('Absolute Error (vs. MATLAB SQRT Reference Results)');
```



Example

Compute the square root of 16-bit fixed-point input data with a large positive range using CORDIC. Compare the CORDIC-based algorithm results to the floating-point MATLAB reference results over the same input range.

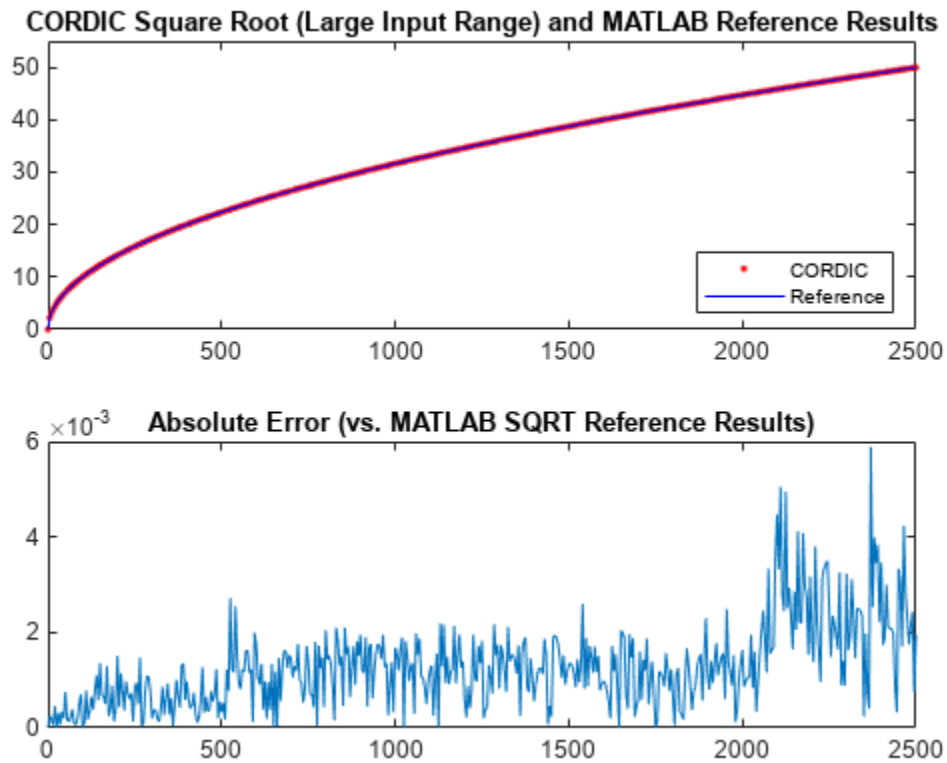
```
u_ref = 0:5:2500; % Input array (larger range of values)
u_in_arb = fi(u_ref,0,16); % 16-bit unsigned fixed-point input data values
u_len = numel(u_ref);
sqrt_ref = sqrt(double(u_in_arb)); % MATLAB sqrt reference results
niter = 16;
results = zeros(u_len, 2);
results(:,2) = sqrt_ref(:);
```

Compute the equivalent real-world value (RWV) result for plotting. Plot the RWV of CORDIC and MATLAB reference results.

```
x_out = cordicsqrt(u_in_arb, niter);
results(:,1) = double(x_out);
figure;
subplot(211);
plot(u_ref, results(:,1), 'r.', u_ref, results(:,2), 'b-');
legend('CORDIC', 'Reference', 'Location', 'SouthEast');
title('CORDIC Square Root (Large Input Range) and MATLAB Reference Results');
axis([0 2500 0 55]);
subplot(212);
absErr = abs(results(:,2) - results(:,1));
```



```
plot(u_ref, absErr);  
title('Absolute Error (vs. MATLAB SQRT Reference Results)');
```



References

- 1 Jack E. Volder, "The CORDIC Trigonometric Computing Technique," IRE Transactions on Electronic Computers, Volume EC-8, September 1959, pp. 330-334.
- 2 J.S. Walther, "A Unified Algorithm for Elementary Functions," Conference Proceedings, Spring Joint Computer Conference, May 1971, pp. 379-385.

Normalize Data for Lookup Tables

This example shows how to normalize data for use in lookup tables.

Lookup tables are a very efficient way to write computationally-intense functions for fixed-point embedded devices. For example, you can efficiently implement logarithm, sine, cosine, tangent, and square-root using lookup tables. You normalize the inputs to these functions to produce a smaller lookup table, and then you scale the outputs by the normalization factor. This example shows how to implement the normalization function that is used in examples Implement Fixed-Point Square Root Using Lookup Table and Implement Fixed-Point Log2 Using Lookup Table.

Setup

To assure that this example does not change your preferences or settings, this code stores the original state.

```
originalFormat = get(0, 'format'); format long g
originalWarningState = warning('off', 'fixed:fi:underflow');
originalFiprefState = get(fipref); reset(fipref)
```

You will restore this state at the end of the example.

Example

Use the normalization function `fi_normalize_unsigned_8_bit_byte`, defined below, to normalize the data in `u`.

```
u = fi(0.3, 1, 16, 8);
```

In binary, $u = 00000000.01001101_2 = 0.30078125$ (the fixed-point value is not exactly 0.3 because of roundoff to 8 bits). The goal is to normalize such that

$$u = 1.00110100000000_2 * 2^{(-2)} = x * 2^n.$$

Start with `u` represented as an unsigned integer.

High byte Low byte

00000000 01001101 Start: `u` as unsigned integer.

The high byte is `0 = 00000000_2`. Add 1 to make an index out of it: `index = 0 + 1 = 1`. The number-of-leading-zeros lookup table at index 1 indicates that there are 8 leading zeros: `NLZLUT(1) = 8`. Left shift by this many bits.

High byte Low byte

01001101 00000000 Left-shifted by 8 bits.

Iterate once more to remove the leading zeros from the next byte.

The high byte is `77 = 01001101_2`. Add 1 to make an index out of it: `index = 77 + 1 = 78`. The number-of-leading-zeros lookup table at index 78 indicates that there is 1 leading zero: `NLZLUT(78) = 1`. Left shift by this many bits.

High byte Low byte

100110100 0000000 Left-shifted by 1 additional bit, for a total of 9.

Reinterpret these bits as unsigned fixed-point with 15 fractional bits.

$x = 1.001101000000000_2 = 1.203125$

The value for n is the word-length of u , minus the fraction length of u , minus the number of left shifts, minus 1.

$n = 16 - 8 - 9 - 1 = -2$

And so your result is:

```
[x,n] = fi_normalize_unsigned_8_bit_byte(u)
```

```
x =
    1.203125
    DataTypeMode: Fixed-point: binary point scaling
    Signedness: Unsigned
    WordLength: 16
    FractionLength: 15
```

```
n = int8
    -2
```

Comparing binary values, you can see that x has the same bits as u , left-shifted by 9 bits.

```
binary_representation_of_u = bin(u)
```

```
binary_representation_of_u =
'0000000001001101'
```

```
binary_representation_of_x = bin(x)
```

```
binary_representation_of_x =
'1001101000000000'
```

Cleanup

Restore original state.

```
set(0, 'format', originalFormat);
warning(originalWarningState);
fipref(originalFiprefState);
%#ok< *NOPTS>
```

Function to Normalize Unsigned Data

This algorithm normalizes unsigned data with 8-bit bytes. Given input $u > 0$, the output x is normalized such that

$u = x \cdot 2^n$

where $1 \leq x < 2$ and n is an integer. Note that n may be positive, negative, or zero.

Function `fi_normalize_unsigned_8_bit_byte` looks at the 8 most-significant-bits of the input at a time, and left shifts the bits until the most-significant bit is a 1. The number of bits to shift for each 8-bit byte is read from the number-of-leading-zeros lookup table, `NLZLUT`.

```
function [x,n] = fi_normalize_unsigned_8_bit_byte(u)
    assert(isscalar(u),'Input must be scalar');
    assert(all(u>0),'Input must be positive. ');
    assert(isfi(u) && isfixed(u),'Input must be a fi object with fixed-point data type. ');
    u = removefimath(u);
    NLZLUT = number_of_leading_zeros_look_up_table();
    word_length = u.WordLength;
    u_fraction_length = u.FractionLength;
    B = 8;
    leftshifts=int8(0);
    % Reinterpret the input as an unsigned integer.
    T_unsigned_integer = numerictype(0, word_length, 0);
    v = reinterpretcast(u,T_unsigned_integer);
    F = fimath('OverflowAction','Wrap',...
              'RoundingMethod','Floor',...
              'SumMode','KeepLSB',...
              'SumWordLength',v.WordLength);
    v = setfimath(v,F);
    % Unroll the loop in generated code so there will be no branching.
    for k = coder.unroll(1:ceil(word_length/B))
        % For each iteration, see how many leading zeros are in the high
        % byte of V, and shift them out to the left. Continue with the
        % shifted V for as many bytes as it has.
        %
        % The index is the high byte of the input plus 1 to make it a
        % one-based index.
        index = int32(bitsra(v,word_length-B) + uint8(1));
        % Index into the number-of-leading-zeros lookup table. This lookup
        % table takes in a byte and returns the number of leading zeros in the
        % binary representation.
        shiftamount = NLZLUT(index);
        % Left-shift out all the leading zeros in the high byte.
        v = bitsll(v,shiftamount);
        % Update the total number of left-shifts
        leftshifts = leftshifts+shiftamount;
    end
    % The input has been left-shifted so the most-significant-bit is a 1.
    % Reinterpret the output as unsigned with one integer bit, so
    % that 1 <= x < 2.
    T_x = numerictype(0,word_length,word_length-1);
    x = reinterpretcast(v,T_x);
    x = removefimath(x);
    % Let Q = int(u). Then u = Q*2^(-u_fraction_length),
    % and x = Q*2^leftshifts * 2^(1-word_length). Therefore,
    % u = x*2^n, where n is defined as:
    n = word_length - u_fraction_length - leftshifts - 1;
end
```

Number-of-Leading-Zeros Lookup Table

Function `number_of_leading_zeros_look_up_table` is used by `fi_normalize_unsigned_8_bit_byte` and returns a table of the number of leading zero bits in an 8-bit word.

Implement Fixed-Point Log2 Using Lookup Table

This example shows how to implement fixed-point \log_2 using a lookup table. Lookup tables generate efficient code for embedded devices.

Setup

To ensure that this example does not change your preferences or settings, this code stores the original state.

```
originalFormat = get(0,'format'); format long g
originalWarningState = warning('off','fixed:fi:underflow');
originalFiprefState = get(fipref); reset(fipref)
```

You will restore this state at the end of the example.

Log2 Implementation

The \log_2 algorithm, implemented in the function `fi_log2lookup_8_bit_byte` below, is summarized here.

- 1 Declare the number of bits in a byte, B , as a constant. In this example, $B=8$.
- 2 Use function `fi_normalize_unsigned_8_bit_byte()` described in example [Normalize Data for Lookup Tables](#) to normalize the input $u>0$ such that $u = x \cdot 2^n$ and $1 \leq x < 2$.
- 3 Extract the upper B -bits of x . Let x_B denote the upper B -bits of x .
- 4 Generate lookup table, `LOG2LUT`, such that the integer $i = x_B - 2^{(B-1)} + 1$ is used as an index to `LOG2LUT` so that $\log_2(x_B)$ can be evaluated by looking up the index $\log_2(x_B) = \text{LOG2LUT}(i)$.
- 5 Use the remainder, $r = x - x_B$, interpreted as a fraction, to linearly interpolate between `LOG2LUT(i)` and the next value in the table `LOG2LUT(i+1)`. The remainder, r , is created by extracting the lower $w - B$ bits of x , where w denotes the word length of x . It is interpreted as a fraction by using function `reinterpretcast()`.
- 6 Finally, compute the output using the lookup table and linear interpolation:

$$\begin{aligned} \log_2(u) &= \log_2(x \cdot 2^n) \\ &= n + \log_2(x) \\ &= n + \text{LOG2LUT}(i) + r \cdot (\text{LOG2LUT}(i+1) - \text{LOG2LUT}(i)) \end{aligned}$$

Example

Use `fi_log2lookup_8_bit_byte()` to compute the fixed-point \log_2 using a lookup table. Compare the fixed-point lookup table result to the logarithm calculated using `log2` and double precision.

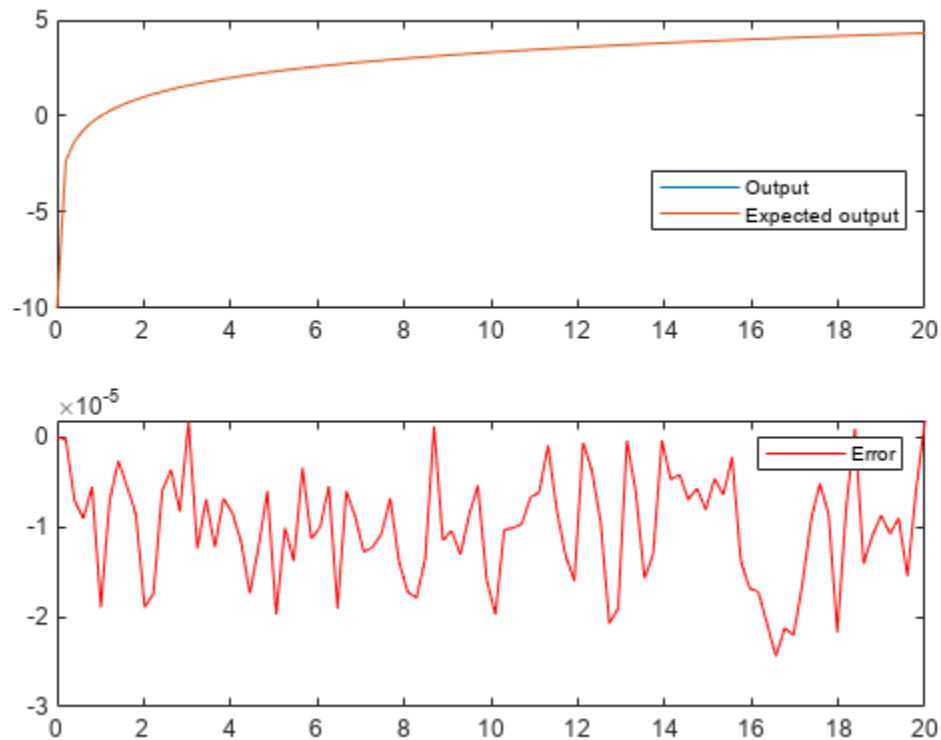
```
u = fi(linspace(0.001,20,100));
y = fi_log2lookup_8_bit_byte(u);
y_expected = log2(double(u));
```

Plot the results.

```
clf
subplot(211)
```

```
plot(u,y,u,y_expected)
legend('Output','Expected output','Location','Best')
```

```
subplot(212)
plot(u,double(y)-y_expected,'r')
legend('Error')
```



```
figure(gcf)
```

Cleanup

Restore the original state.

```
set(0,'format',originalFormat);
warning(originalWarningState);
fipref(originalFiprefState);
```

fi_log2lookup_8_bit_byte Function Definition

```
function y = fi_log2lookup_8_bit_byte(u)
    % Load the lookup table
    LOG2LUT = log2_lookup_table();
    % Remove fimath from the input to insulate this function from math
    % settings declared outside this function.
    u = removefimath(u);
    % Declare the output
    y = coder.nullcopy(fi(zeros(size(u)),numerictype(LOG2LUT),fimath(LOG2LUT)));
    B = 8; % Number of bits in a byte
```

```

w = u.WordLength;
for k = 1:numel(u)
    assert(u(k)>0,'Input must be positive. ');
    % Normalize the input such that u = x*2^n and 1 <= x < 2
    [x,n] = fi_normalize_unsigned_8_bit_byte(u(k));
    % Extract the high byte of x
    high_byte = storedInteger(bitsliceget(x, w, w - B + 1));
    % Convert the high byte into an index for LOG2LUT
    i = high_byte - 2^(B-1) + 1;
    % Interpolate between points.
    % The upper byte was used for the index into LOG2LUT
    % The remaining bits make up the fraction between points.
    T_unsigned_fraction = numericType(0, w-B, w-B);
    r = reinterpretcast(bitsliceget(x,w-B,1), T_unsigned_fraction);
    y(k) = n + LOG2LUT(i) + ...
        r*(LOG2LUT(i+1) - LOG2LUT(i)) ;
end
% Remove fimath from the output to insulate the caller from math settings
% declared inside this function.
y = removefimath(y);
end

```

Log2 Lookup Table

The function `log2_lookup_table` loads the lookup table of log2 values. You can create the table by running:

```
B = 8;
```

```
log2_table = log2((2^(B-1):2^(B))/2^(B-1))
```

```

function LOG2LUT = log2_lookup_table()
    B = 8; % Number of bits in a byte
    % log2_table = log2((2^(B-1) : 2^(B)) / 2^(B - 1))
    log2_table = [0.0000000000000000    0.011227255423254    0.022367813028454    0.033423001537450
0.044394119358453    0.055282435501190    0.066089190457773    0.076815597050831
0.087462841250339    0.098032082960527    0.108524456778169    0.118941072723507
0.129283016944966    0.139551352398794    0.149747119504682    0.159871336778389
0.169925001442312    0.179909090014934    0.189824558880017    0.199672344836364
0.209453365628950    0.219168520462162    0.228818690495881    0.238404739325079
0.247927513443586    0.257387842692652    0.266786540694901    0.276124405274238
0.285402218862248    0.294620748891627    0.303780748177103    0.312882955284355
0.321928094887362    0.330916878114617    0.339850002884625    0.348728154231078
0.357552004618084    0.366322214245816    0.375039431346925    0.383704292474052
0.392317422778760    0.400879436282184    0.409390936137702    0.417852514885898
0.426264754702098    0.434628227636725    0.442943495848728    0.451211111832329
0.459431618637297    0.467605550082997    0.475733430966398    0.483815777264256
0.491853096329675    0.499845887083205    0.507794640198696    0.515699838284042
0.523561956057013    0.531381460516312    0.539158811108031    0.546894459887637
0.554588851677637    0.562242424221073    0.569855608330948    0.577428828035749
0.584962500721156    0.592457037268080    0.599912842187128    0.607330313749611
0.614709844115208    0.622051819456376    0.629356620079610    0.636624620543649
0.643856189774725    0.651051691178929    0.658211482751795    0.665335917185176
0.672425341971496    0.679480099505446    0.686500527183218    0.693486957499325
0.700439718141092    0.707359132080883    0.714245517666123    0.721099188707185
0.727920454563199    0.734709620225838    0.741466986401147    0.748192849589460
0.754887502163469    0.761551232444479    0.768184324776926    0.774787059601173
0.781359713524660    0.787902559391432    0.794415866350106    0.800899899920305

```



```

0.807354922057604    0.813781191217037    0.820178962415188    0.826548487290915
0.832890014164742    0.839203788096944    0.845490050944375    0.851749041416058
0.857980995127572    0.864186144654280    0.870364719583405    0.876516946565000
0.882643049361841    0.888743248898259    0.894817763307943    0.900866807980749
0.906890595608518    0.912889336229962    0.918863237274595    0.924812503605781
0.930737337562886    0.936637939002571    0.942514505339240    0.948367231584678
0.954196310386875    0.960001932068081    0.965784284662087    0.971543553950772
0.977279923499916    0.982993574694310    0.988684686772166    0.994353436858858
1.000000000000000];
% Cast to fixed point with the most accurate rounding method
WL = 4*B; % Word length
FL = 2*B; % Fraction length
LOG2LUT = fi(log2_table,1,WL,FL,'RoundingMethod','Nearest');
% Set fimath for the most efficient math operations
F = fimath('OverflowAction','Wrap',...
'RoundingMethod','Floor',...
'SumMode','SpecifyPrecision',...
'SumWordLength',WL,...
'SumFractionLength',FL,...
'ProductMode','SpecifyPrecision',...
'ProductWordLength',WL,...
'ProductFractionLength',2*FL);
LOG2LUT = setfimath(LOG2LUT,F);
end

```

Implement Fixed-Point Square Root Using Lookup Table

This example shows how to implement fixed-point square root using a lookup table. Lookup tables generate efficient code for embedded devices.

Setup

To ensure that this example does not change your preferences or settings, this code stores the original state.

```
originalFormat = get(0,'format'); format long g
originalWarningState = warning('off','fixed:fi:underflow');
originalFiprefState = get(fipref); reset(fipref)
```

You will restore this state at the end of the example.

Square Root Implementation

The square root algorithm is summarized here.

- 1 Declare the number of bits in a byte, B , as a constant. In this example, $B = 8$.
- 2 Use the function `fi_normalize_unsigned_8_bit_byte()`, described in the example “Normalize Data for Lookup Tables” on page 54-96, to normalize the input $u > 0$ such that $u = x \cdot 2^n$, $0.5 \leq x < 2$, and n is even.
- 3 Extract the upper B bits of x . Let x_B denote the upper B bits of x .
- 4 Generate a lookup table, `SQRTLUT`, such that the integer $i = x_B - 2^{(B-2)} + 1$ is used as an index to `SQRTLUT` so that $\sqrt{x_B}$ can be evaluated by looking up the index $\sqrt{x_B} = \text{SQRTLUT}(i)$.
- 5 Use the remainder, $r = x - x_B$, interpreted as a fraction, to linearly interpolate between `SQRTLUT(i)` and the next value in the table `SQRTLUT(i+1)`. The remainder, r , is created by extracting the lower $w - B$ bits of x , where w denotes the wordlength of x . It is interpreted as a fraction by using function `reinterpcast()`.
- 6 Finally, compute the output using the lookup table and linear interpolation:

$$\begin{aligned} \sqrt{u} &= \sqrt{x \cdot 2^n} \\ &= \sqrt{x} \cdot 2^{(n/2)} \\ &= (\text{SQRTLUT}(i) + r \cdot (\text{SQRTLUT}(i+1) - \text{SQRTLUT}(i))) \cdot 2^{(n/2)} \end{aligned}$$

The function `fi_sqrtlookup_8_bit_byte()`, defined at the end of this example, implements this algorithm.

Example

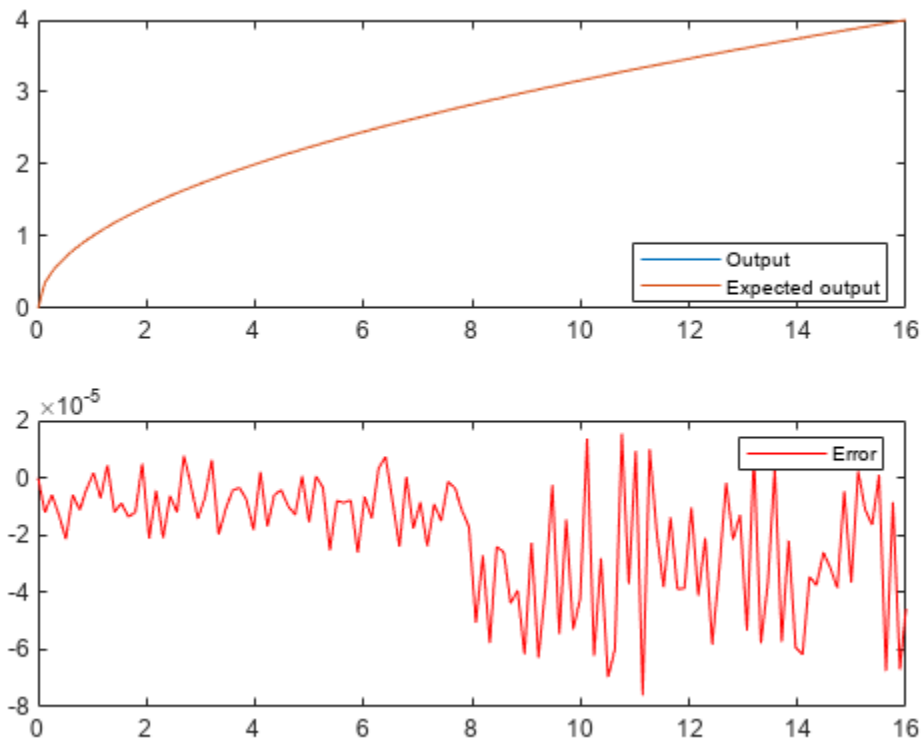
Use `fi_sqrtlookup_8_bit_byte()` to compute the fixed-point square root using a lookup table. Compare the fixed-point lookup table result to the square root calculated using `sqrt` and double precision.

```
u = fi(linspace(0,128,1000),0,16,12);
y = fi_sqrtlookup_8_bit_byte(u);
y_expected = sqrt(double(u));
```

Plot the results.

```
clf
subplot(211)
plot(u,y,u,y_expected)
legend('Output','Expected output','Location','Best')

subplot(212)
plot(u,double(y)-y_expected,'r')
legend('Error')
```



```
figure(gcf)
```

Cleanup

Restore the original state.

```
set(0,'format',originalFormat);
warning(originalWarningState);
fipref(originalFiprefState);
```

sqrt_lookup_table Function Definition

The function `sqrt_lookup_table` loads the lookup table of square-root values. You can create the table by running:

```
sqrt_table = sqrt((2^(B-2):2^(B))/2^(B-1));
```

```

function SQRTLUT = sqrt_lookup_table()
    B = 8; % Number of bits in a byte
    % sqrt_table = sqrt((2^(B-2):2^(B))/2^(B-1))
    sqrt_table = [0.707106781186548    0.712609640686961    0.718070330817254    0.723489806424389
0.728868986855663    0.734208757779421    0.739509972887452    0.744773455488312
0.750000000000000    0.755190373349661    0.760345316287277    0.765465544619743
0.770551750371122    0.775604602874429    0.780624749799800    0.785612818123533
0.790569415042095    0.795495128834866    0.800390529679106    0.805256170420320
0.810092587300983    0.814900300650331    0.819679815537750    0.824431622392057
0.829156197588850    0.833854004007896    0.838525491562421    0.843171097702003
0.847791247890659    0.852386356061616    0.856956825050130    0.861503047005639
0.866025403784439    0.870524267324007    0.875000000000000    0.879452954966893
0.883883476483184    0.888291900221993    0.892678553567856    0.897043755900458
0.901387818865997    0.905711046636840    0.910013736160065    0.914296177395487
0.918558653543692    0.922801441264588    0.927024810886958    0.931229026609459
0.935414346693485    0.939581023648307    0.943729304408844    0.947859430506444
0.951971638232989    0.956066158798647    0.960143218483576    0.964203038783845
0.968245836551854    0.972271824131503    0.976281209488332    0.980274196334883
0.984250984251476    0.988211768802619    0.992156741649222    0.996086090656827
1.000000000000000    1.003898650263063    1.007782218537319    1.011650878514915
1.015504800579495    1.019344151893756    1.023169096484056    1.026979795322186
1.030776406404415    1.034559084827928    1.038327982864759    1.042083250033317
1.045825033167594    1.049553476484167    1.053268721647045    1.056970907830485
1.060660171779821    1.064336647870400    1.068000468164691    1.071651762467640
1.075290658380328    1.078917281352004    1.082531754730548    1.086134199811423
1.089724735885168    1.093303480283494    1.096870548424015    1.100426053853688
1.103970108290981    1.107502821666834    1.111024302164449    1.114534656257938
1.118033988749895    1.121522402807898    1.125000000000000    1.128466880329237
1.131923142267177    1.135368882786559    1.138804197393037    1.142229180156067
1.145643923738960    1.149048519428140    1.152443057161611    1.155827625556683
1.159202311936963    1.162567202358642    1.165922381636102    1.169267933366857
1.172603939955857    1.175930482639174    1.179247641507075    1.182555495526531
1.185854122563142    1.189143599402528    1.192424001771182    1.195695404356812
1.198957880828180    1.202211503854459    1.205456345124119    1.208692475363357
1.211919964354082    1.215138880951474    1.218349293101120    1.221551267855754
1.224744871391589    1.227930169024281    1.231107225224513    1.234276103633219
1.237436867076458    1.240589577579950    1.243734296383275    1.246871083953750
1.250000000000000    1.253121103485214    1.256234452640111    1.259340104975618
1.262438117295260    1.265528545707287    1.268611445636527    1.271686871835988
1.274754878398196    1.277815518766305    1.280868845744950    1.283914911510884
1.286953767623375    1.289985465034393    1.293010054098575    1.296027584582983
1.299038105676658    1.302041665999979    1.305038313613819    1.308028096028522
1.311011060212689    1.313987252601790    1.316956719106592    1.319919505121430
1.322875655532295    1.325825214724777    1.328768226591831    1.331704734541407
1.334634781503914    1.337558409939543    1.340475661845451    1.343386578762792
1.346291201783626    1.349189571557681    1.352081728298996    1.354967711792425
1.357847561400027    1.360721316067327    1.363589014329464    1.366450694317215
1.369306393762915    1.372156150006259    1.375000000000000    1.377837980315538
1.380670127148408    1.383496476323666    1.386317063301177    1.389131923180804
1.391941090707505    1.394744600276337    1.397542485937369    1.400334781400505
1.403121520040228    1.405902734900249    1.408678458698081    1.411448723829527
1.414213562373095];
% Cast to fixed point with the most accurate rounding method
WL = 4*B; % Word length
FL = 2*B; % Fraction length
SQRTLUT = fi(sqrt_table,1,WL,FL,'RoundingMethod','Nearest');
% Set fimath for the most efficient math operations
F = fimath('OverflowAction','Wrap',...

```

```

        'RoundingMethod','Floor',...
        'SumMode','KeepLSB',...
        'SumWordLength',WL,...
        'ProductMode','KeepLSB',...
        'ProductWordLength',WL);
    SQRTLUT = setfimath(SQRTLUT,F);
end

fi_sqrtlookup_8_bit_byte() Function Definition

function y = fi_sqrtlookup_8_bit_byte(u)
    % Load the lookup table
    SQRTLUT = sqrt_lookup_table();
    % Remove fimath from the input to insulate this function from math
    % settings declared outside this function.
    u = removefimath(u);
    % Declare the output
    y = coder.nullcopy(fi(zeros(size(u)),numerictype(SQRTLUT),fimath(SQRTLUT)));
    B = 8; % Number of bits in a byte
    w = u.WordLength;
    for k = 1:numel(u)
        assert(u(k)>=0,'Input must be non-negative. ');
        if u(k)==0
            y(k)=0;
        else
            % Normalize the input such that  $u = x \cdot 2^n$  and  $0.5 \leq x < 2$ 
            [x,n] = fi_normalize_unsigned_8_bit_byte(u(k));
            isodd = storedInteger(bitand(fi(1,1,8,0),fi(n)));
            x = bitsra(x,isodd);
            n = n + isodd;
            % Extract the high byte of x
            high_byte = storedInteger(bitsliceget(x,w,w-B+1));
            % Convert the high byte into an index for SQRTLUT
            i = high_byte - 2^(B-2) + 1;
            % The upper byte was used for the index into SQRTLUT.
            % The remainder, r, interpreted as a fraction, is used to
            % linearly interpolate between points.
            T_unsigned_fraction = numerictype(0,w-B,w-B);
            r = reinterpretcast(bitsliceget(x,w-B,1),T_unsigned_fraction);
            y(k) = bitshift((SQRTLUT(i) + r*(SQRTLUT(i+1) - SQRTLUT(i))),...
                bitsra(n,1));
        end
    end
    % Remove fimath from the output to insulate the caller from math settings
    % declared inside this function.
    y = removefimath(y);
end

```

Convert Fast Fourier Transform (FFT) to Fixed Point

This example shows how to convert a textbook version of the Fast Fourier Transform (FFT) algorithm into fixed-point MATLAB® code.

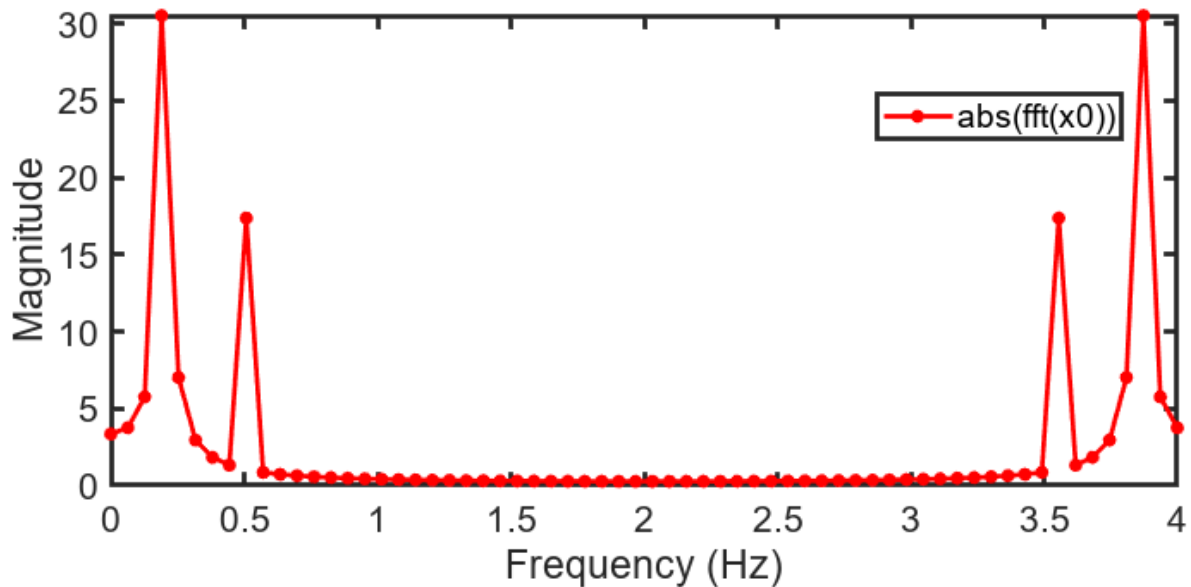
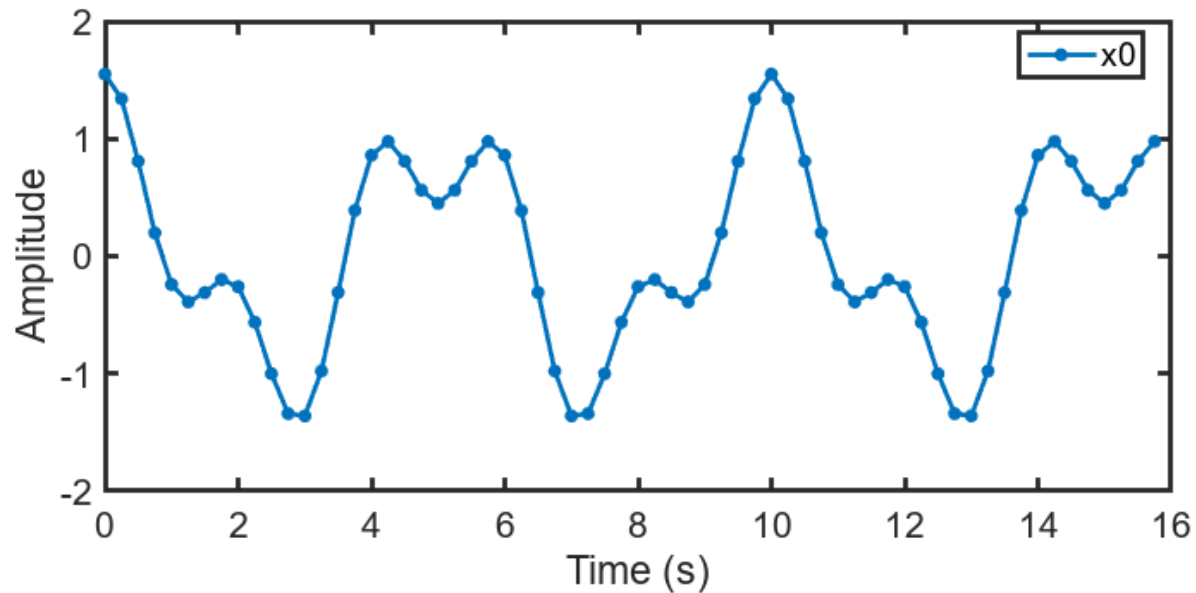
Run the following code to capture current states, and reset the global states.

```
FIPREF_STATE = get(fipref);
reset(fipref)
```

Textbook FFT Algorithm

FFT is a complex-valued linear transformation from the time domain to the frequency domain. For example, if you construct a vector as the sum of two sinusoids and transform it with the FFT, you can see the peaks of the frequencies in the FFT magnitude plot.

```
n = 64; % Number of points
Fs = 4; % Sampling frequency in Hz
t = (0:(n-1))/Fs; % Time vector
f = linspace(0,Fs,n); % Frequency vector
f0 = 0.2; f1 = 0.5; % Frequencies, in Hz
x0 = cos(2*pi*f0*t) + 0.55*cos(2*pi*f1*t); % Time-domain signal
x0 = complex(x0); % The textbook algorithm requires the input to be complex
y0 = fft(x0); % Frequency-domain transformation fft() is a MATLAB function
fi_fft_demo_ini_plot(t,x0,f,y0); % Plot the results from fft and time-domain signal
```



The peaks at 0.2 and 0.5 Hz in the frequency plot correspond to the two sinusoids of the time-domain signal at those frequencies.

Note the reflected peaks at 3.5 and 3.8 Hz. When the input to an FFT is real-valued, as it is in this case, then the output y is the conjugate-symmetric:

$$y(k) = \text{conj}(y(n - k))$$

There are many different implementations of the FFT, each having its own costs and benefits. You may find that a different algorithm is better for your application than the one given here. This algorithm provides you with an example of how you can begin your own exploration.

This example uses the decimation-in-time unit-stride FFT shown in Algorithm 1.6.2 on page 45 of the book *Computational Frameworks for the Fast Fourier Transform* by Charles Van Loan.

In pseudo-code, the algorithm in the textbook is as follows:

Algorithm 1.6.2. If x is a complex vector of length n and $n = 2^t$, then the following algorithm overwrites x with $F_n x$.

```

    x = P_n x
w = w_n^{(long)} (See Van Loan section 1.4.11.)
for q = 1:t
    L = 2^q; r = n/L; L* = L/2;
for k = 0:r - 1
for j = 0:L* - 1
    tau = w(L* - 1 + j) * x(kL + j + L*)
    x(kL + j + L*) = x(kL + j) - tau
    x(kL + j) = x(kL + j) + tau
end
end
end

```

The textbook algorithm uses zero-based indexing. F_n is an n -by- n Fourier-transform matrix, P_n is an n -by- n bit-reversal permutation matrix, and w is a complex vector of twiddle factors. The twiddle factors, w , are complex roots of unity computed by the following algorithm:

```

function w = fi_radix2twiddles(n)
%FI_RADIX2TWIDDLES Twiddle factors for radix-2 FFT example.
% W = FI_RADIX2TWIDDLES(N) computes the length N-1 vector W of
% twiddle factors to be used in the FI_M_RADIX2FFT example code.
%
% See also FI_RADIX2FFT_DEMO.
%
% Reference:
%
% Twiddle factors for Algorithm 1.6.2, p. 45, Charles Van Loan,
% Computational Frameworks for the Fast Fourier Transform, SIAM,
% Philadelphia, 1992.
%
% Copyright 2003-2011 The MathWorks, Inc.
%

```

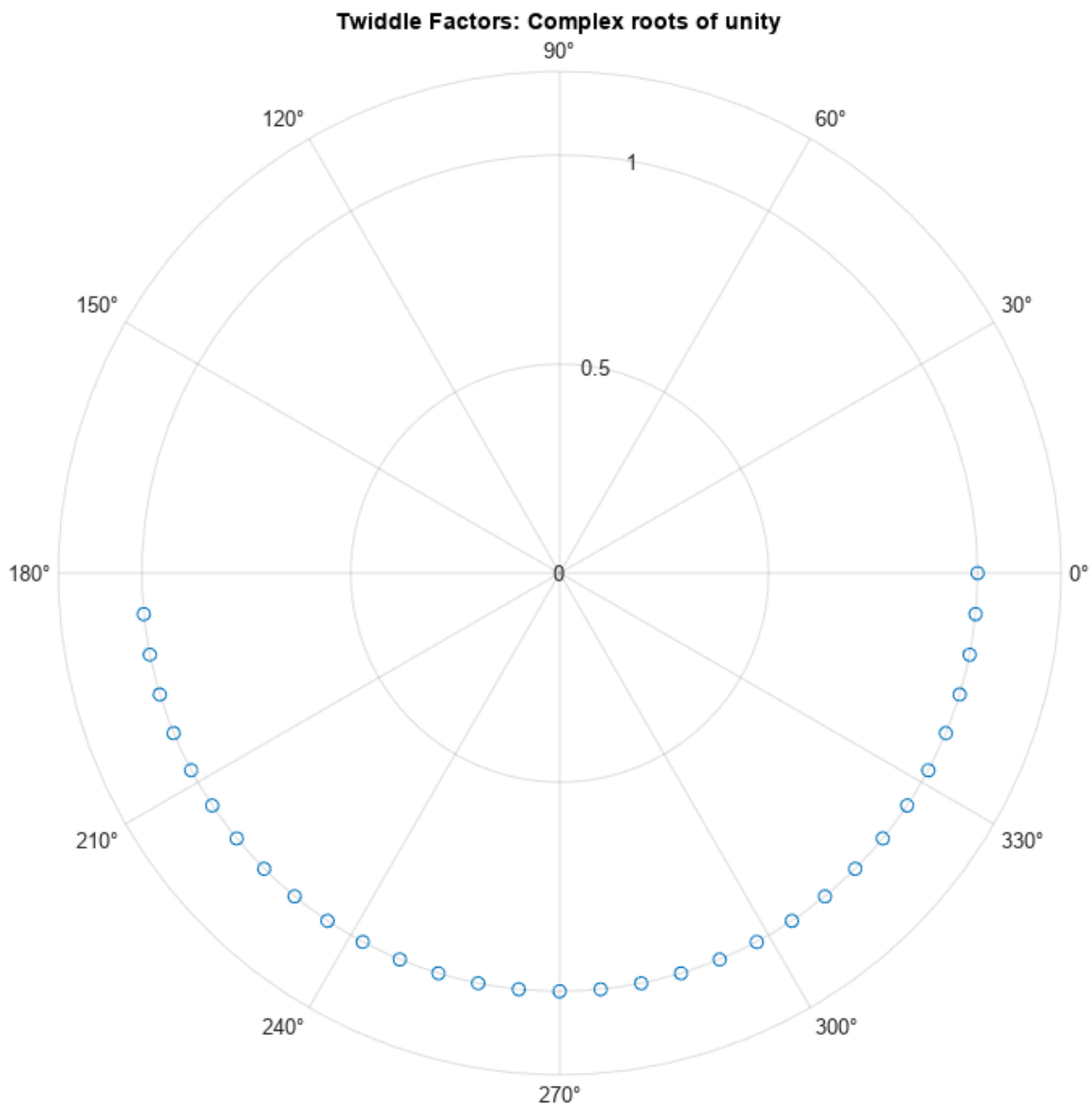


```
t = log2(n);
if floor(t) ~= t
    error('N must be an exact power of two.');
```

end

```
w = zeros(n-1,1);
k = 1;
L = 2;
% Equation 1.4.11, p. 34
while L <= n
    theta = 2*pi/L;
    % Algorithm 1.4.1, p. 23
    for j = 0:(L/2 - 1)
        w(k) = complex(cos(j*theta), -sin(j*theta));
        k = k + 1;
    end
    L = L*2;
end

figure(gcf)
clf
w0 = fi_radix2twiddles(n);
polarplot(angle(w0),abs(w0),'o')
title('Twiddle Factors: Complex roots of unity')
```



Verify Floating-Point Code

To implement the algorithm in MATLAB®, you can use the `fi_bitreverse` function to bit-reverse the input sequence. You must add one to the indices to convert them from zero-based to one-based.

```
function x = fi_m_radix2fft_algorithm1_6_2(x, w)
%FI_M_RADIX2FFT_ALGORITHM1_6_2 Radix-2 FFT example.
% Y = FI_M_RADIX2FFT_ALGORITHM1_6_2(X, W) computes the radix-2 FFT of
% input vector X with twiddle-factors W. Input X is assumed to be
```

```

% complex.
%
% The length of vector X must be an exact power of two.
% Twiddle-factors W are computed via
%   W = fi_radix2twiddles(N)
% where N = length(X).
%
% This version of the algorithm has no scaling before the stages.
%
% See also FI_RADIX2FFT_DEMO, FI_M_RADIX2FFT_WITHSCALING.
%
% Reference:
%   Charles Van Loan, Computational Frameworks for the Fast Fourier
%   Transform, SIAM, Philadelphia, 1992, Algorithm 1.6.2, p. 45.
%
% Copyright 2004-2015 The MathWorks, Inc.

n = length(x); t = log2(n);
x = fi_bitreverse(x,n);
for q=1:t
    L = 2^q; r = n/L; L2 = L/2;
    for k = 0:(r-1)
        for j = 0:(L2-1)
            temp = w(L2-1+j+1) * x(k*L+j+L2+1);
            x(k*L+j+L2+1) = x(k*L+j+1) - temp;
            x(k*L+j+1) = x(k*L+j+1) + temp;
        end
    end
end
end
end

```

Visualization

To verify that you correctly implemented the algorithm in MATLAB®, run a known signal through it and compare the results to the results produced by the MATLAB® FFT function.

As seen in the plot below, the error is within tolerance of the MATLAB® built-in FFT function, verifying that you have correctly implemented the algorithm.

```

y = fi_m_radix2fft_algorithm1_6_2(x0, w0);

fi_fft_demo_plot(real(x0),y,y0,Fs,'Double data', ...
    {'FFT Algorithm 1.6.2','Built-in FFT'});

```

Convert Functions to Use Types Tables

To separate data types from the algorithm:

- 1 Create a table of data type definitions.
- 2 Modify the algorithm code to use data types from that table.

This example shows the iterative steps by creating different files. In practice, you can make the iterative changes to the same file.

Original types table

Create a types table using a structure with prototypes for the variables set to their original types. Use the baseline types to validate that you made the initial conversion correctly, and to programmatically

toggle your function between floating point and fixed point types. The index variables are automatically converted to integers by MATLAB® Coder™, so you don't need to specify their types in the table.

Specify the prototype values as empty ([]) since the data types are used, but not the values.

```
function T = fi_m_radix2fft_original_types()
%FI_M_RADIX2FFT_ORIGINAL_TYPES Types Table Example

% Copyright 2015 The MathWorks, Inc.

    T.x = double([]);
    T.w = double([]);
    T.n = double([]);

end
```

Type-aware algorithm function

Add types table T as an input to the function and use it to cast variables to a particular type, while keeping the body of the algorithm unchanged.

```
function x = fi_m_radix2fft_algorithm1_6_2_typed(x, w, T)
%FI_M_RADIX2FFT_ORIGINAL_TYPED Radix-2 FFT example.
% Y = FI_M_RADIX2FFT_ALGORITHM1_6_2_TYPED(X, W, T) computes the radix-2
% FFT of input vector X with twiddle-factors W. Input X is assumed to be
% complex.
%
% The length of vector X must be an exact power of two.
% Twiddle-factors W are computed via
%     W = fi_radix2twiddles(N)
% where N = length(X).
%
% T is a types table to cast variables to a particular type, while keeping
% the body of the algorithm unchanged.
%
% This version of the algorithm has no scaling before the stages.
%
% See also FI_RADIX2FFT_DEMO, FI_M_RADIX2FFT_WITHSCALING.
%
% Reference:
% Charles Van Loan, Computational Frameworks for the Fast Fourier
% Transform, SIAM, Philadelphia, 1992, Algorithm 1.6.2, p. 45.
%
% Copyright 2015 The MathWorks, Inc.
%
%#codegen

    n = length(x);
    t = log2(n);
    x = fi_bitreverse_typed(x,n,T);
    LL = cast(2.^(1:t), 'like', T.n);
    rr = cast(n./LL, 'like', T.n);
    LL2 = cast(LL./2, 'like', T.n);
    for q=1:t
        L = LL(q);
        r = rr(q);
        L2 = LL2(q);
```

```

    for k = 0:(r-1)
        for j = 0:(L2-1)
            temp = w(L2-1+j+1) * x(k*L+j+L2+1);
            x(k*L+j+L2+1) = x(k*L+j+1) - temp;
            x(k*L+j+1) = x(k*L+j+1) + temp;
        end
    end
end
end
end

```

Type-aware bit-reversal function

Add types table T as an input to the function and use it to cast variables to a particular type, while keeping the body of the algorithm unchanged.

```

function x = fi_bitreverse_typed(x,n0,T)
%FI_BITREVERSE_TYPED Bit-reverse the input.
% X = FI_BITREVERSE_TYPED(x,n,T) bit-reverse the input sequence X, where
% N=length(X).
%
% T is a types table to cast variables to a particular type, while keeping
% the body of the algorithm unchanged.
%
% See also FI_RADIX2FFT_DEMO.

% Copyright 2004-2015 The MathWorks, Inc.
%
%#codegen
n = cast(n0,'like',T.n);
nv2 = bitsra(n,1);
j = cast(1,'like',T.n);
for i = 1:(n-1)
    if i < j
        temp = x(j);
        x(j) = x(i);
        x(i) = temp;
    end
    k = nv2;
    while k < j
        j(:) = j-k;
        k = bitsra(k,1);
    end
    j(:) = j+k;
end
end

```

Validate modified function

Every time you modify your function, validate that the results still match your baseline. Since you used the original types in the types table, the outputs should be identical. This validates that you made the conversion to separate the types from the algorithm correctly.

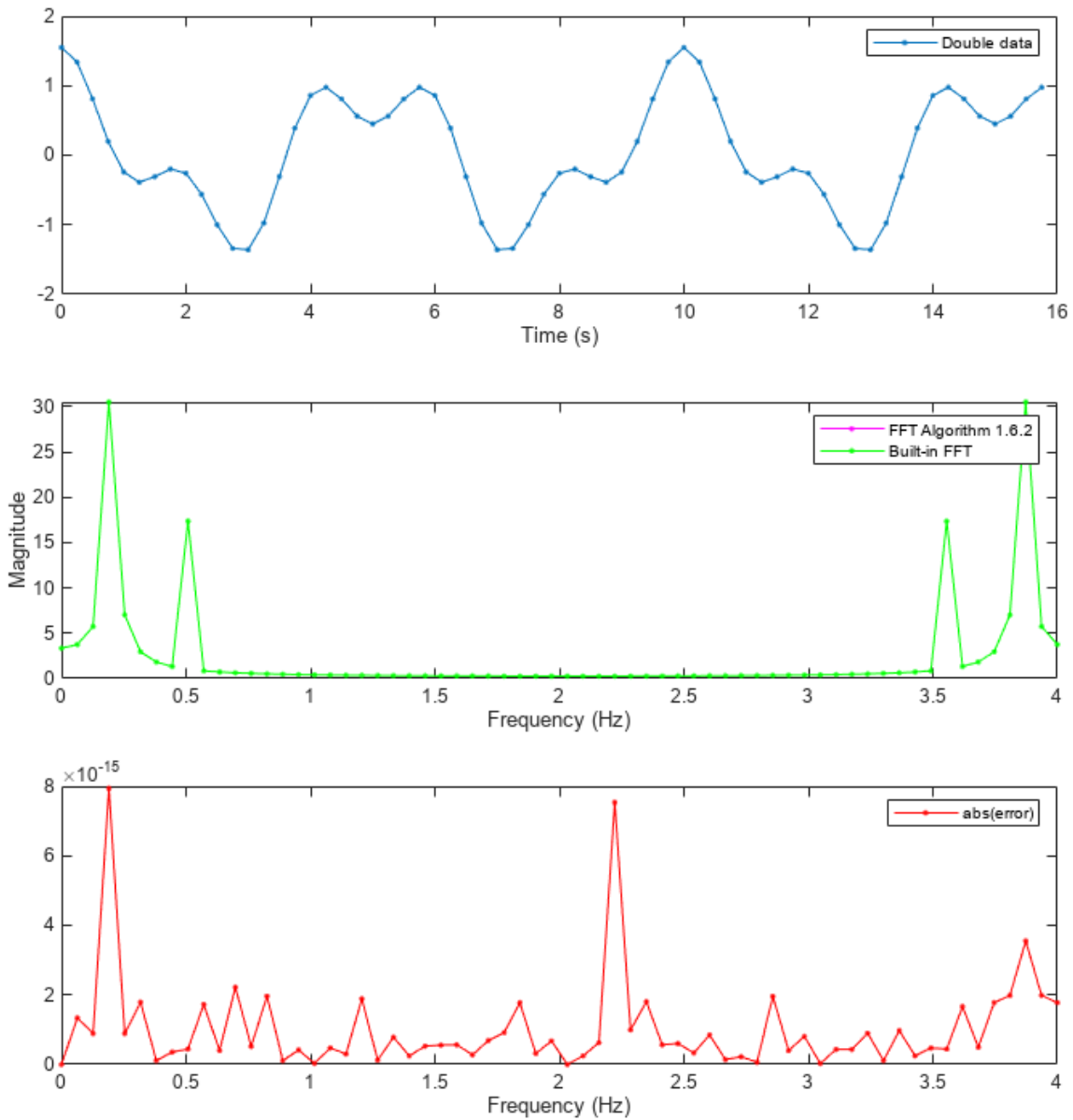
```

T1 = fi_m_radix2fft_original_types(); % Get original data types declared in table
x = cast(x0,'like',T1.x);
w = cast(w0,'like',T1.w);

y = fi_m_radix2fft_algorithm1_6_2_typed(x, w, T1);

```

```
fi_fft_demo_plot(real(x),y,y0,Fs,'Double data', ...
    {'FFT Algorithm 1.6.2','Built-in FFT'});
```



Create a Fixed-Point Types Table

Create a fixed-point types table using a structure with prototypes for the variables. Specify the prototype values as empty (`[]`) since the data types are used, but not the values.

```

function T = fi_m_radix2fft_fixed_types()
%FI_M_RADIX2FFT_FIXED_TYPES Example function

% Copyright 2015 The MathWorks, Inc.

    T.x = fi([],1,16,14); % Picked the following types to ensure that the
    T.w = fi([],1,16,14); % inputs have maximum precision and will not
                        % overflow
    T.n = int32([]);     % Picked int32 as n is an index

end

```

Identify Fixed-Point Issues

Now, try converting the input data to fixed-point and see if the algorithm still looks good. In this first pass, you use all the defaults for signed fixed-point data by using the `fi` constructor.

```

T2 = fi_m_radix2fft_fixed_types(); % Get fixed point data types declared in table

x = cast(x0, 'like', T2.x);
w = cast(w0, 'like', T2.w);

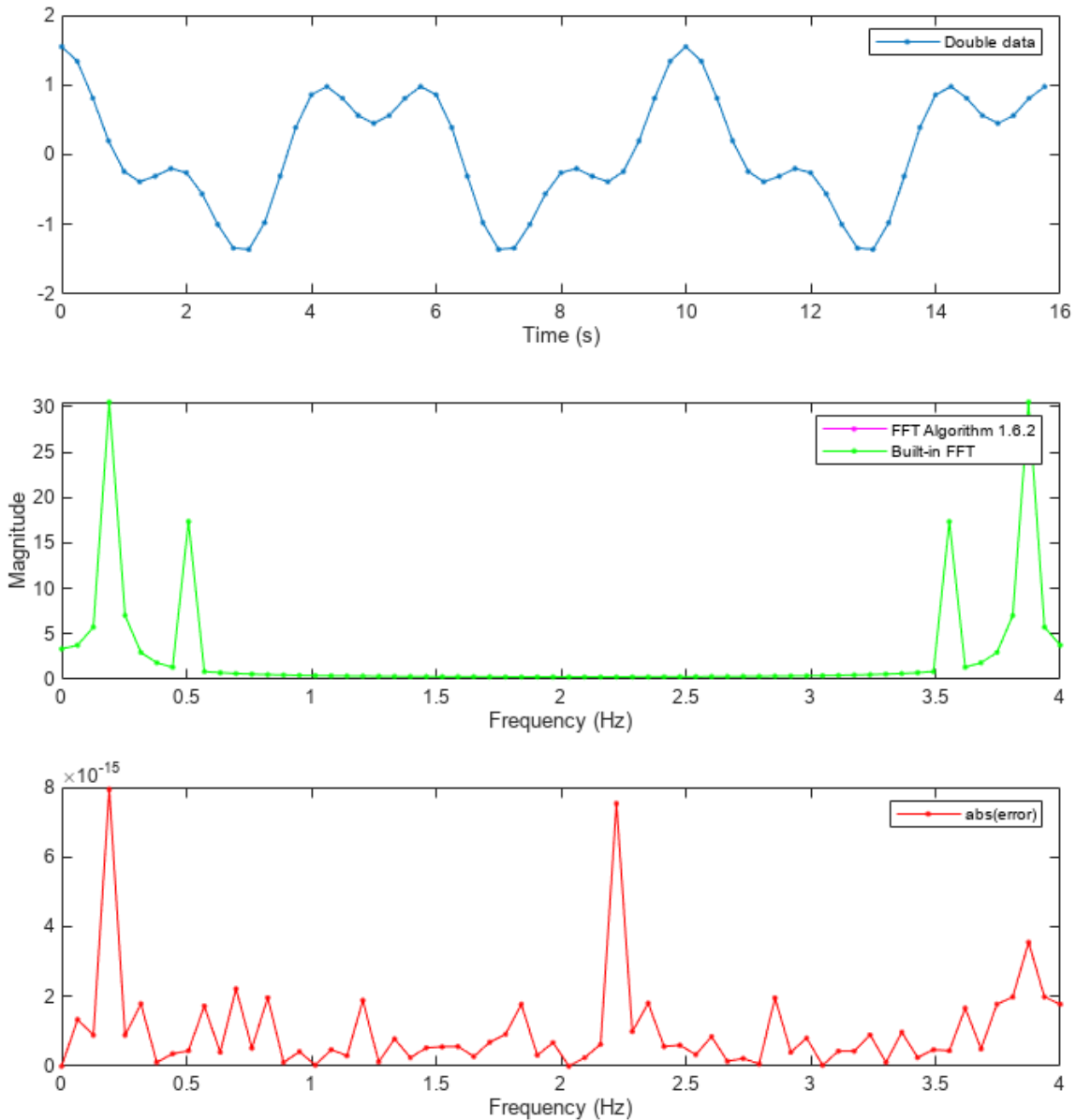
```

Re-run the same algorithm with the fixed-point inputs.

```

y = fi_m_radix2fft_algorithm1_6_2_typed(x,w,T2);
fi_fft_demo_plot(real(x),y,y0,Fs,'Fixed-point data', ...
    {'Fixed-point FFT Algorithm 1.6.2','Built-in FFT'});

```



Note that the magnitude plot (center) of the fixed-point FFT does not resemble the plot of the built-in FFT. The error (bottom plot) is much larger than what you would expect to see for round off error, so it is likely that overflow has occurred.

Use Min/Max Instrumentation to Identify Overflows

To instrument the MATLAB® code, create a MEX function from the MATLAB® function using the `buildInstrumentedMex` command. The inputs to `buildInstrumentedMex` are the same as the

inputs to `fiaccel`, but `buildInstrumentedMex` has no `fi` object restrictions. The output of `buildInstrumentedMex` is a MEX function with instrumentation inserted, so when the MEX function is run, the simulated minimum and maximum values are recorded for all named variables and intermediate values.

The `'-o'` option is used to name the MEX function that is generated. If the `'-o'` option is not used, then the MEX function is the name of the MATLAB® function with `'_mex'` appended. You can also name the MEX function the same as the MATLAB® function, but you need to remember that MEX functions take precedence over MATLAB® functions and so changes to the MATLAB® function will not run until either the MEX function is regenerated, or the MEX function is deleted and cleared.

Create the input with a scaled double data type so its values will attain full range and you can identify potential overflows.

```
function T = fi_m_radix2fft_scaled_fixed_types()
%FI_M_RADIX2FFT_SCALDED_FIXED_TYPES_Example function

% Copyright 2015 The MathWorks, Inc.

DT = 'ScaledDouble';           % Data type to be used for fi
                                % constructor
T.x = fi([],1,16,14,'DataType',DT); % Picked the following types to
T.w = fi([],1,16,14,'DataType',DT); % ensure that the inputs have
                                % maximum precision and will not
                                % overflow
T.n = int32([]);                % Picked int32 as n is an index

end

T3 = fi_m_radix2fft_scaled_fixed_types(); % Get fixed point data types declared in table

x_scaled_double = cast(x0,'like',T3.x);
w_scaled_double = cast(w0,'like',T3.w);

buildInstrumentedMex fi_m_radix2fft_algorithm1_6_2_typed ...
    -o fft_instrumented -args {x_scaled_double w_scaled_double T3}

Run the instrumented MEX function to record min/max values.

y_scaled_double = fft_instrumented(x_scaled_double,w_scaled_double,T3);

Show the instrumentation results.

showInstrumentationResults fft_instrumented
```

The screenshot shows the MATLAB Instrumentation Results window. The top pane displays the source code for the function `fi_m_radix2fft_algorithm1_6_2_typed`. The code includes comments and a loop for `q = 1:t`. The bottom pane shows a table of variables with the following data:

| Name | Type | Size | Class | DT Mode | Signedness: WL | FL | Percent of Current Range | Always Whole Number | Sim Min | Sim Max | |
|------|-------|--------|---------------------|--------------|----------------|----|--------------------------|---------------------|---------|--------------------|--------------------|
| x | I/O | 1 × 64 | complex embedded fi | ScaledDouble | Signed | 16 | 14 | 1263 | No | -17.04520175389448 | 25.320138438385612 |
| T | Input | 1 × 1 | struct | | | | | No | | | |
| w | Input | 63 × 1 | complex embedded fi | ScaledDouble | Signed | 16 | 14 | 51 | No | -1 | 1 |
| j | Local | 1 × 1 | int32 | | | | | Yes | 0 | 31 | |
| k | Local | 1 × 1 | int32 | | | | | Yes | 0 | 31 | |
| L | Local | 1 × 1 | int32 | | | | | Yes | 2 | 64 | |
| L2 | Local | 1 × 1 | int32 | | | | | Yes | 1 | 32 | |
| LL | Local | 1 × 6 | int32 | | | | | Yes | 2 | 64 | |
| LL2 | Local | 1 × 6 | int32 | | | | | Yes | 1 | 32 | |

You can see from the instrumentation results that there were overflows when assigning into the variable `x`.

Modify the Algorithm to Address Fixed-Point Issues

The magnitude of an individual bin in the FFT grows, at most, by a factor of n , where n is the length of the FFT. Hence, by scaling your data by $1/n$, you can prevent overflow from occurring for any input. When you scale only the input to the first stage of a length- n FFT by $1/n$, you obtain a signal-to-noise ratio proportional to n^2 [Oppenheim & Schaffer 1989, equation 9.101], [Welch 1969]. However, if you scale the input to each of the stages of the FFT by $1/2$, you can obtain an overall scaling of $1/n$ and produce a signal-to-noise ratio proportional to n [Oppenheim & Schaffer 1989, equation 9.105], [Welch 1969].

An efficient way to scale by $1/2$ in fixed-point is to right-shift the data. To do this, you use the bit shift right arithmetic function `bitsra`. After scaling each stage of the FFT, and optimizing the index variable computation, your algorithm becomes:

```
function x = fi_m_radix2fft_withscaling_typed(x, w, T)
%FI_M_RADIX2FFT_WITHSCALING Radix-2 FFT example with scaling at each stage.
% Y = FI_M_RADIX2FFT_WITHSCALING_TYPED(X, W, T) computes the radix-2 FFT of
% input vector X with twiddle-factors W with scaling by 1/2 at each stage.
% Input X is assumed to be complex.
%
% The length of vector X must be an exact power of two.
% Twiddle-factors W are computed via
% W = fi_radix2twiddles(N)
```

```

% where N = length(X).
%
% T is a types table to cast variables to a particular type, while keeping
% the body of the algorithm unchanged.
%
% This version of the algorithm has no scaling before the stages.
%
% See also FI_RADIX2FFT_DEMO.

% Reference:
% Charles Van Loan, Computational Frameworks for the Fast Fourier
% Transform, SIAM, Philadelphia, 1992, Algorithm 1.6.2, p. 45.
%
% Copyright 2004-2015 The MathWorks, Inc.
%
%#codegen

n = length(x); t = log2(n);
x = fi_bitreverse(x,n);

% Generate index variables as integer constants so they are not computed in
% the loop.
LL = cast(2.^(1:t), 'like', T.n);
rr = cast(n./LL, 'like', T.n);
LL2 = cast(LL./2, 'like', T.n);
for q = 1:t
    L = LL(q); r = rr(q); L2 = LL2(q);
    for k = 0:(r-1)
        for j = 0:(L2-1)
            temp = w(L2-1+j+1) * x(k*L+j+L2+1);
            x(k*L+j+L2+1) = bitsra(x(k*L+j+1) - temp, 1);
            x(k*L+j+1) = bitsra(x(k*L+j+1) + temp, 1);
        end
    end
end
end

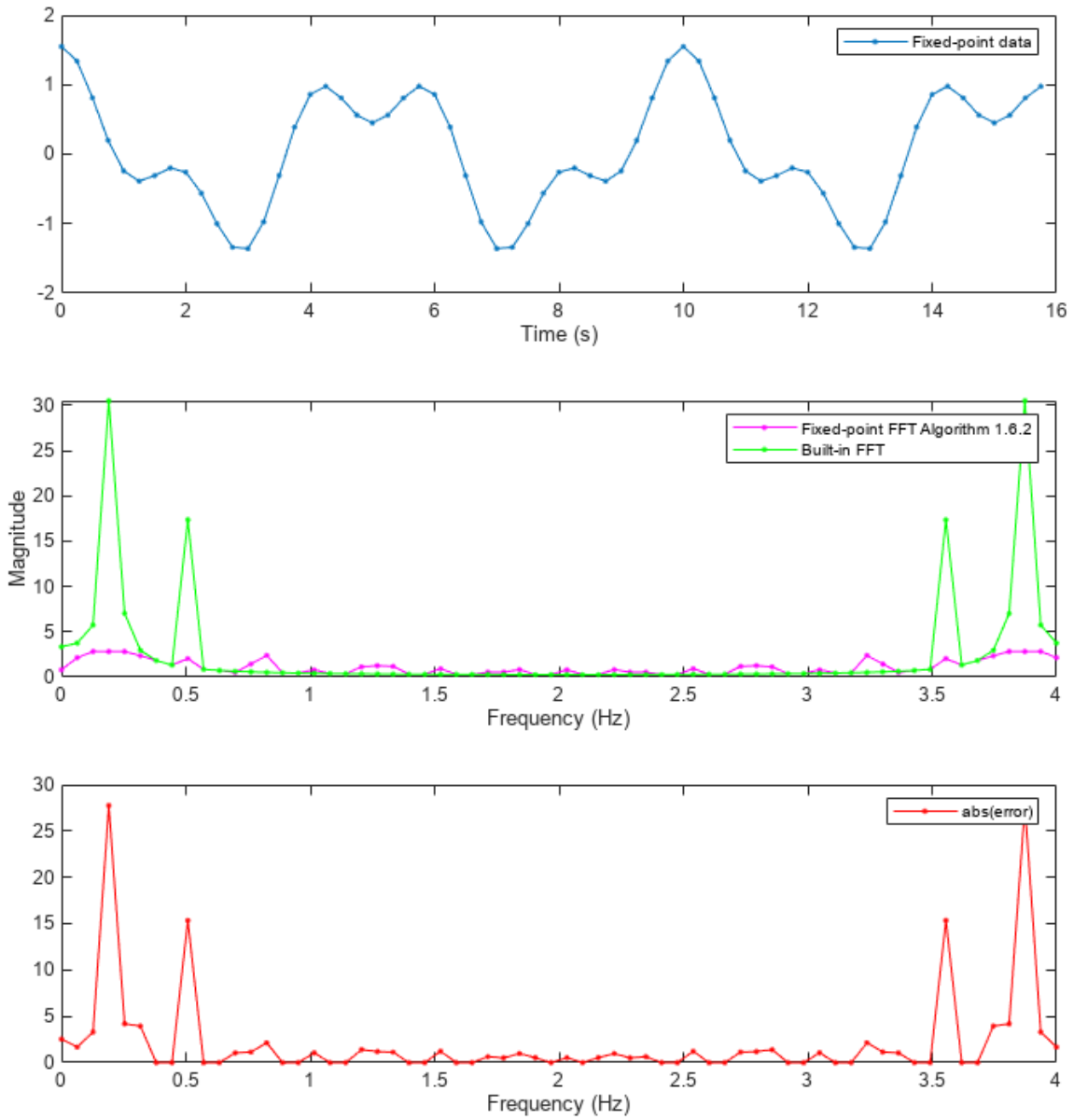
Run the scaled algorithm with fixed-point data.

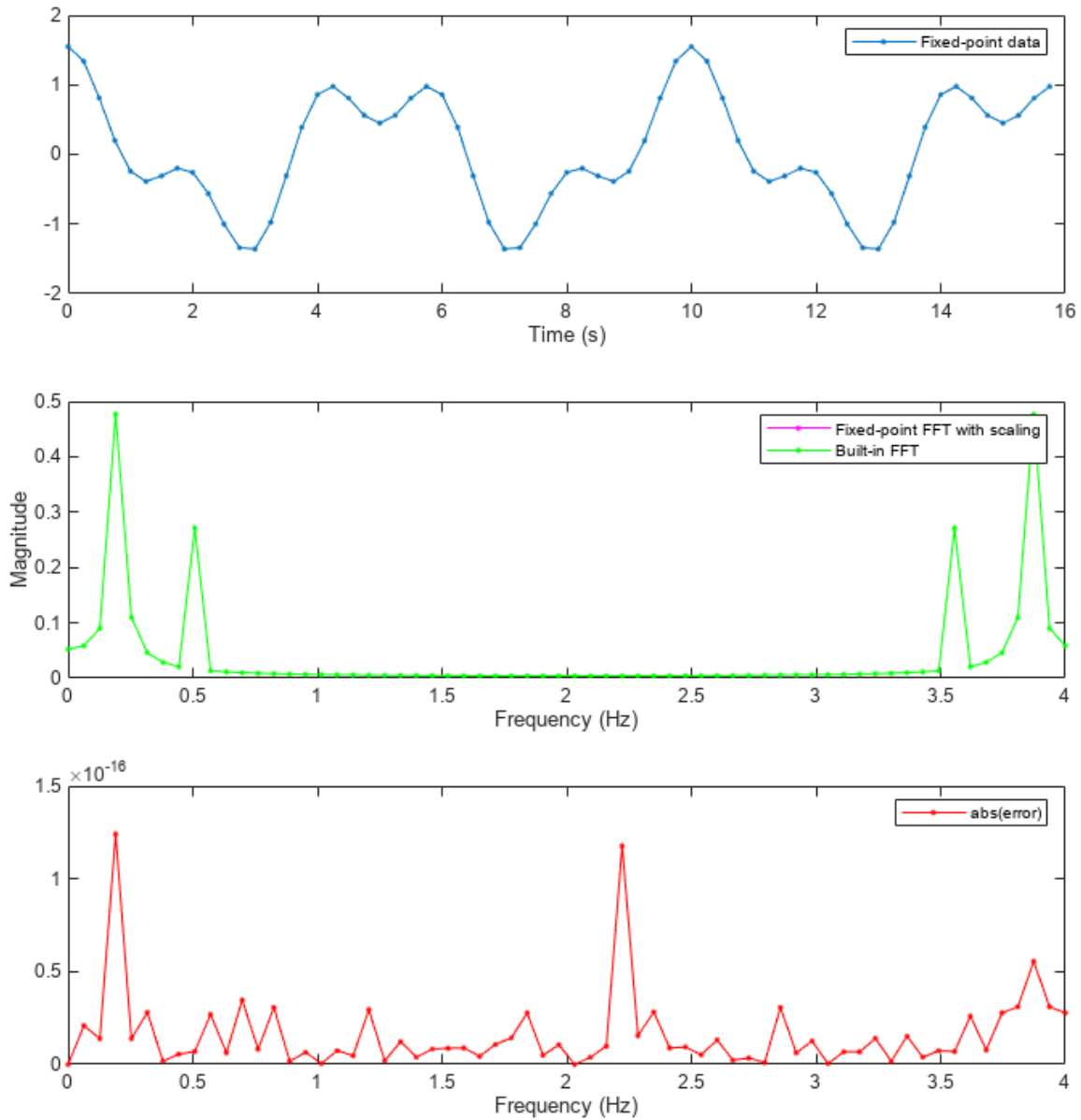
x = cast(x0, 'like', T3.x);
w = cast(w0, 'like', T3.w);

y = fi_m_radix2fft_withscaling_typed(x,w,T3);

fi_fft_demo_plot(real(x), y, y0/n, Fs, 'Fixed-point data', ...
    {'Fixed-point FFT with scaling', 'Built-in FFT'});

```





You can see that the scaled fixed-point FFT algorithm now matches the built-in FFT to a tolerance that is expected for 16-bit fixed-point data.

Clean Up

Run the following code to restore the global states.

```
fipref(FIPREF_STATE);
```

References

Charles Van Loan, *Computational Frameworks for the Fast Fourier Transform*, SIAM, 1992.

Cleve Moler, *Numerical Computing with MATLAB*, SIAM, 2004, Chapter 8 Fourier Analysis.

Alan V. Oppenheim and Ronald W. Schaffer, *Discrete-Time Signal Processing*, Prentice Hall, 1989.

Peter D. Welch, "A Fixed-Point Fast Fourier Transform Error Analysis," *IEEE® Transactions on Audio and Electroacoustics*, Vol. AU-17, No. 2, June 1969, pp. 151-157.

Set Data Types Using Min/Max Instrumentation

In this example, you set fixed-point data types by instrumenting MATLAB® code for min/max logging then use the tools to propose data types. You use the `buildInstrumentedMex` function to build the MEX function with instrumentation enabled, then use the `showInstrumentationResults` to show instrumentation results and the `clearInstrumentationResults` function to clear instrumentation results.

Define the Unit Under Test

The function that you convert to fixed-point in this example is a second-order direct-form 2 transposed filter. You can substitute your own function in place of this one to reproduce these steps in your own work.

```
function [y,z] = fi_2nd_order_df2t_filter(b,a,x,y,z)
    for i=1:length(x)
        y(i) = b(1)*x(i) + z(1);
        z(1) = b(2)*x(i) + z(2) - a(2) * y(i);
        z(2) = b(3)*x(i)          - a(3) * y(i);
    end
end
```

For a MATLAB® function to be instrumented, it must be suitable for code generation. For information on code generation, see the `buildInstrumentedMex` reference page. A MATLAB® Coder™ license is not required to use the `buildInstrumentedMex` function.

In the function `fi_2nd_order_df2t_filter`, the variables `y` and `z` are used as both inputs and outputs. This is an important pattern because:

- You can set the data type of `y` and `z` outside the function, thus allowing you to re-use the function for both fixed-point and floating-point types.
- The generated C code will create `y` and `z` as references in the function argument list.

Use Design Requirements to Determine Data Types

In this example, the requirements of the design determine the data type of input `x`. These requirements are signed, 16-bit, and fractional.

```
clearvars
N = 256;
x = fi(zeros(N,1),1,16,15);
```

The requirements of the design also determine the fixed-point math for a DSP target with a 40-bit accumulator. This example uses floor rounding and wrap on overflow to produce efficient generated code.

```
F = fimath('RoundingMethod','Floor',...
          'OverflowAction','Wrap',...
          'ProductMode','KeepLSB',...
          'ProductWordLength',40,...
          'SumMode','KeepLSB',...
          'SumWordLength',40);
```

The following coefficients correspond to a second-order lowpass filter created by

```
[num,den] = butter(2,0.125)
```

The values of the coefficients influence the range of the values that will be assigned to the filter output and states.

```
num = [0.0299545822080925  0.0599091644161849  0.0299545822080925];
den = [1                    -1.4542435862515900  0.5740619150839550];
```

The data type of the coefficients, determined by the requirements of the design, are specified as 16-bit word length and scaled to best-precision. To create `fi` objects from constant coefficients:

1. Cast the coefficients to `fi` objects using the default round-to-nearest and saturate on overflow settings, which gives the coefficients better accuracy.

```
b = fi(num,1,16);
a = fi(den,1,16);
```

2. Attach `fimath` with floor rounding and wrap on overflow settings to control arithmetic, which leads to more efficient C code.

```
b = setfimath(b,F);
a = setfimath(a,F);
```

Use Values of the Coefficients and Inputs to Determine Data Types

The values of the coefficients and values of the inputs determine the data types of output `y` and state vector `z`. Create them with a scaled double datatype so their values will attain full range and you can identify potential overflows and propose data types.

```
yisd = fi(zeros(N,1),1,16,15,'DataType','ScaledDouble','fimath',F);
zisd = fi(zeros(2,1),1,16,15,'DataType','ScaledDouble','fimath',F);
```

Instrument the MATLAB Function as a Scaled-Double MEX Function

To instrument the MATLAB® code, you create a MEX function from the MATLAB® function using the `buildInstrumentedMex` function. The inputs to `buildInstrumentedMex` are the same as the inputs to `fiaccel`, but `buildInstrumentedMex` has no `fi`-object restrictions. The output of `buildInstrumentedMex` is a MEX function with instrumentation inserted. When the MEX function is run, the simulated minimum and maximum values are recorded for all named variables and intermediate values.

Use the `'-o'` option to name the MEX function that is generated. If you do not use the `'-o'` option, then the MEX function is the name of the MATLAB® function with `'_mex'` appended. You can also name the MEX function the same as the MATLAB® function, but you need to remember that MEX functions take precedence over MATLAB® functions and so changes to the MATLAB® function will not run until either the MEX function is re-generated, or the MEX function is deleted and cleared.

Hard-code the filter coefficients into the implementation of this filter by passing them as constants to the `buildInstrumentedMex` function.

```
buildInstrumentedMex fi_2nd_order_df2t_filter ...
    -o filter_scaled_double ...
    -args {coder.Constant(b),coder.Constant(a),x,yisd,zisd}
```


Test Bench with Chirp Input

The test bench for this system is set up to run chirp and step signals. In general, test benches for systems should cover a wide range of input signals.

The first test bench uses a chirp input. A chirp signal is a good representative input because it covers a wide range of frequencies.

```
t = linspace(0,1,N);      % Time vector from 0 to 1 second
f1 = N/2;                % Target frequency of chirp set to Nyquist
xchirp = sin(pi*f1*t.^2); % Linear chirp from 0 to Fs/2 Hz in 1 second
x(:) = xchirp;           % Cast the chirp to fixed-point
```

Run the Instrumented MEX Function to Record Min/Max Values

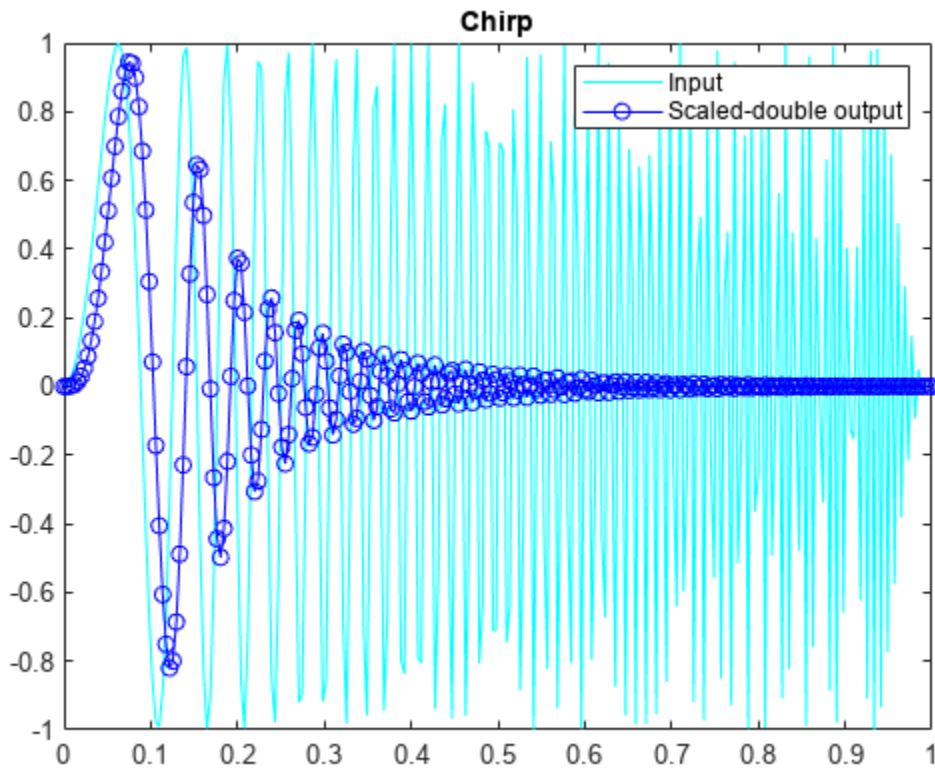
The instrumented MEX function must be run to record minimum and maximum values for that simulation run. Subsequent runs accumulate the instrumentation results until they are cleared with `clearInstrumentationResults`.

Note that the numerator and denominator coefficients were compiled as constants so they are not provided as input to the generated MEX function.

```
ychirp = filter_scaled_double(b,a,x,yisd,zisd);
```

The plot of the filtered chirp signal shows the lowpass behavior of the filter with these particular coefficients. Low frequencies are passed through and higher frequencies are attenuated.

```
plot(t,x,'c',t,ychirp,'bo-')
title('Chirp')
legend('Input','Scaled-double output')
```



Show Instrumentation Results with Proposed Fraction Lengths for Chirp

The `showInstrumentationResults` function displays the code generation report with instrumented values. The input to the `showInstrumentationResults` function is the name of the instrumented MEX function for which you wish to show results.

Potential overflows are only displayed for `fi` objects with Scaled Double data type.

This particular design is for a DSP where the word lengths are fixed, so use the `-proposeFL` flag to propose fraction lengths.

```
showInstrumentationResults filter_scaled_double -proposeFL
```

Hover over expressions or variables in the instrumented code generation report to see the simulation minimum and maximum values. In this design, the inputs fall between -1 and +1, and the values of all variables and intermediate results also fall between -1 and +1. This suggests that the data types can all be fractional (fraction length one bit less than the word length). However, this will not always be true for this function for other kinds of inputs and it is important to test many types of inputs before setting final fixed-point data types.

The screenshot shows the MATLAB Instrumentation Reporter interface. The top part displays the MATLAB source code for a function named `fi_2nd_order_df2t_filter`. The code is as follows:

```

1 function [y,z] = fi_2nd_order_df2t_filter(b,a,x)
2 %FI_2ND_ORDER_DF2T_FILTER Second-order Direct-Form II
3
4 % Copyright 2011 The MathWorks, Inc.
5 for i=1:length(x)
6     y(i) = b(1)*x(i) + z(1);
7     z(1) = b(2)*x(i) + z(2) - a(2) * y(i);
8     z(2) = b(3)*x(i) - a(3) * y(i);
9 end
10 end
11
12

```

The bottom part of the screenshot shows the 'INSTRUMENTATION RESULTS' table. The table has columns for Name, Type, Size, Class, and various instrumentation metrics. The data is as follows:

| Name | Type | Size | Class | Percent of Current Range | Proposed FL | Percent of Current Range | Always Whole Number | Sim Min | Sim Max |
|------|-------|---------|-------------|---------------------------|-------------|--------------------------|---------------------|---------|---------|
| y | I/O | 256 x 1 | embedded fi | ScaledDouble Signed 16 15 | 15 | 95 | No | -0.81 | 0.944 |
| z | I/O | 2 x 1 | embedded fi | ScaledDouble Signed 16 15 | 15 | 93 | No | -0.80 | 0.922 |
| a | input | 1 x 3 | embedded fi | Signed 16 14 | - | 73 | No | -1.45 | 1 |
| b | input | 1 x 3 | embedded fi | Signed 16 19 | - | 96 | No | 0.029 | 0.056 |
| x | input | 256 x 1 | embedded fi | Signed 16 15 | - | 100 | No | -0.99 | 0.996 |
| i | Local | 1 x 1 | double | - | - | - | Yes | 1 | 256 |

Test Bench with Step Input

The next test bench is run with a step input. A step input is a good representative input because it is often used to characterize the behavior of a system.

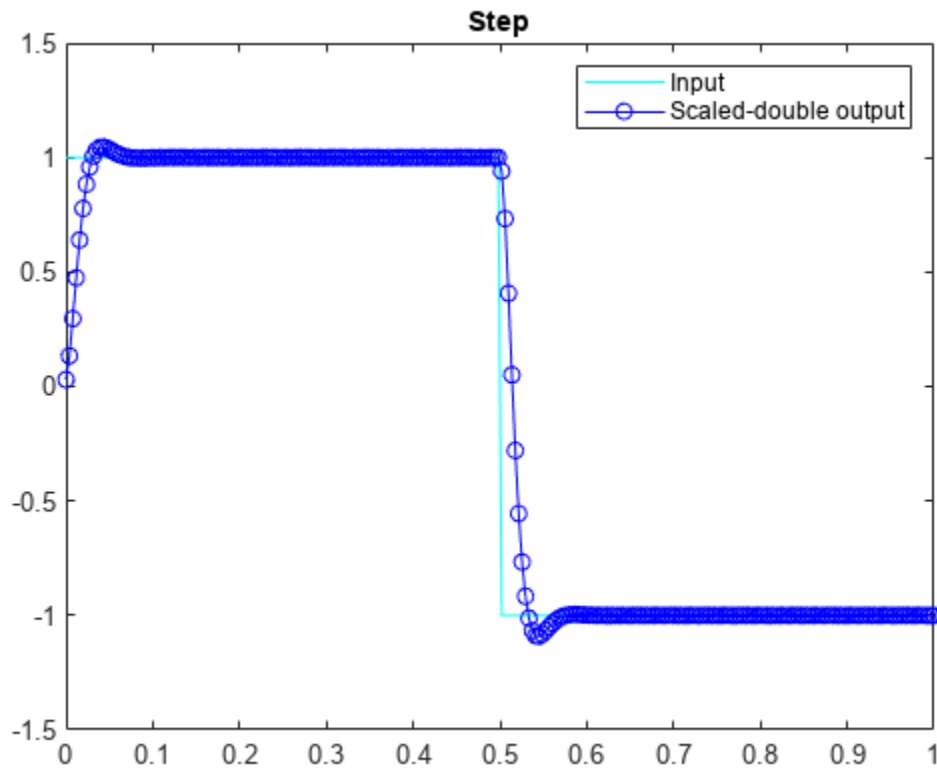
```
xstep = [ones(N/2,1); -ones(N/2,1)];
x(:) = xstep;
```

Run the Instrumented MEX Function with Step Input

The instrumentation results are accumulated until they are cleared with `clearInstrumentationResults`.

```
ystep = filter_scaled_double(b,a,x,yisd,zisd);
```

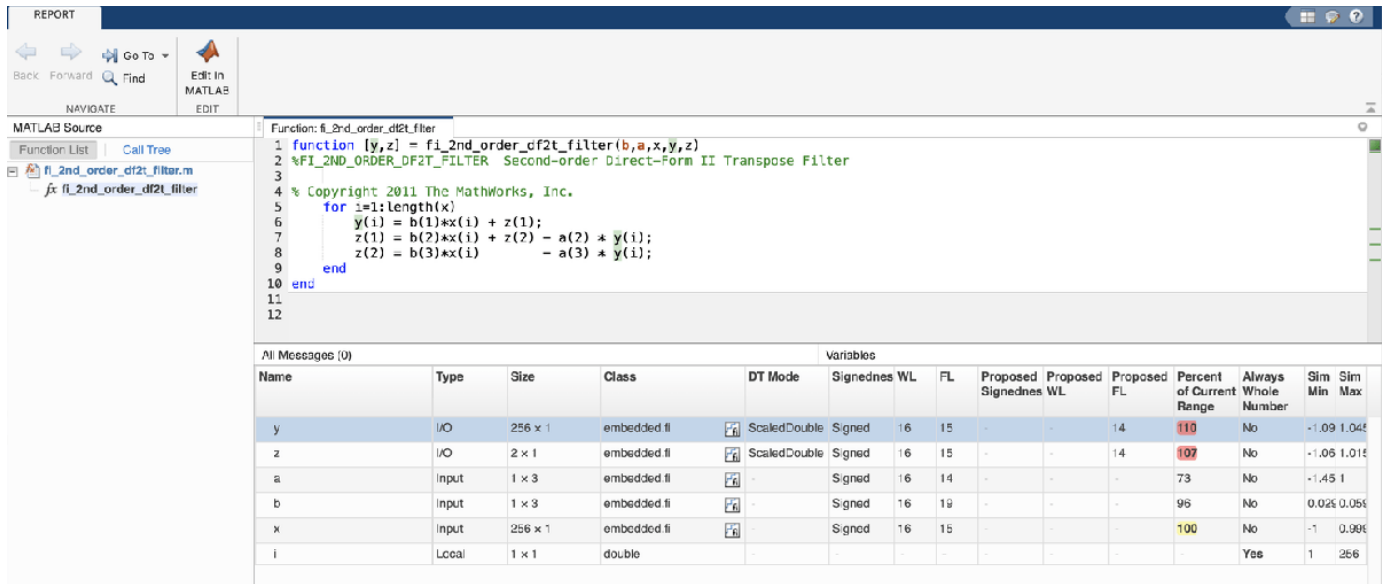
```
plot(t,x,'c',t,ystep,'bo-')
title('Step')
legend('Input', 'Scaled-double output')
```



Show Accumulated Instrumentation Results

Even though the inputs for step and chirp inputs are both full range as indicated by x at 100 percent current range in the instrumented code generation report, the step input causes overflow while the chirp input did not. This is an illustration of the necessity to have many different inputs for your test bench. For the purposes of this example, only two inputs were used, but real test benches should be more thorough.

```
showInstrumentationResults filter_scaled_double -proposeFL
```



Function: fi_2nd_order_df2t_filter

```

1 function [y,z] = fi_2nd_order_df2t_filter(b,a,x,y,z)
2 %FI_2ND_ORDER_DF2T_FILTER Second-order Direct-Form II Transpose Filter
3
4 % Copyright 2011 The MathWorks, Inc.
5     for i=1:length(x)
6         y(i) = b(1)*x(i) + z(1);
7         z(1) = b(2)*x(i) + z(2) - a(2) * y(i);
8         z(2) = b(3)*x(i) - a(3) * y(i);
9     end
10 end
11
12

```

| All Messages (0) | | | | | | | | | | | Variables | | | |
|------------------|-------|---------|-------------|--------------|-----------|----|----|--------------------|-------------|-------------|--------------------------|---------------------|---------|---------|
| Name | Type | Size | Class | DT Mode | Signednes | WL | FL | Proposed Signednes | Proposed WL | Proposed FL | Percent of Current Range | Always Whole Number | Sim Min | Sim Max |
| y | I/O | 256 x 1 | embedded fi | ScaledDouble | Signed | 16 | 15 | - | - | 14 | 110 | No | -1.09 | 1.04 |
| z | I/O | 2 x 1 | embedded fi | ScaledDouble | Signed | 16 | 15 | - | - | 14 | 107 | No | -1.05 | 1.015 |
| a | Input | 1 x 3 | embedded fi | - | Signed | 16 | 14 | - | - | - | 73 | No | -1.45 | 1 |
| b | Input | 1 x 3 | embedded fi | - | Signed | 16 | 14 | - | - | - | 95 | No | 0.025 | 0.05 |
| x | Input | 256 x 1 | embedded fi | - | Signed | 16 | 15 | - | - | - | 100 | No | -1 | 0.99 |
| i | Local | 1 x 1 | double | - | - | - | - | - | - | - | - | Yes | 1 | 256 |

Apply Proposed Fixed-Point Properties

To prevent overflow, set proposed fixed-point properties based on the proposed fraction lengths of 14-bits for y and z from the instrumented code generation report.

At this point in the workflow, you use true fixed-point types (as opposed to the scaled double types that were used in the earlier step of determining data types).

```

yi = fi(zeros(N,1),1,16,14,'fimath',F);
zi = fi(zeros(2,1),1,16,14,'fimath',F);

```

Instrument the MATLAB Function as a Fixed-Point MEX Function

Create an instrumented fixed-point MEX function by using fixed-point inputs and the `buildInstrumentedMex` function.

```

buildInstrumentedMex fi_2nd_order_df2t_filter ...
    -o filter_fixed_point ...
    -args {coder.Constant(b),coder.Constant(a),x,yi,zi}

```

Validate the Fixed-Point Algorithm

After converting to fixed-point, run the test bench again with fixed-point inputs to validate the design.

Validate with Chirp Input

Run the fixed-point algorithm with a chirp input to validate the design.

```

x(:) = xchirp;
[y,z] = filter_fixed_point(b,a,x,yi,zi);
[ysd,zsd] = filter_scaled_double(b,a,x,yisd,zisd);
err = double(y) - double(ysd);

```

Compare the fixed-point outputs to the scaled-double outputs to verify that they meet your design criteria.

```

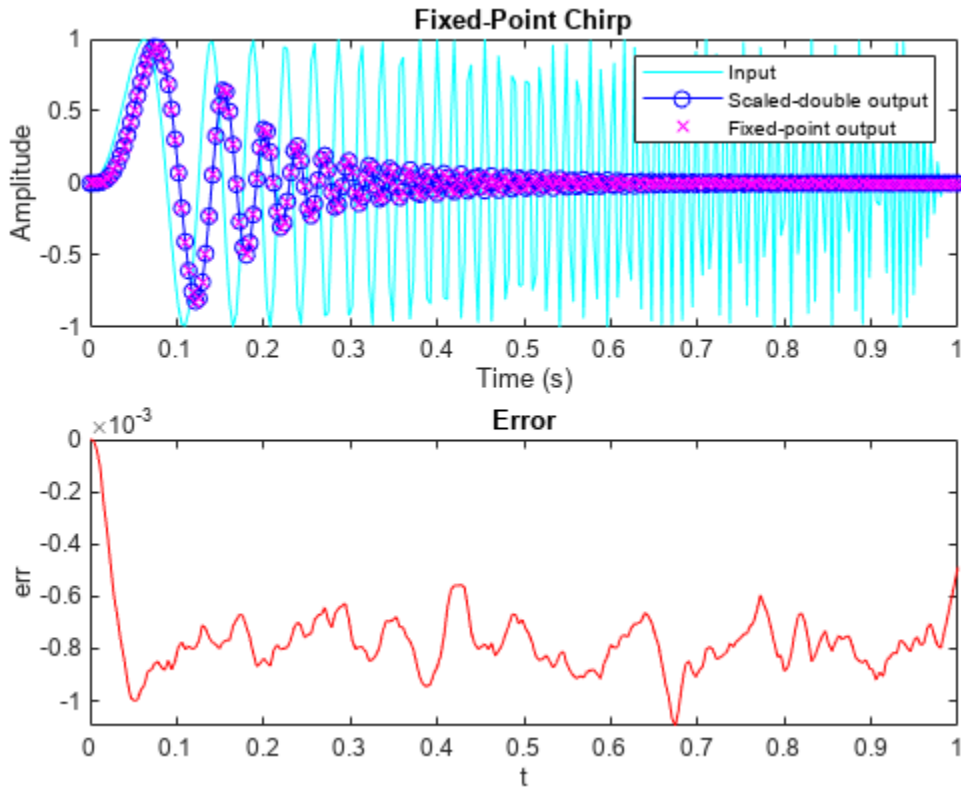
subplot(211);
plot(t,x,'c',t,ysd,'bo-',t,y,'mx')

```

```

xlabel('Time (s)');
ylabel('Amplitude');
legend('Input','Scaled-double output','Fixed-point output');
title('Fixed-Point Chirp')
subplot(212);
plot(t,err,'r');title('Error');xlabel('t'); ylabel('err');

```



Inspect the variables and intermediate results to ensure that the min/max values are within range.

showInstrumentationResults `filter_fixed_point`

The screenshot shows the MATLAB IDE with the following components:

- REPORT** window at the top.
- NAVIGATE** section with Back, Forward, Go To, Find, and Edit In MATLAB buttons.
- MATLAB Source** window showing the function `fi_2nd_order_df2t_filter` with the following code:


```

1 function [y,z] = fi_2nd_order_df2t_filter(b,a,x,y,z)
2 %FI_2ND_ORDER_DF2T_FILTER Second-Order Direct-Form II Transpose Filter
3
4 % Copyright 2011 The MathWorks, Inc.
5 for i=1:length(x)
6     y(i) = b(1)*x(i) + z(1);
7     z(1) = b(2)*x(i) + z(2) - a(2) * y(i);
8     z(2) = b(3)*x(i) - a(3) * y(i);
9 end
10 end
11
12
```
- All Messages (0)** section containing a table of variables:

| Name | Type | Size | Class | DT Mode | Signednes | WL | FL | Percent of Current Range | Always Whole Number | Sim Min | Sim Max |
|------|-------|---------|-------------|---------|-----------|----|----|--------------------------|---------------------|---------------|---------------|
| y | I/O | 256 x 1 | embedded.fi | | Signed | 16 | 14 | 48 | Nc | -0.8206176757 | 0.94379662106 |
| z | I/O | 2 x 1 | embedded.fi | | Signed | 16 | 14 | 47 | Nc | -0.8108079101 | 0.82193603516 |
| a | Input | 1 x 3 | embedded.fi | | Signed | 16 | 14 | 73 | Nc | -1.454226328 | 1 |
| b | Input | 1 x 3 | embedded.fi | | Signed | 16 | 19 | 96 | Nc | 0.02955491027 | 0.05990982056 |
| x | Input | 256 x 1 | embedded.fi | | Signed | 16 | 15 | 100 | Nc | -0.9993694824 | 0.9993896484 |
| i | Local | 1 x 1 | double | | | | | | Yes | 1 | 256 |

Validate with Step Inputs

Run the fixed-point algorithm with a step input to validate the design.

Run the following code to clear the previous instrumentation results to see only the effects of running the step input.

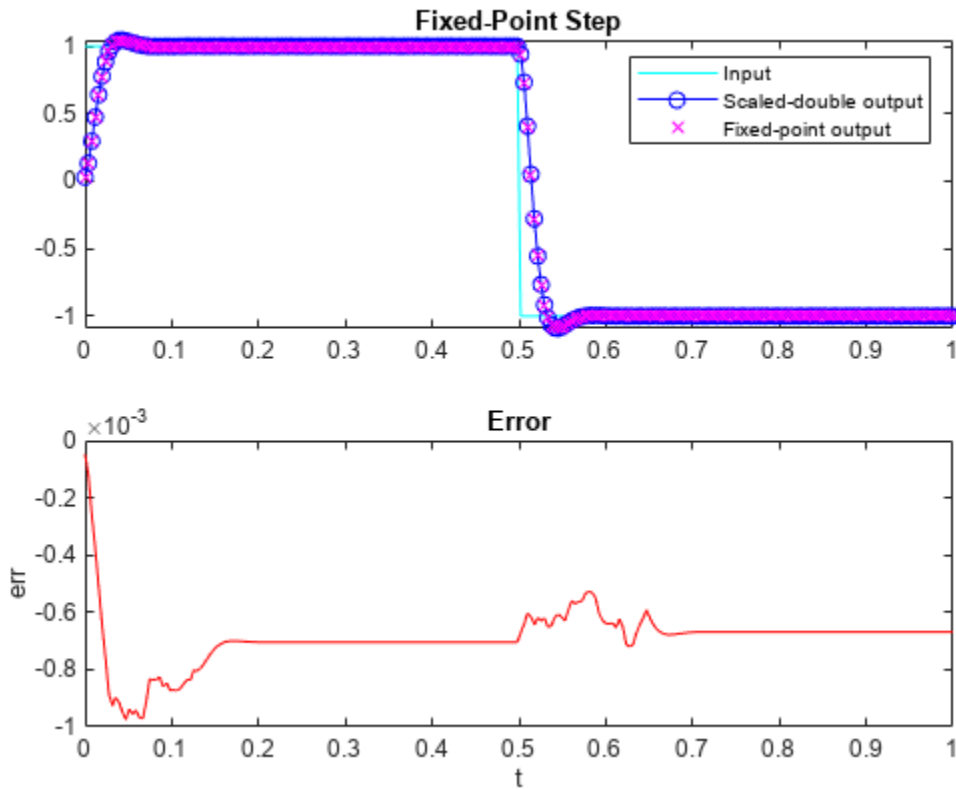
```
clearInstrumentationResults filter_fixed_point
```

Run the step input through the fixed-point filter and compare with the output of the scaled double filter.

```
x(:) = xstep;
[y,z] = filter_fixed_point(b,a,x,yi,zi);
[ysd,zsd] = filter_scaled_double(b,a,x,yisd,zisd);
err = double(y) - double(ysd);
```

Plot the fixed-point outputs against the scaled-double outputs to verify that they meet your design criteria.

```
subplot(211);
plot(t,x,'c',t,ysd,'bo-',t,y,'mx')
title('Fixed-Point Step');
legend('Input', 'Scaled-double output', 'Fixed-point output')
subplot(212);
plot(t,err,'r');title('Error');xlabel('t'); ylabel('err');
```



Inspect the variables and intermediate results to ensure that the min/max values are within range.

showInstrumentationResults filter_fixed_point

REPORT

Back Forward Go To Find Edit In MATLAB

NAVIGATE EDIT

MAT-LAB Source

Function List Call Tree

fi_2nd_order_df2t_filter.m

fx fi_2nd_order_df2t_filter

```

1 function [y,z] = fi_2nd_order_df2t_filter(b,a,x,y,z)
2 %FI_2ND_ORDER_DF2T_FILTER Second-order Direct-Form II Transpose Filter
3
4 % Copyright 2011 The MathWorks, Inc.
5 for i=1:length(x)
6     y(i) = b(1)*x(i) + z(1);
7     z(1) = b(2)*x(i) + z(2) - a(2) * y(i);
8     z(2) = b(3)*x(i) - a(3) * y(i);
9 end
10 end
11
12

```

All Messages (0)

| Name | Type | Size | Class | DT Mode | Signedness | WL | FL | Percent of Current Range | Always Whole Number | Sim Min | Sim Max |
|------|-------|---------|-------------|---------|------------|----|----|--------------------------|---------------------|---------------|--------------|
| y | I/O | 256 x 1 | embedded fi | | Signed | 15 | 14 | 55 | No | -1.0917358398 | 1.0446166992 |
| z | I/O | 2 x 1 | embedded fi | | Signed | 15 | 14 | 64 | No | -1.0617675781 | 1.0147094726 |
| a | Input | 1 x 3 | embedded fi | | Signed | 15 | 14 | 73 | No | -1.454220328 | 1 |
| b | Input | 1 x 3 | embedded fi | | Signed | 15 | 19 | 96 | No | 0.0259549102 | 0.0599098205 |
| x | Input | 256 x 1 | embedded fi | | Signed | 15 | 15 | 100 | No | -1 | 0.9999894824 |
| i | Local | 1 x 1 | double | | - | - | - | - | Yes | 1 | 255 |

Suppress Code Analyzer warnings.

```
%#ok<*ASGLU>
```

Implement FIR Filter Algorithm for Floating-Point and Fixed-Point Types Using cast and zeros

This example shows you how to convert a finite impulse-response (FIR) filter to fixed point by separating the fixed-point type specification from the algorithm code.

Separating data type specification from algorithm code allows you to:

- Re-use your algorithm code with different data types
- Keep your algorithm uncluttered with data type specification and switch statements for different data types
- Keep your algorithm code more readable
- Switch between fixed point and floating point to compare baselines
- Switch between variations of fixed-point settings without changing the algorithm code

Original FIR Filter Algorithm

This example converts MATLAB® code for a finite impulse response (FIR) filter to fixed point.

The formula for the n th output, $|y(n)$, of an FIR filter given filter coefficients b and input x is:

$$|y(n) = b(1)*x(n) + b(2)*x(n-1) + \dots + b(\text{end})*x(n-\text{length}(b)+1)|$$

Linear Buffer Implementation

There are several different ways to write an FIR filter. One way is with a linear buffer like in the following function, where b is a row vector and z is a column vector the same length as b .

```
function [y,z] = fir_filt_linear_buff(b,x,z)
    y = zeros(size(x));
    for n=1:length(x)
        z = [x(n); z(1:end-1)];
        y(n) = b * z;
    end
end
```

The vector z is made up of the current and previous samples of x .

$$z = [x(n); x(n-1); \dots ; x(n-\text{length}(b)+1)]$$

The vector z is called a state vector. It is used as both an input and an output. When it is an input, it is the initial state of the filter. When it is an output, it is the final state of the filter. Defining the state vector outside the function allows you to continuously stream data through the filter in a real-time system, and to reuse the filter algorithm within the same project. The state vector is often named z because it is synonymous with the delays z^{-1} associated with the Z-transform. The Z-transform was named in honor of Lotfi Zadeh for his seminal work in this area [1].

The linear buffer implementation takes advantage of MATLAB's convenient matrix syntax and is easy to read and understand. However, it introduces a full copy of the state buffer for every sample of the input.

Circular Buffer Implementation

To implement the FIR filter more efficiently, you can store the states in a circular buffer, z , whose elements are $z(p) = x(n)$, where $p = \text{mod}(n-1, \text{length}(b))+1$, for $n = 1, 2, 3, \dots$

For example, let $\text{length}(b) = 3$, and initialize p and z to zero.

$$p = 0, z = [0 \ 0 \ 0]$$

Start with the first sample and fill the state buffer z in a circular manner.

$$n = 1, p = 1, z(1) = x(1), z = [x(1) \ 0 \ 0] \\ y(1) = b(1)*z(1) + b(2)*z(3) + b(3)*z(2)$$

$$n = 2, p = 2, z(2) = x(2), z = [x(1) \ x(2) \ 0] \\ y(2) = b(1)*z(2) + b(2)*z(1) + b(3)*z(3)$$

$$n = 3, p = 3, z(3) = x(3), z = [x(1) \ x(2) \ x(3)] \\ y(3) = b(1)*z(3) + b(2)*z(2) + b(3)*z(1)$$

$$n = 4, p = 1, z(1) = x(4), z = [x(4) \ x(2) \ x(3)] \\ y(4) = b(1)*z(1) + b(2)*z(3) + b(3)*z(2)$$

$$n = 5, p = 2, z(2) = x(5), z = [x(4) \ x(5) \ x(3)] \\ y(5) = b(1)*z(2) + b(2)*z(1) + b(3)*z(3)$$

$$n = 6, p = 3, z(3) = x(6), z = [x(4) \ x(5) \ x(6)] \\ y(6) = b(1)*z(3) + b(2)*z(2) + b(3)*z(1)$$

...

You can implement the FIR filter using a circular buffer like the following MATLAB function.

```
function [y,z,p] = fir_filt_circ_buff_original(b,x,z,p)
    y = zeros(size(x));
    nx = length(x);
    nb = length(b);
    for n=1:nx
        p=p+1; if p>nb, p=1; end
        z(p) = x(n);
        acc = 0;
        k = p;
        for j=1:nb
            acc = acc + b(j)*z(k);
            k=k-1; if k<1, k=nb; end
        end
        y(n) = acc;
    end
end
```

Create a Test File

Create a test file to validate that the floating-point algorithm works as expected before converting it to fixed point. You can use the same test file to propose fixed-point data types and to compare fixed-point results to the floating-point baseline after the conversion.

The test vectors should represent realistic inputs that exercise the full range of values expected by your system. Realistic inputs are impulses, sums of sinusoids, and chirp signals, for which you can verify that the outputs are correct using linear theory. Signals that produce maximum output are useful for verifying that your system does not overflow.

Setup

Run the following code to capture and reset the current state of global fixed-point math settings and fixed-point preferences.

```
resetglobalfimath;  
FIPREF_STATE = get(fipref);  
resetfipref;
```

Filter Coefficients

Use the following low-pass filter coefficients that were computed using the `fir1` function from Signal Processing Toolbox™.

```
b = fir1(11,0.25);  
  
b = [-0.004465461051254  
     -0.004324228005260  
     +0.012676739550326  
     +0.074351188907780  
     +0.172173206073645  
     +0.249588554524763  
     +0.249588554524763  
     +0.172173206073645  
     +0.074351188907780  
     +0.012676739550326  
     -0.004324228005260  
     -0.004465461051254]';
```

Time Vector

Use this time vector to create the test signals.

```
nx = 256;  
t = linspace(0,10*pi,nx)';
```

Impulse Input

The response of an FIR filter to an impulse input is the filter coefficients themselves.

```
x_impulse = zeros(nx,1); x_impulse(1) = 1;
```

Signal that Produces the Maximum Output

The maximum output of a filter occurs when the signs of the inputs line up with the signs of the filter's impulse response.

```
x_max_output = sign(fliplr(b))';  
x_max_output = repmat(x_max_output,ceil(nx/length(b)),1);  
x_max_output = x_max_output(1:nx);
```

The maximum magnitude of the output is the 1-norm of its impulse response, which is `norm(b,1) = sum(abs(b))`.

```
maximum_output_magnitude = norm(b,1) %#ok<*NOPTS>
```

```
maximum_output_magnitude = 1.0352
```

Sum of Sinusoids

A sum of sinusoids is a typical input for a filter. You can easily see the high frequencies filtered out in the plot.

```
f0 = 0.1; f1 = 2;
x_sines = sin(2*pi*t*f0) + 0.1*sin(2*pi*t*f1);
```

Chirp

A chirp gives a good visual of the low-pass filter action of passing the low frequencies and attenuating the high frequencies.

```
f_chirp = 1/16; % Target frequency
x_chirp = sin(pi*f_chirp*t.^2); % Linear chirp

titles = {'Impulse', 'Max output', 'Sum of sines', 'Chirp'};
x = [x_impulse, x_max_output, x_sines, x_chirp];
```

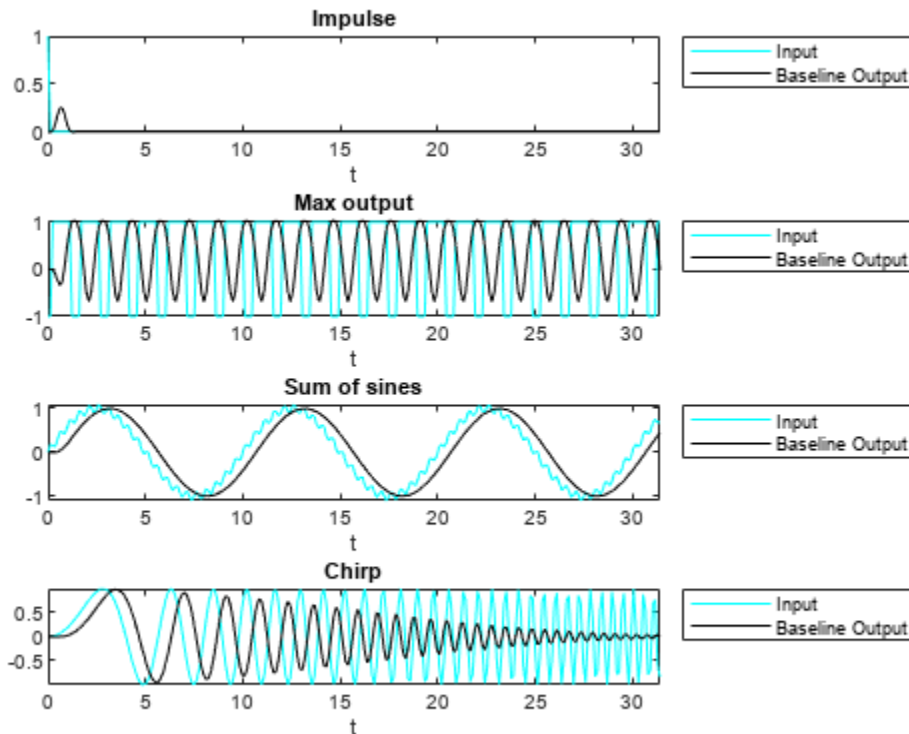
Call Original Function

Before starting the conversion to fixed point, call your original function with the test file inputs to establish a baseline to compare to subsequent outputs.

```
y0 = zeros(size(x));
for i=1:size(x,2)
    % Initialize the states for each column of input
    p = 0;
    z = zeros(size(b));
    y0(:,i) = fir_filt_circ_buff_original(b,x(:,i),z,p);
end
```

Baseline Output

```
fir_filt_circ_buff_plot(1,titles,t,x,y0)
```



Prepare for Instrumentation and Code Generation

The first step after the algorithm works in MATLAB is to prepare it for instrumentation, which requires code generation. Before the conversion, you can use the `coder.screener` function to analyze your code and identify unsupported functions and language features.

Entry-Point Function

When doing instrumentation and code generation, it is convenient to have an entry-point function that calls the function to be converted to fixed point. You can cast the FIR filter's inputs to different data types, and add calls to different variations of the filter for comparison. By using an entry-point function you can run both fixed-point and floating-point variants of your filter, and also different variants of fixed point. This allows you to iterate on your code more quickly to arrive at the optimal fixed-point design.

```
function y = fir_filt_circ_buff_original_entry_point(b,x,reset)
    if nargin<3, reset = true; end
    % Define the circular buffer z and buffer position index p.
    % They are declared persistent so the filter can be called in a streaming
    % loop, each section picking up where the last section left off.
    persistent z p
    if isempty(z) || reset
        p = 0;
        z = zeros(size(b));
    end
    [y,z,p] = fir_filt_circ_buff_original(b,x,z,p);
end
```

Test File

Your test file calls the compiled entry-point function.

```
function y = fir_filt_circ_buff_test(b,x)
    y = zeros(size(x));
    for i=1:size(x,2)
        reset = true;
        y(:,i) = fir_filt_circ_buff_original_entry_point_mex(b,x(:,i),reset);
    end
end
```

Build Original Function

Compile the original entry-point function with `buildInstrumentedMex`. This instruments your code for logging so you can collect minimum and maximum values from the simulation and get proposed data types.

```
reset = true;
buildInstrumentedMex fir_filt_circ_buff_original_entry_point -args {b, x(:,1), reset}
```

Run Original Function

Run your test file inputs through the algorithm to log minimum and maximum values.

```
y1 = fir_filt_circ_buff_test(b,x);
```

Show Types

Use `showInstrumentationResults` to view the data types of all your variables and the minimum and maximum values that were logged during the test file run. Look at the maximum value logged for the output variable `y` and accumulator variable `acc` and note that they attained the theoretical maximum output value that you calculated previously.

```
showInstrumentationResults fir_filt_circ_buff_original_entry_point_mex
```

To see these results in the instrumented **Code Generation Report**:

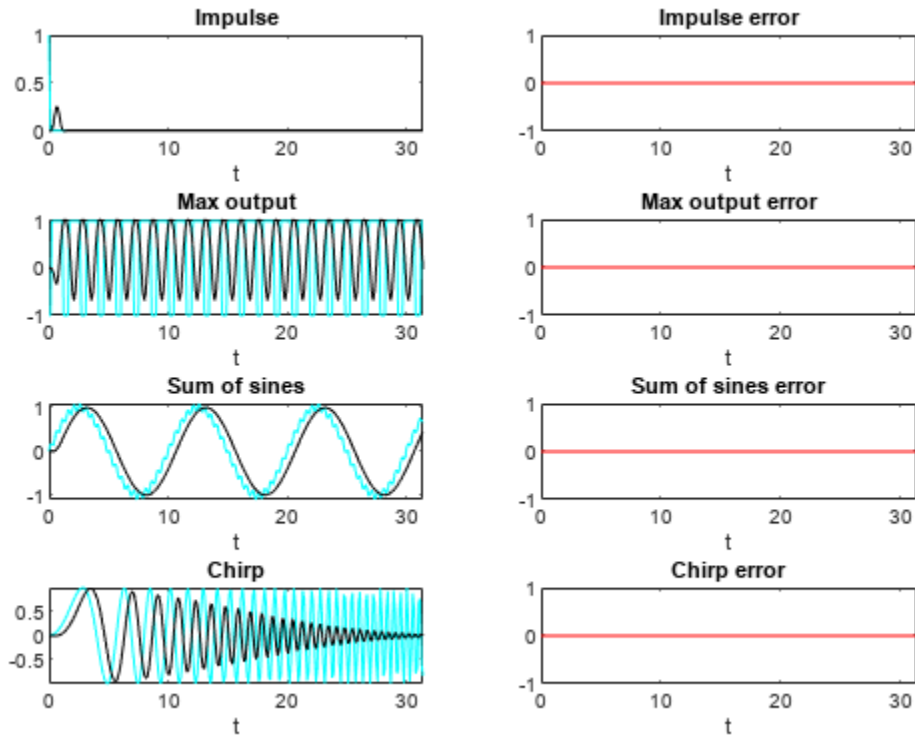
- Select function `fir_filt_circ_buff_original`
- Select the **Variables** tab

| Variable | Type | Size | Class | Complex | Always Whole Number | SimMin | SimMax |
|----------|--------|---------|--------|---------|---------------------|---------------------|--------------------|
| acc | Local | 1 x 1 | double | No | No | -1.0045281391112986 | 1.035158756226056 |
| b | Input | 1 x 12 | double | No | No | -0.004465461051254 | 0.249588554524763 |
| j | Local | 1 x 1 | double | No | Yes | 1 | 12 |
| k | Local | 1 x 1 | double | No | Yes | 0 | 12 |
| n | Local | 1 x 1 | double | No | Yes | 1 | 256 |
| nb | Local | 1 x 1 | double | No | Yes | 12 | 12 |
| nx | Local | 1 x 1 | double | No | Yes | 256 | 256 |
| p | I/O | 1 x 1 | double | No | Yes | 0 | 13 |
| x | Input | 256 x 1 | double | No | No | -1.0966086573451541 | 1.0874250377226369 |
| y | Output | 256 x 1 | double | No | No | -0.9959923266057762 | 1.035158756226056 |
| z | I/O | 1 x 12 | double | No | No | -1.0966086573451541 | 1.0874250377226369 |

Validate Original Function

Every time you modify your function, validate that the results still match your baseline.

```
fir_filt_circ_buff_plot2(2,titles,t,x,y0,y1)
```



Convert Functions to Use Types Tables

To separate data types from the algorithm, you:

- 1 Create a table of data type definitions.
- 2 Modify the algorithm code to use data types from that table.

This example shows the iterative steps by creating different files. In practice, you can make the iterative changes to the same file.

Original Types Table

Create a types table using a structure with prototypes for the variables set to their original types. Use the baseline types to validate that you made the initial conversion correctly, and also use it to programmatically toggle your function between floating-point and fixed-point types. The index variables j , k , n , nb , nx are automatically converted to integers by MATLAB Coder™, so you don't need to specify their types in the table.

Specify the prototype values as empty (`[]`) since the data types are used, but not the values.

```
function T = fir_filt_circ_buff_original_types()
    T.acc=double([]);
    T.b=double([]);
    T.p=double([]);
    T.x=double([]);
    T.y=double([]);
    T.z=double([]);
end
```

Type-Aware Filter Function

Prepare the filter function and entry-point function to be type-aware by using the `cast` and `zeros` functions and the types table.

Use subscripted assignment `acc(:)=...`, `p(:)=1`, and `k(:)=nb` to preserve data types during assignment. See the “Cast fi Objects” on page 2-10 for more details about subscripted assignment and preserving data types.

The function call `y = zeros(size(x), 'like', T.y)` creates an array of zeros the same size as `x` with the properties of variable `T.y`. Initially, `T.y` is a double defined in function `fir_filt_circ_buff_original_types`, but it is re-defined as a fixed-point type later in this example.

The function call `acc = cast(0, 'like', T.acc)` casts the value `0` with the same properties as variable `T.acc`. Initially, `T.acc` is a double defined in function `fir_filt_circ_buff_original_types`, but it is re-defined as a fixed-point type later in this example.

```
function [y,z,p] = fir_filt_circ_buff_typed(b,x,z,p,T)
    y = zeros(size(x), 'like', T.y);
    nx = length(x);
    nb = length(b);
    for n=1:nx
        p(:)=p+1; if p>nb, p(:)=1; end
        z(p) = x(n);
        acc = cast(0, 'like', T.acc);
```

```

        k = p;
        for j=1:nb
            acc(:) = acc + b(j)*z(k);
            k(:)=k-1; if k<1, k(:)=nb; end
        end
        y(n) = acc;
    end
end

```

Type-Aware Entry-Point Function

The function call `p1 = cast(0, 'like', T1.p)` casts the value 0 with the same properties as variable `T1.p`. Initially, `T1.p` is a double defined in function `fir_filt_circ_buff_original_types`, but it is re-defined as an integer type later in this example.

The function call `z1 = zeros(size(b), 'like', T1.z)` creates an array of zeros the same size as `b` with the properties of variable `T1.z`. Initially, `T1.z` is a double defined in function `fir_filt_circ_buff_original_types`, but it is re-defined as a fixed-point type later in this example.

```

function y1 = fir_filt_circ_buff_typed_entry_point(b,x,reset)
    if nargin<3, reset = true; end
    %
    % Baseline types
    %
    T1 = fir_filt_circ_buff_original_types();
    % Each call to the filter needs to maintain its own states.
    persistent z1 p1
    if isempty(z1) || reset
        p1 = cast(0, 'like', T1.p);
        z1 = zeros(size(b), 'like', T1.z);
    end
    b1 = cast(b, 'like', T1.b);
    x1 = cast(x, 'like', T1.x);
    [y1,z1,p1] = fir_filt_circ_buff_typed(b1,x1,z1,p1,T1);
end

```

Validate Modified Function

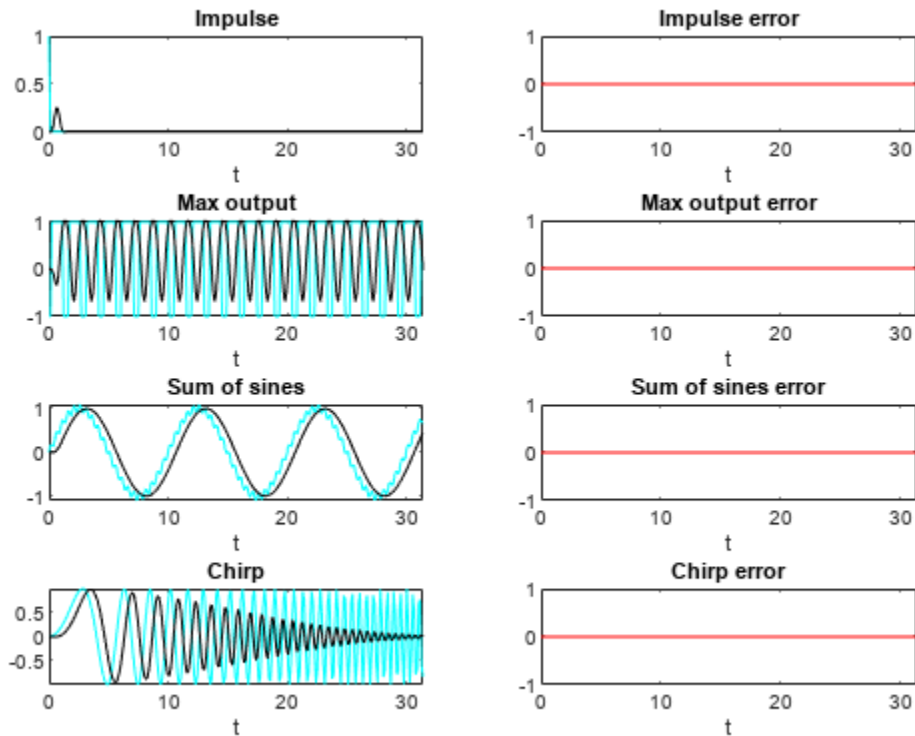
Every time you modify your function, validate that the results still match your baseline. Since you used the original types in the types table, the outputs should be identical. This validates that you made the conversion to separate the types from the algorithm correctly.

```

buildInstrumentedMex fir_filt_circ_buff_typed_entry_point -args {b, x(:,1), reset}

y1 = fir_filt_circ_buff_typed_test(b,x);
fir_filt_circ_buff_plot2(3,titles,t,x,y0,y1)

```



Propose Data Types from Simulation Min/Max Logs

Use the `showInstrumentationResults` function to propose fixed-point fraction lengths, given a default signed fixed-point type and 16-bit word length.

```
showInstrumentationResults fir_filt_circ_buff_original_entry_point_mex ...
    -defaultDT numerictype(1,16) -proposeFL
```

In the instrumented **Code Generation Report**, select function `fir_filt_circ_buff_original` and the **Variables** tab to see these results.

| Variable | Type | Size | Class | Complex | Proposed Signedness | Proposed WL | Proposed FL | Always Whole Number | SimMin | SimMax |
|----------|--------|---------|--------|---------|---------------------|-------------|-------------|---------------------|---------------------|--------------------|
| acc | Local | 1 x 1 | double | No | Signed | 16 | 14 | No | -1.0045281391112986 | 1.035158756226056 |
| b | Input | 1 x 12 | double | No | Signed | 16 | 17 | No | -0.004465461051254 | 0.249588554524763 |
| j | Local | 1 x 1 | double | No | Signed | 16 | 0 | Yes | 1 | 12 |
| k | Local | 1 x 1 | double | No | Signed | 16 | 0 | Yes | 0 | 12 |
| n | Local | 1 x 1 | double | No | Signed | 16 | 0 | Yes | 1 | 256 |
| nb | Local | 1 x 1 | double | No | Signed | 16 | 0 | Yes | 12 | 12 |
| nx | Local | 1 x 1 | double | No | Signed | 16 | 0 | Yes | 256 | 256 |
| p | I/O | 1 x 1 | double | No | Signed | 16 | 0 | Yes | 0 | 13 |
| x | Input | 256 x 1 | double | No | Signed | 16 | 14 | No | -1.0966086573451541 | 1.0874250377226369 |
| y | Output | 256 x 1 | double | No | Signed | 16 | 14 | No | -0.9959923266057762 | 1.035158756226056 |
| z | I/O | 1 x 12 | double | No | Signed | 16 | 14 | No | -1.0966086573451541 | 1.0874250377226369 |

Create a Fixed-Point Types Table

Use the proposed types from the **Code Generation Report** to guide you in choosing fixed-point types and create a fixed-point types table using a structure with prototypes for the variables.

Use your knowledge of the algorithm to improve on the proposals. For example, you are using the `acc` variable as an accumulator, so make it 32-bits. From the **Code Generation Report**, you can see that `acc` needs at least 2 integer bits to prevent overflow, so set the fraction length to 30.

Variable `p` is used as an index, so you can make it a builtin 16-bit integer.

Specify the prototype values as empty (`[]`) since the data types are used, but not the values.

```
function T = fir_filt_circ_buff_fixed_point_types()
    T.acc=fi([],true,32,30);
    T.b=fi([],true,16,17);
    T.p=int16([]);
    T.x=fi([],true,16,14);
    T.y=fi([],true,16,14);
    T.z=fi([],true,16,14);
end
```

Add Fixed Point to Entry-Point Function

Add a call to the fixed-point types table in the entry-point function:

```
T2 = fir_filt_circ_buff_fixed_point_types();
persistent z2 p2
if isempty(z2) || reset
    p2 = cast(0,'like',T2.p);
    z2 = zeros(size(b),'like',T2.z);
end
b2 = cast(b,'like',T2.b);
x2 = cast(x,'like',T2.x);
[y2,z2,p2] = fir_filt_circ_buff_typed(b2,x2,z2,p2,T2);
```

Build and Run Algorithm with Fixed-Point Data Types

```
buildInstrumentedMex fir_filt_circ_buff_typed_entry_point -args {b, x(:,1), reset}

[y1,y2] = fir_filt_circ_buff_typed_test(b,x);

showInstrumentationResults fir_filt_circ_buff_typed_entry_point_mex
```

To see these results in the instrumented **Code Generation Report**:

- Select the entry-point function, `fir_filt_circ_buff_typed_entry_point`
- Select `fir_filt_circ_buff_typed` in the following line of code:

```
[y2,z2,p2] = fir_filt_circ_buff_typed(b2,x2,z2,p2,T2);
```

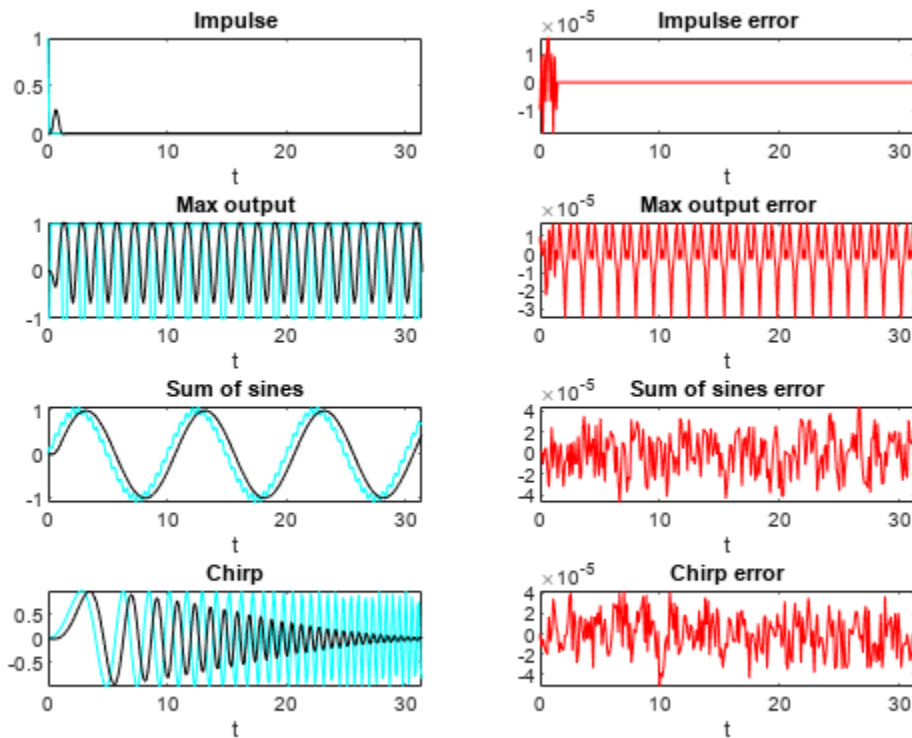
- Select the **Variables** tab

| Variable | Type | Size | Class | Complex | Signedness | WL | FL | Percent of Current Range | Always Whole Number | SimMin | SimMax |
|----------|--------|---------|-------------|---------|------------|----|----|--------------------------|---------------------|----------------------|--------------------|
| acc | Local | 1 x 1 | embedded.fi | No | Signed | 32 | 30 | 52 | No | -1.004551738500595 | 1.03515625 |
| b | Input | 1 x 12 | embedded.fi | No | Signed | 16 | 17 | 100 | No | -0.00446319580078125 | 0.2495880126953125 |
| j | Local | 1 x 1 | double | No | - | - | - | - | Yes | 1 | 12 |
| k | Local | 1 x 1 | int16 | No | - | - | - | - | Yes | 0 | 12 |
| n | Local | 1 x 1 | double | No | - | - | - | - | Yes | 1 | 256 |
| nb | Local | 1 x 1 | double | No | - | - | - | - | Yes | 12 | 12 |
| nx | Local | 1 x 1 | double | No | - | - | - | - | Yes | 256 | 256 |
| p | I/O | 1 x 1 | int16 | No | - | - | - | - | Yes | 0 | 13 |
| T | Input | 1 x 1 | struct | - | - | - | - | - | - | - | - |
| x | Input | 256 x 1 | embedded.fi | No | Signed | 16 | 14 | 55 | No | -1.09661865234375 | 1.08740234375 |
| y | Output | 256 x 1 | embedded.fi | No | Signed | 16 | 14 | 52 | No | -0.9959716796875 | 1.03515625 |
| z | I/O | 1 x 12 | embedded.fi | No | Signed | 16 | 14 | 55 | No | -1.09661865234375 | 1.08740234375 |

16-bit Word Length, Full Precision Math

Validate that the results are within an acceptable tolerance of your baseline.

```
fir_filt_circ_buff_plot2(4,titles,t,x,y1,y2);
```



Your algorithm has now been converted to fixed-point MATLAB code. If you also want to convert to C-code, then proceed to the next section.

Generate C-Code

This section describes how to generate efficient C-code from the fixed-point MATLAB code from the previous section.

Required Products

You need MATLAB Coder to generate C-code, and you need Embedded Coder® for the hardware implementation settings used in this example.

Algorithm Tuned for Most-Efficient C-Code

The output variable `y` is initialized to zeros, and then completely overwritten before it is used. Therefore, filling `y` with all zeros is unnecessary. You can use the `coder.nullcopy` function to declare a variable without actually filling it with values, which makes the code in this case more efficient. However, you have to be very careful when using `coder.nullcopy` because if you access an element of a variable before it is assigned, then you are accessing uninitialized memory and its contents are unpredictable.

A rule of thumb for when to use `coder.nullcopy` is when the initialization takes significant time compared to the rest of the algorithm. If you are not sure, then the safest thing to do is to not use it.

```
function [y,z,p] = fir_filt_circ_buff_typed_codegen(b,x,z,p,T)
    % Use coder.nullcopy only when you are certain that every value of
    % the variable is overwritten before it is used.
    y = coder.nullcopy(zeros(size(x),'like',T.y));
    nx = length(x);
    nb = length(b);
    for n=1:nx
        p(:)=p+1; if p>nb, p(:)=1; end
        z(p) = x(n);
        acc = cast(0,'like',T.acc);
        k = p;
        for j=1:nb
            acc(:) = acc + b(j)*z(k);
            k(:)=k-1; if k<1, k(:)=nb; end
        end
        y(n) = acc;
    end
end
```

Native C-Code Types

You can set the fixed-point math properties to match the native actions of C. This generates the most efficient C-code, but this example shows that it can create problems with overflow and produce less accurate results which are corrected in the next section. It doesn't always create problems, though, so it is worth trying first to see if you can get the cleanest possible C-code.

Set the fixed-point math properties to use floor rounding and wrap overflow because those are the default actions in C.

Set the fixed-point math properties of products and sums to match native C 32-bit integer types, and to keep the least significant bits (LSBs) of math operations.

Add these settings to a fixed-point types table.

```

function T = fir_filt_circ_buff_dsp_types()
    F = fimath('RoundingMethod','Floor',...
             'OverflowAction','Wrap',...
             'ProductMode','KeepLSB',...
             'ProductWordLength',32,...
             'SumMode','KeepLSB',...
             'SumWordLength',32);
    T.acc=fi([],true,32,30,F);
    T.p=int16([]);
    T.b=fi([],true,16,17,F);
    T.x=fi([],true,16,14,F);
    T.y=fi([],true,16,14,F);
    T.z=fi([],true,16,14,F);
end

```

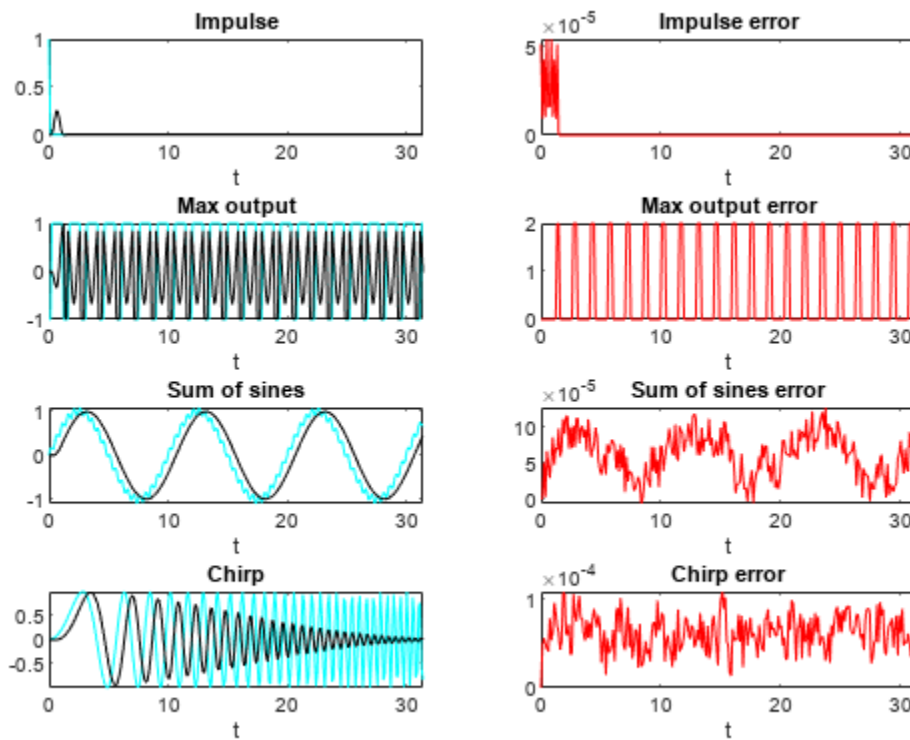
Test the Native C-Code Types

Add a call to the types table in the entry-point function and run the test file.

```
[y1,y2,y3] = fir_filt_circ_buff_typed_test(b,x); %#ok<*ASGLU>
```

In the second row of plots, you can see that the maximum output error is twice the size of the input, indicating that a value that should have been positive overflowed to negative. You can also see that the other outputs did not overflow. This is why it is important to have your test file exercise the full range of values in addition to other typical inputs.

```
fir_filt_circ_buff_plot2(5,titles,t,x,y1,y3);
```



Use Scaled Double Types to Find Overflows

Scaled double variables store their data in double-precision floating-point, so they carry out arithmetic in full range. They also retain their fixed-point settings, so they are able to report when a computation goes out of the range of the fixed-point type.

Change the data types to scaled double, and add these settings to a scaled-double types table.

```
function T = fir_filt_circ_buff_scaled_double_types()
    F = fimath('RoundingMethod','Floor',...
             'OverflowAction','Wrap',...
             'ProductMode','KeepLSB',...
             'ProductWordLength',32,...
             'SumMode','KeepLSB',...
             'SumWordLength',32);
    DT = 'ScaledDouble';
    T.acc=fi([],true,32,30,F,'DataType',DT);
    T.p=int16([]);
    T.b=fi([],true,16,17,F,'DataType',DT);
    T.x=fi([],true,16,14,F,'DataType',DT);
    T.y=fi([],true,16,14,F,'DataType',DT);
    T.z=fi([],true,16,14,F,'DataType',DT);
end
```

Add a call to the scaled-double types table to the entry-point function and run the test file.

```
[y1,y2,y3,y4] = fir_filt_circ_buff_typed_test(b,x); %#ok<*NASGU>
```

Show the instrumentation results with the scaled-double types.

```
showInstrumentationResults fir_filt_circ_buff_typed_entry_point_mex
```

To see these results in the instrumented **Code Generation Report**:

- Select the entry-point function, `fir_filt_circ_buff_typed_entry_point`
- Select `fir_filt_circ_buff_typed_codegen` in the following line of code:

```
[y4,z4,p4] = fir_filt_circ_buff_typed_codegen(b4,x4,z4,p4,T4);
```

- Select the **Variables** tab.
- Look at the variables in the table. None of the variables overflowed, which indicates that the overflow occurred as the result of an operation.
- Hover over the operators in the report (+, -, *, =).
- Hover over the + in this line of MATLAB code in the instrumented **Code Generation Report**:

```
acc(:) = acc + b(j)*z(k);
```

The report shows that the sum overflowed:


```

% Use coder.nullcopy only when you are certain that every value
% the variable will be overwritten before it is used.
y = coder.nullcopy(zeros(size(x), 'like', T.y));
nx = length(x);
nb = length(b);
for n=1:nx
    p(:)=p+1; if p>nb, p(:)=1; end
    z(p) = x(n);
    acc = cast(0, 'like', T.acc);
    k = p;
    for j=1:nb
        acc(:) = acc + b(j)*z(k);
        k(:)=k-1; if k
    end
    y(n) = acc;
end

```

| Order | Variable | Type | Direction |
|-------|----------|--------|-----------|
| | y | Output | Output |
| | z | Input | Input |
| | p | Input | Input |

| Information for the selected expression: | |
|--|---------------------|
| Size | 1 x 1 |
| Class | embedded.fi |
| Complex | No |
| DT Mode | ScaledDouble |
| Signedness | Signed |
| WL | 32 |
| FL | 31 |
| Percent of Current Range | 104 |
| Always Whole Number | No |
| SimMin | -1.0045281391112986 |
| SimMax | 1.035158756226056 |

The reason the sum overflowed is that a full-precision product for $b(j)*z(k)$ produces a `numerictype(true, 32, 31)` because b has `numerictype(true, 16, 17)` and z has `numerictype(true, 16, 14)`. The sum type is set to "keep least significant bits" (KeepLSB), so the sum has `numerictype(true, 32, 31)`. However, 2 integer bits are necessary to store the minimum and maximum simulated values of -1.0045 and $+1.035$, respectively.

Adjust to Avoid the Overflow

Set the fraction length of b to 16 instead of 17 so that $b(j)*z(k)$ is `numerictype(true, 32, 30)`, and so the sum is also `numerictype(true, 32, 30)` following the KeepLSB rule for sums.

Leave all other settings the same, and set

```
T.b=fi([], true, 16, 16, F);
```

Then the sum in this line of MATLAB code no longer overflows:

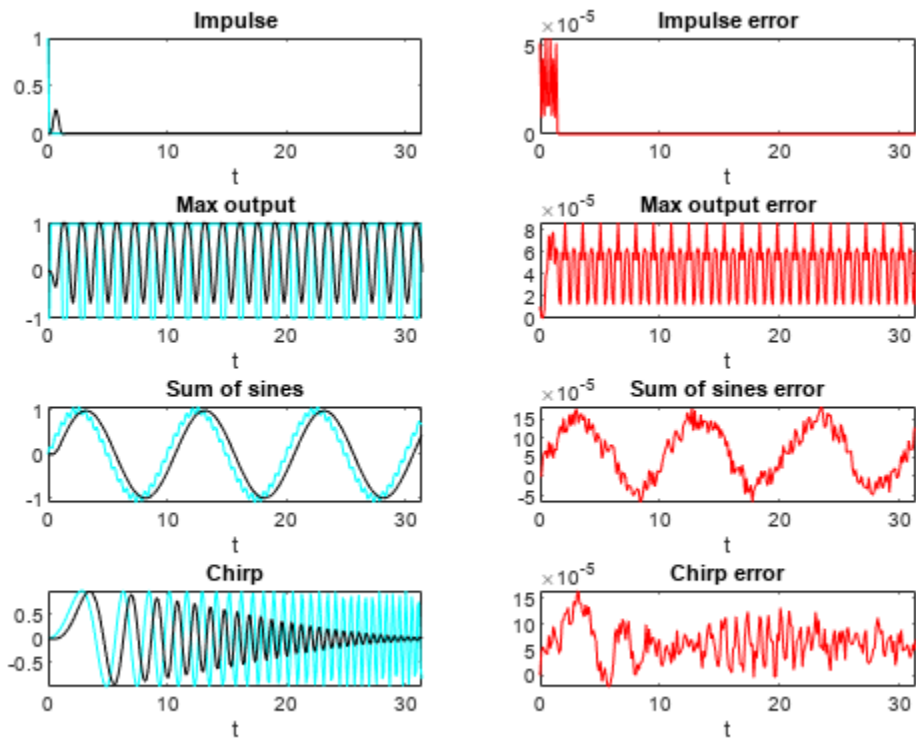
```
acc(:) = acc + b(j)*z(k);
```

Run the test file with the new settings and plot the results.

```
[y1,y2,y3,y4,y5] = fir_filt_circ_buff_typed_test(b,x);
```

You can see that the overflow has been avoided. However, the plots show a bias and a larger error due to using C's natural floor rounding. If this bias is acceptable to you, then you can stop here and the generated C-code is very clean.

```
fir_filt_circ_buff_plot2(6,titles,t,x,y1,y5);
```



Eliminate the Bias

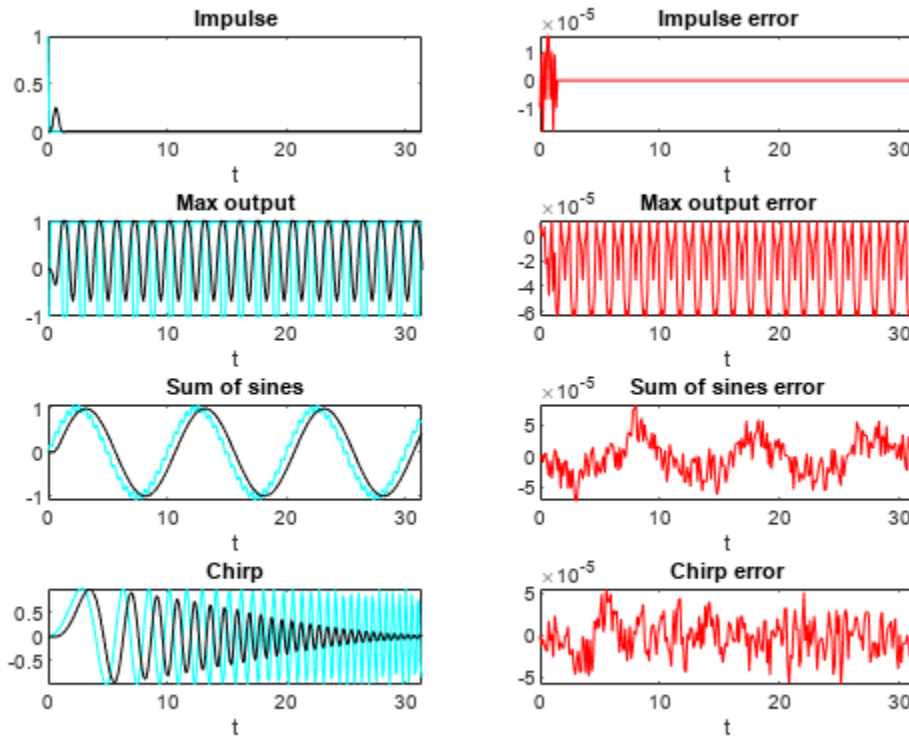
If the bias is not acceptable in your application, then change the rounding method to 'Nearest' to eliminate the bias. Rounding to nearest generates slightly more complicated C-code, but it may be necessary for you if you want to eliminate the bias and have a smaller error.

The final fixed-point types table with nearest rounding and adjusted coefficient fraction length is:

```
function T = fir_filt_circ_buff_dsp_nearest_types()
    F = fimath('RoundingMethod','Nearest',...
             'OverflowAction','Wrap',...
             'ProductMode','KeepLSB',...
             'ProductWordLength',32,...
             'SumMode','KeepLSB',...
             'SumWordLength',32);
    T.acc=fi([],true,32,30,F);
    T.p=int16([]);
    T.b=fi([],true,16,16,F);
    T.x=fi([],true,16,14,F);
    T.y=fi([],true,16,14,F);
    T.z=fi([],true,16,14,F);
end
```

Call this types table from the entry-point function and run and plot the output.

```
[y1,y2,y3,y4,y5,y6] = fir_filt_circ_buff_typed_test(b,x);
fir_filt_circ_buff_plot2(7,titles,t,x,y1,y6);
```



Run Code Generation Command

Run this build function to generate C-code. It is a best practice to create a build function so you can generate C-code for your core algorithm without the entry-point function or test file so the C-code for the core algorithm can be included in a larger project.

```
function fir_filt_circ_buff_build_function()
%
% Declare input arguments
%
T = fir_filt_circ_buff_dsp_nearest_types();
b = zeros(1,12,'like',T.b);
x = zeros(256,1,'like',T.x);
z = zeros(size(b),'like',T.z);
p = cast(0,'like',T.p);
%
% Code generation configuration
%
h = coder.config('lib');
h.PurelyIntegerCode = true;
h.SaturateOnIntegerOverflow = false;
h.SupportNonFinite = false;
h.HardwareImplementation.ProdBitPerShort = 8;
h.HardwareImplementation.ProdBitPerInt = 16;
```

```

h.HardwareImplementation.ProdBitPerLong = 32;
%
% Generate C-code
%
codegen fir_filt_circ_buff_typed_codegen -args {b,x,z,p,T} -config h -launchreport
end

```

Generated C-Code

Using these settings, MATLAB Coder generates the following C-code:

```

void fir_filt_circ_buff_typed_codegen(const int16_T b[12], const int16_T x[256],
  int16_T z[12], int16_T *p, int16_T y[256])
{
  int16_T n;
  int32_T acc;
  int16_T k;
  int16_T j;
  for (n = 0; n < 256; n++) {
    (*p)++;
    if (*p > 12) {
      *p = 1;
    }
    z[*p - 1] = x[n];
    acc = 0L;
    k = *p;
    for (j = 0; j < 12; j++) {
      acc += (int32_T)b[j] * z[k - 1];
      k--;
      if (k < 1) {
        k = 12;
      }
    }
    y[n] = (int16_T)((acc >> 16) + ((acc & 32768L) != 0L));
  }
}

```

Run the following code to restore the global states.

```

fipref(FIPREF_STATE);
clearInstrumentationResults fir_filt_circ_buff_original_entry_point_mex
clearInstrumentationResults fir_filt_circ_buff_typed_entry_point_mex
clear fir_filt_circ_buff_original_entry_point_mex
clear fir_filt_circ_buff_typed_entry_point_mex

```

References:

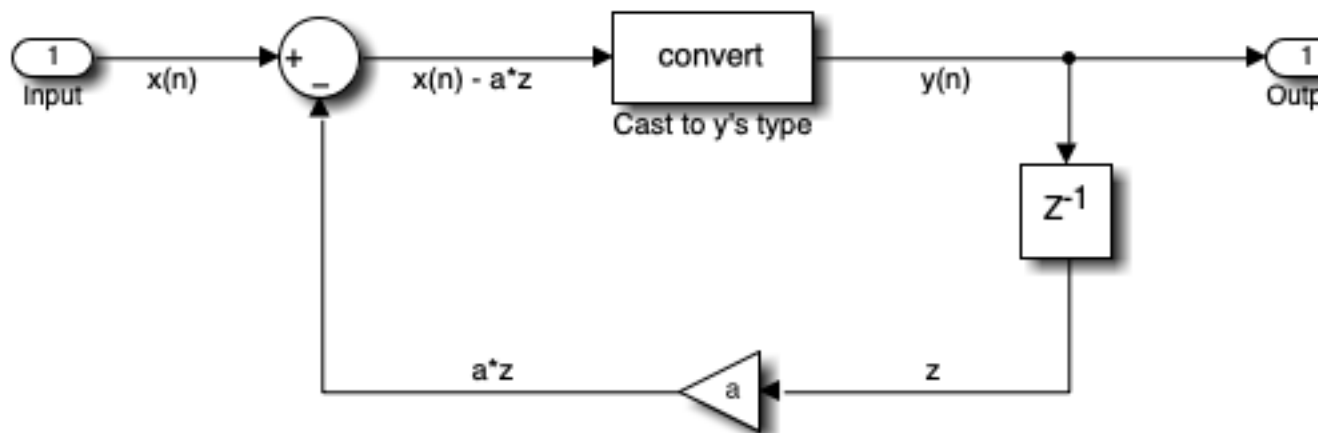
[1] J. R. Ragazzini and L. A. Zadeh. "The analysis of sampled-data systems". In: Transactions of the American Institute of Electrical Engineers 71 (II) (1952), pp. 225-234.

Accelerate Fixed-Point Simulation

This example shows how to accelerate fixed-point algorithms using the `fiaccel` function. Code acceleration provides optimizations for accelerating fixed-point algorithms through MEX file generation. Fixed-Point Designer™ provides a convenience function `fiaccel` to convert your MATLAB code to a MEX function, which can greatly accelerate the execution speed of your fixed-point algorithms. In this example, you generate a MEX function from MATLAB® code, run the generated MEX function, and compare the execution speed with MATLAB code simulation.

Description of the Example

This example uses a first-order feedback loop. Casting to the output-signal type prevents infinite bit growth. The output signal is delayed by one sample and fed back to dampen the input signal.



Inspect the MATLAB® Feedback Function Code

The MATLAB function that performs the feedback loop is in the file `fiaccelFeedback.m`. Subscripted assignment into the output `y` casts to `y`'s type and prevents infinite bit growth.

```
function [y,z] = fiaccelFeedback(x,a,y,z)
    for n = 1:length(x)
        y(n) = x(n) - a*z;
        z(:) = y(n);
    end
end
```

The following variables are used in this function:

- `x` is the input signal vector.
- `y` is the output signal vector.
- `a` is the feedback gain.
- `z` is the unit-delayed output signal.

Create the Input Signal and Initialize Variables

```
clearvars
```

Put the settings of the random number generator to its default value.

```
rng('default');
```

Input signal.

```
x = fi(2*rand(1000,1)-1,true,16,15);
```

Feedback gain.

```
a = fi(0.9,true,16,15);
```

Initialize output. Fraction length is chosen to prevent overflow.

```
y = fi(zeros(size(x)),true,16,12);
```

Initialize delayed output.

```
z = cast(0,'like',y);
```

Run Interpreted MATLAB and Time

```
tic  
y1 = fiaccelFeedback(x,a,y,z);  
t1 = toc;
```

Build the MEX Version of the Feedback Code

Declare feedback gain parameter a constant for code generation.

```
fiaccel fiaccelFeedback -args {x,coder.Constant(a),y,z} -o fiaccelFeedback_mex
```

Run the MEX Version and Time

Run once to load the MEX file in memory.

```
fiaccelFeedback_mex(x,a,y1,z);
```

Run again to time.

```
tic  
y2 = fiaccelFeedback_mex(x,a,y,z);  
t2 = toc;
```

Acceleration Ratio

Compare the MEX execution speed with MATLAB code simulation.

```
ratio_of_speed_up = t1/t2
```

```
ratio_of_speed_up =  
245.3719
```

Verify that Fixed-Point Interpreted MATLAB and MEX Outputs are Identical

```
isequal(y1,y2)
```

```
ans =
```

```
logical
```

```
1
```

Suppress Code Analyzer warnings.

```
 %#ok<*NOPTS>
```

See Also

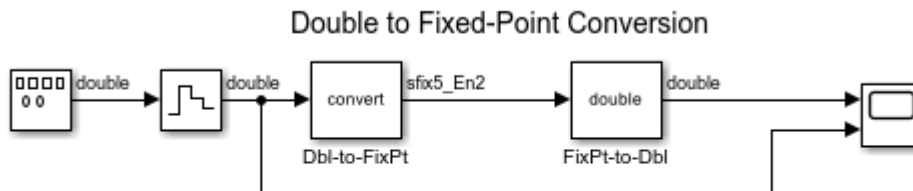
```
fiaccel
```

Double to Fixed-Point Conversion

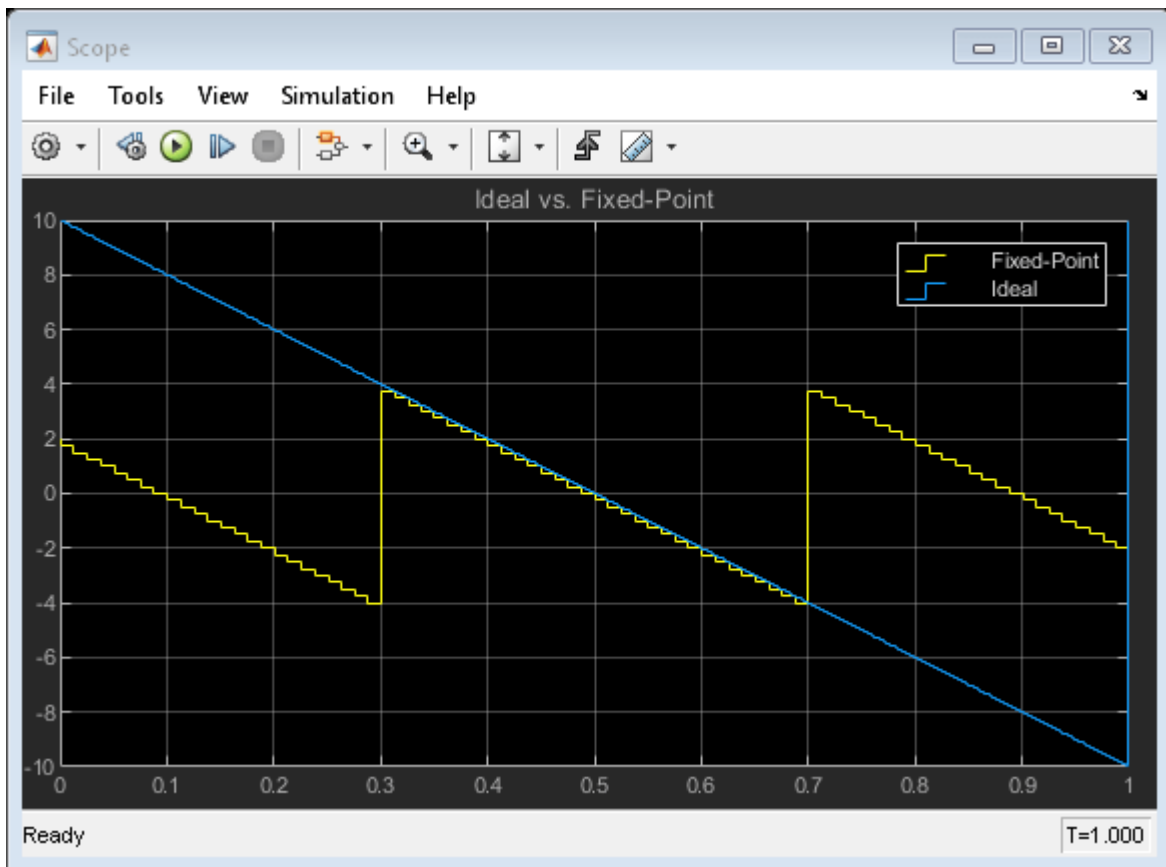
This model shows how to convert signals between built-in and fixed-point data types and illustrates how fixed-point data types affect the representable precision and range. The fixed-point data type used in this model is `fixdt(1,5,2)`, which is a signed 5-bit number with 2 bits to the right of the binary point:

$$\begin{aligned} \text{Precision} &= (2^{-2}) = 0.25 \\ \text{Representable minimum} &= -(2^{-2}) * (2^4) = -4 \\ \text{Representable maximum} &= (2^{-2}) * (2^4 - 1) = 3.75 \end{aligned}$$

Open the Data Type Conversion block `Dbl-to-FixPt` to modify the attributes of the fixed-point data type and see how they impact the range and precision of the resulting fixed-point signal.



Copyright 1990-2009 The MathWorks, Inc.

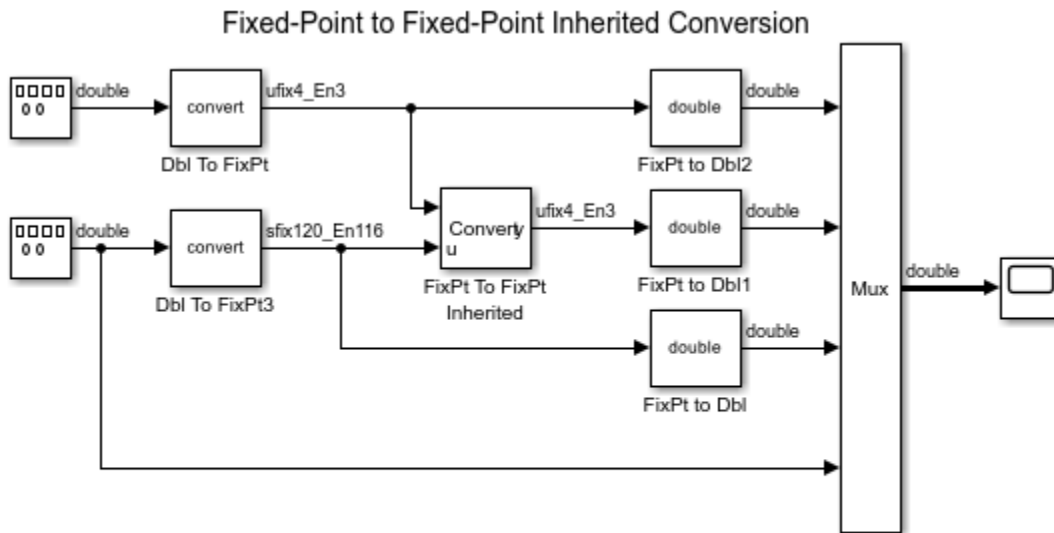


Tips for specifying fixed-point data types:

- Click the >> button to open the Data Type Assistant.
- You can specify the output minimum/maximum and use these values to calculate the best-precision scaling, which maximizes the precision while covering the specified range.
- Once you are familiar with the syntax for `fixdt`, you can enter the expression directly into the data type parameter without using the Data Type Assistant.

Fixed-Point to Fixed-Point Inherited Conversion

This example shows how to use the Data Type Conversion block `FixPt To FixPt Inherited`. Because Simulink® propagates data types throughout a block diagram, fixed-point utility modeling can be templated for multiple use scenarios. The data type inheritance capability can be given additional information through the use of various relative-type specifications. The `FixPt To FixPt Inherited` block provides a way to specify that "this signal should have the same fixed-point data type as that signal" by wiring it up to the signals of interest.



NOTE: Bit width of 120 works for Simulation. It is unlikely to work with Simulink Coder because few compilers currently support more than 64 bits.

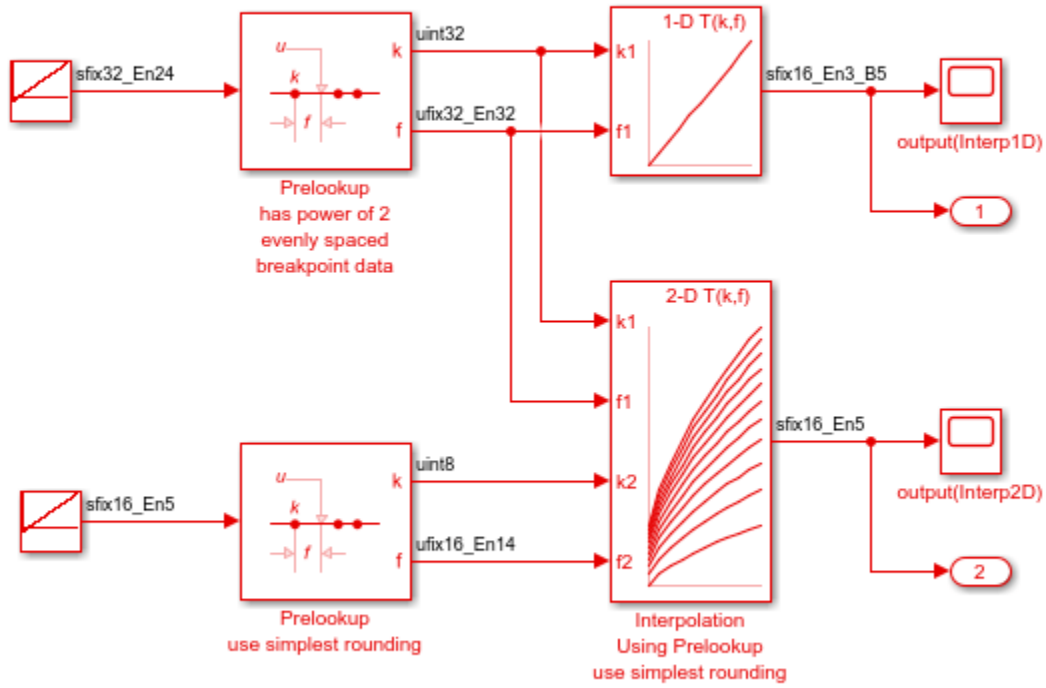
Copyright 1990-2012 The MathWorks, Inc.

Fixed-Point Prelookup and Interpolation

This example demonstrates some of the features of Prelookup and Interpolation Using Prelookup blocks.

- Prelookup and Interpolation Using Prelookup blocks support both floating-point data types and fixed-point data types.
- The algorithms that perform fraction calculation and evenly-spaced index search do not saturate. Therefore, the Prelookup block does not have a saturation parameter. Press **Ctrl+b** to generate code for the example. Observe the saturation-free algorithms in the generated code.
- Even if the **Saturate on integer overflow** parameter is checked, the algorithms that perform interpolation will saturate only when the **Intermediate results data type** cannot hold the intermediate results, or the **Output data type** cannot hold the result. Press **Ctrl+b** to generate code for the example model. Observe the saturation-free algorithms in the generated code.
- Prelookup and Interpolation Using Prelookup blocks support all rounding modes, including **Simplest** rounding mode. Double-click the blocks to open their dialogs and specify the rounding modes.
- When evenly-spaced index search is used with breakpoints spaced apart by a power of two in Prelookup blocks, the division needed to calculate indices is optimized using an efficient shift operation in the generated code.
- Simulink® always checks dimensional consistency between the **Breakpoint data** parameter of the Prelookup block and the **Table data** parameter of the Interpolation Using Prelookup block.
- Prelookup and Interpolation Using Prelookup blocks support two different indexing conventions, specified by the **Use last breakpoint for input at or above upper limit** parameter in the Prelookup block and the **Valid index input may reach last index** parameter in the Interpolation Using Prelookup block. Simulink® always checks the consistency of indexing conventions between these blocks.

Fixed-Point Prelookup and Interpolation



Copyright 2006-2022 The MathWorks, Inc.

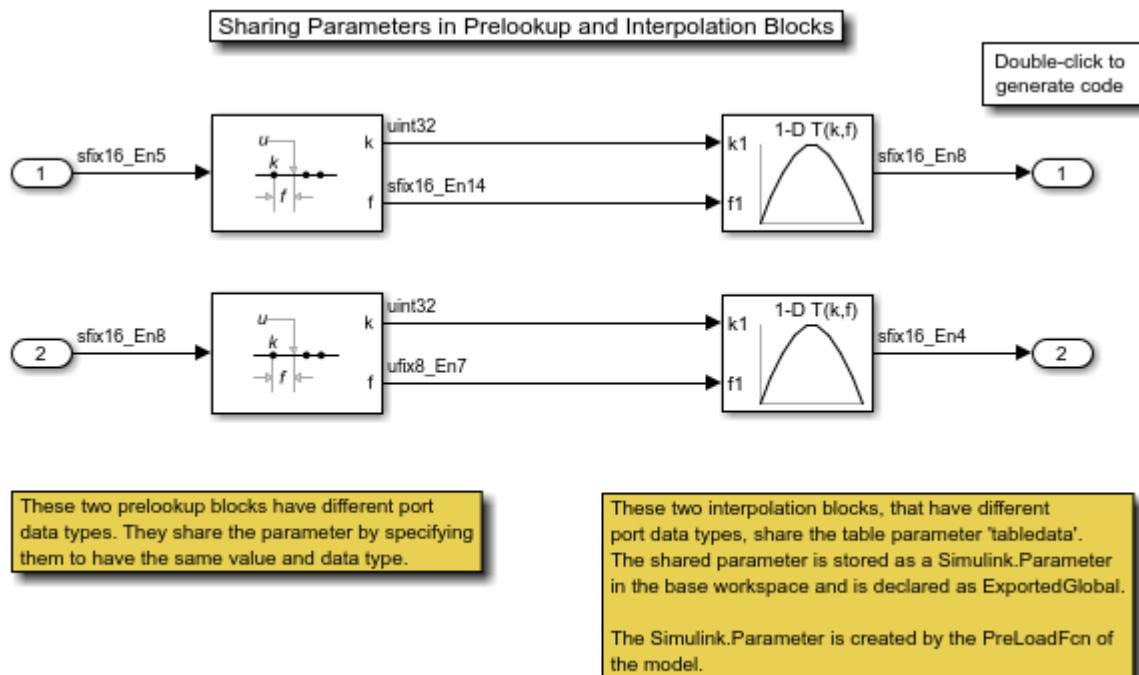
Sharing Parameters in Prelookup and Interpolation Blocks

This example shows how Prelookup and Interpolation blocks share their parameter data in generated code.

The Prelookup and Interpolation Using Prelookup blocks have support for specifying the data type for the breakpoints and table parameters. This makes it possible for blocks that have different port data types to share their parameter data in the generated code. One way of doing this is to specify the parameter with the same values and data types in multiple blocks. This is done in the Prelookup blocks in this model. Another way is to use a Simulink.Parameter object to define the blocks' shared parameter. This is done for the Interpolation blocks. The **Table data type** must be set to: **Inherit : Inherit from 'Table data'**.

For parameter sharing to take effect, the **Default parameter behavior** must be **Inlined** in the Optimization pane of the Configuration Parameters dialog box.

To see this in the generated code, open the model and build it.

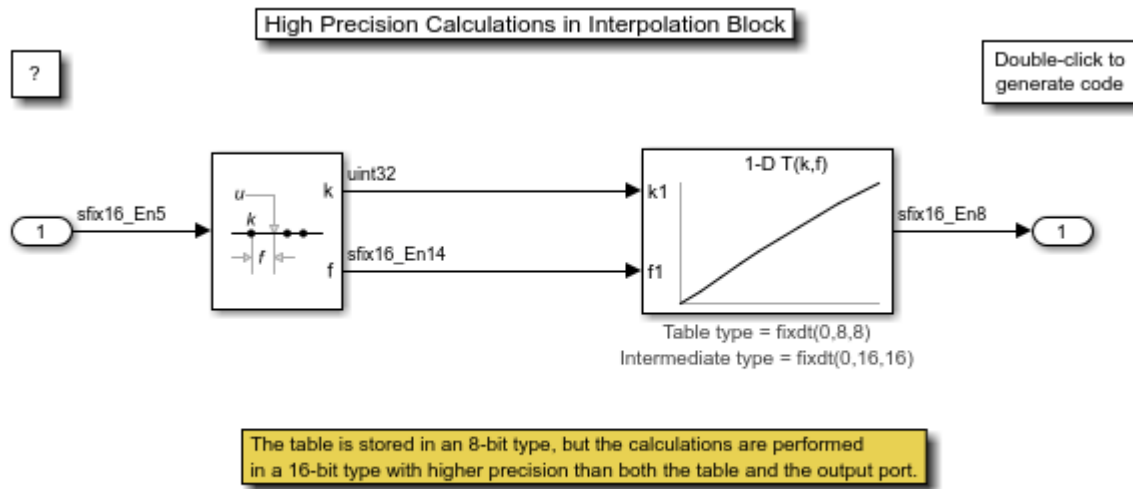


High Precision Calculations in Interpolation Block

This example shows how to perform high precision calculations in the Interpolation Using Prelookup block using internal rules. The Interpolation block allows the data type for intermediate results to be set.

In this model, the table is stored using an 8-bit data type, and the calculations are performed using a 16-bit data type. The default **Intermediate results data type** is `Inherit: Inherit via internal` rule that tries to maximize the precision of the intermediate results.

To see this in the generated code, open the model and build it.

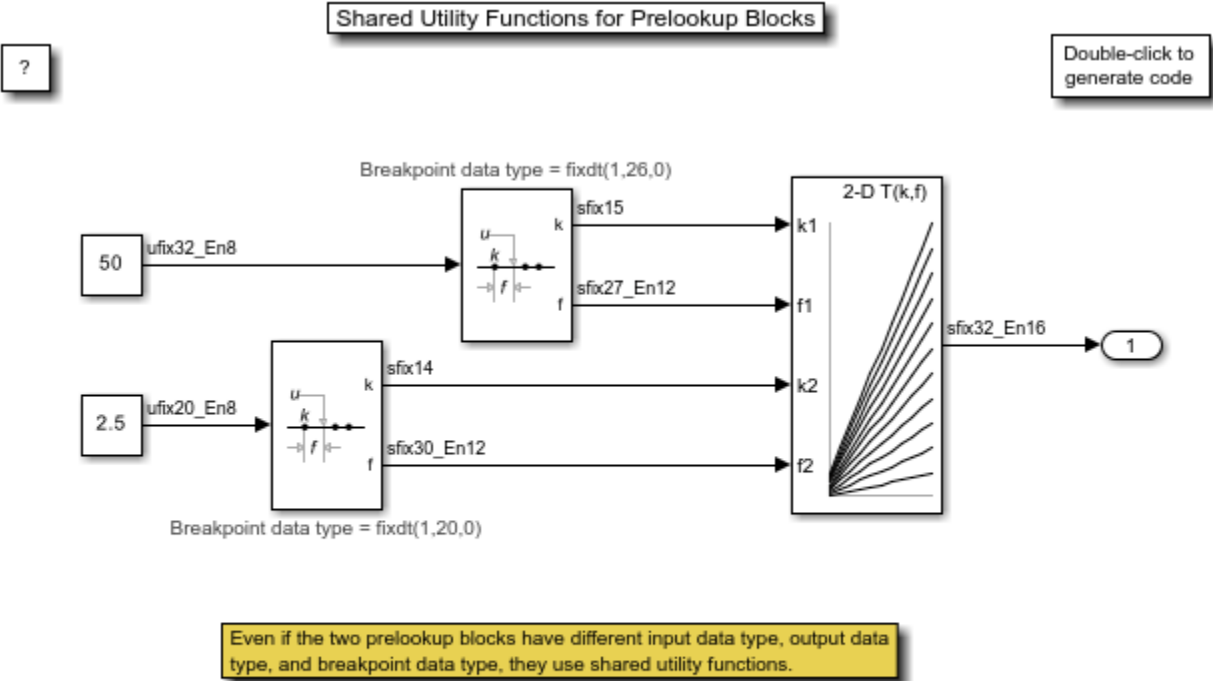


Copyright 2008-2009 The MathWorks, Inc.

Shared Utility Functions for Prelookup Blocks

This example shows how Prelookup blocks share utility functions. The utility functions generated by the Prelookup block are determined by the target data type of the block's inputs, outputs, and breakpoint parameter, as well as the **Index search method** and **Integer rounding mode**. Even if the two Prelookup blocks in this model have different data types, they fulfill the above requirements and share their utility functions.

To see this in the generated code, open the model and build it.



Fixed-Point Multiword Operations In Generated Code

This example shows how to control generation of multiword operations in generated code.

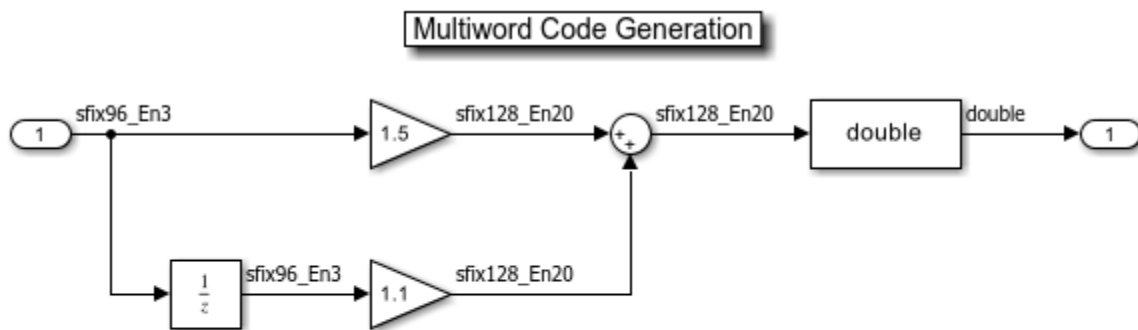
In this example you will learn:

- How to generate code for big data types using multiword operations
- How to prevent multiword code when it is not expected

Simple Multiword Operation

This model shows how wide integer and fixed-point operations become multiword in the generated C code. Multiword code is normally triggered by using parameters or signals with data types wider than C 'long'.

```
open_system('fxpdemo_multiword_example1');
set_param('fxpdemo_multiword_example1','SimulationCommand','Update');
```



Copyright 2012 The MathWorks, Inc.

Generate code for the model:

```
evalc('slbuild(''fxpdemo_multiword_example1'');'); % Suppress output
```

In the generated code, multiword operations are implemented using functions. These functions will have "MultiWord" in their name.

Review one of the multiword functions that was generated: MultiWordAdd()

```
fid = fopen('fxpdemo_multiword_example1_grt_rtw/fxpdemo_multiword_example1.c'); ctext = fread(fid, 'r');
match = regexp(ctext, 'void MultiWordAdd.*?\n\n}', 'match'); disp(match{1});

void MultiWordAdd(const uint32_T u1[], const uint32_T u2[], uint32_T y[],
                 int32_T n)
{
    int32_T i;
    uint32_T carry = 0U;
    uint32_T uli;
```



```

uint32_T yi;
for (i = 0; i < n; i++) {
    u1i = u1[i];
    yi = (u1i + u2[i]) + carry;
    y[i] = yi;
    carry = carry != 0U ? (uint32_T)(yi <= u1i) : (uint32_T)(yi < u1i);
}
}

```

This function implements multiword addition in C. The two operands and the result all have the same number of words, and the addition is performed one word at a time.

```
close_system('fxpdemo_multiword_example1', 0);
```

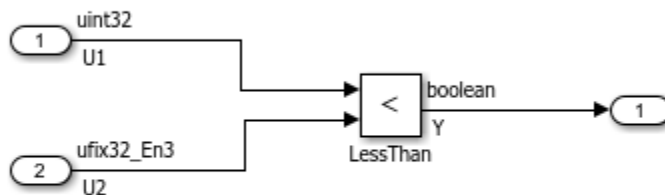
Relational Operator Block

In the Relational Operator Block example below, you would expect to get multiword functions in the generated code. The two input data types are `uint32` and `ufix32_En3`. A good type for comparison is `ufix35_En3`, because this type can represent all the real world values of both operands.

We expect the data type `ufix35_En3` to be implemented using a 64-bit two-word data type.

```
open_system('fxpdemo_multiword_example2');
set_param('fxpdemo_multiword_example2', 'SimulationCommand', 'Update');
```

Relational Operator Block for Multiword Code Generation



Copyright 2012 The MathWorks, Inc.

This model is configured for a CPU with 32-bit C type long. A 64-bit data type will be a multiword type.

```
get_param(bdroot, 'ProdBitPerLong')
```

```
ans =
```

```
32
```

Generate code for the model and review:

```

evalc('slbuild(''fxpdemo_multiword_example2'');'); % Suppress output
fid = fopen('fxpdemo_multiword_example2_grt_rtw/fxpdemo_multiword_example2.c') ; ctext = fread(fid, 'c');
match = regexp(ctext, 'void fxpdemo_multiword_example2_step.*?\n}', 'match'); disp(match{1});

void fxpdemo_multiword_example2_step(void)
{
    /* RelationalOperator: '<Root>/LessThan' incorporates:
    *   Inport: '<Root>/In1'
    *   Inport: '<Root>/In2'
    */
    Y = ((U1 <= 536870911U) && ((U1 << 3) < U2));
}

```

Multiword code was not generated. This code is single-word and uses a comparison data type of `uint32`. As a result, precision loss in the comparison may occur.

Simulink balances the requirements for the internal data type for comparison. In this case, because all data types are single word, it implements an efficient data type that produces small, fast code rather than a more precise, cumbersome computation.

```
close_system('fxpdemo_multiword_example2', 0);
```

To improve the precision of this calculation, do one of the following steps:

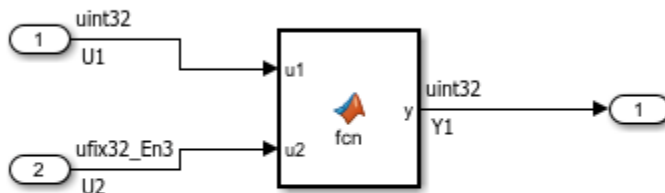
- Pick input data types that can be compared with full precision using a single word comparison type. For example, 16-bit types, or two identical types.
- Force Simulink to use a multiword type (during both simulation and code generation) by specifying a multiword type for at least one of the inputs. This signals Simulink that you want to use multiword operations for this block.
- Configure the model for a 64-bit system.

MATLAB Function Block

The MATLAB Function Block example below showcases an all-single-word calculation. Multiword code is not expected.

```
open_system('fxpdemo_multiword_example3');
set_param('fxpdemo_multiword_example3', 'SimulationCommand', 'Update');
```

MATLAB Function Block for Multiword Code Generation



```

mfb = get_param('fxpdemo_multiword_example3/MATLAB Function','MATLABFunctionConfiguration'); mfb

ans =

    'function y = fcn(u1, u2)
    %#codegen

    y = fi(u1 * u2, 0, 32, 0);'

```

Generate code for the model:

```

evalc('slbuild(''fxpdemo_multiword_example3'');'); % Suppress output
fid = fopen('fxpdemo_multiword_example3_grt_rtw/fxpdemo_multiword_example3.c') ; ctext = fread(fid, 'c');
match = regexp(ctext, 'void fxpdemo_multiword_example3_step.*?\n\n}', 'match'); disp(match{1});

void fxpdemo_multiword_example3_step(void)
{
    uint64m_T tmp;
    uint64m_T tmp_0;

    /* MATLAB Function: '<Root>/MATLAB Function' incorporates:
    *   Inport: '<Root>/In1'
    *   Inport: '<Root>/In2'
    */
    /* MATLAB Function 'MATLAB Function': '<S1>:1' */
    /* '<S1>:1:4' */
    uMultiWordMul(&U1, 1, &U2, 1, &tmp_0.chunks[0U], 2);
    uMultiWordShrNear(&tmp_0.chunks[0U], 2, 3U, &tmp.chunks[0U], 2);

    /* Outport: '<Root>/Out1' incorporates:
    *   MATLAB Function: '<Root>/MATLAB Function'
    */
    fxpdemo_multiword_example3_Y.Out1 = uMultiWord2uLongSat(&tmp.chunks[0U], 2);
}

close_system('fxpdemo_multiword_example3', 0);

```

Even though all the data types in the model were single-word, you still got three calls to multiword functions, and two multiword variables as well.

Fixed-point operations in the MATLAB Function Block are controlled by fimath property settings.

```
fimath
```

```

ans =

    RoundingMethod: Nearest
    OverflowAction: Saturate
    ProductMode: FullPrecision
    SumMode: FullPrecision

```

This fimath specifies full precision **ProductMode**. Therefore multiplications are performed in a way that preserves as much precision as possible. The product data type, which is `uint64_En3`, is implemented as a multiword type.

You can control multiword in code generation for MATLAB code by manipulating the `fimath`. For example:

- Adjusting `fimath` properties to meet efficient code requirements. In this example, set 'ProductMode' to 'KeepLSB' and 'OverflowAction' to 'Wrap'.
- Defining local `fimaths` in the MATLAB Function Block that are tailored for the specific calculations, and not relying on the global `fimath`.

```
load_system('fxpdemo_multiword_example4'); % No need to show this model. Only show the MATLAB code
mfb = get_param('fxpdemo_multiword_example4/MATLAB Function','MATLABFunctionConfiguration'); mfb
```

```
ans =
```

```
'function y = fcn(u1, u2)
    %#codegen
    F = fimath('ProductMode','KeepLSB',...
        'ProductWordLength',32,...
        'OverflowAction','Wrap');
    u1 = setfimath(u1,F);
    u2 = setfimath(u2,F);
    y = fi(u1 * u2,0,32,0);
```

This `fimath` will result in this generated code:

```
evalc('slbuild(''fxpdemo_multiword_example4'');'); % Suppress output
fid = fopen('fxpdemo_multiword_example4_grt_rtw/fxpdemo_multiword_example4.c'); ctext = fread(fid);
match = regexp(ctext, 'void fxpdemo_multiword_example4_step.*?\n}', 'match'); disp(match{1});
```

```
void fxpdemo_multiword_example4_step(void)
{
    uint32_T tmp;

    /* MATLAB Function: '<Root>/MATLAB Function' incorporates:
     * Inport: '<Root>/In1'
     * Inport: '<Root>/In2'
     */
    /* MATLAB Function 'MATLAB Function': '<S1>:1' */
    /* '<S1>:1:6' */
    /* '<S1>:1:7' */
    /* '<S1>:1:8' */
    tmp = U1 * U2;
    Y1 = (uint32_T)((tmp & 4U) != 0U) + (tmp >> 3);
}
```

```
close_system('fxpdemo_multiword_example4', 0);
```

```
clear ctext fid match mfb
clear ans
```

See Also

MATLABFunctionConfiguration | fimath

Fixed-Point Multiplication Helper Functions in Generated Code

This example shows how to control the generation of multiplication helper functions in the generated code.

In this example you will learn:

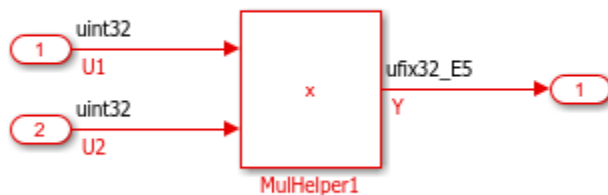
- Under what circumstances these helpers are generated
- What approaches are available to eliminate them

Example 1 - Helper Function Due To Range Of Output

The following model creates helper functions for multiplication because the range of the ideal product exceeds the range of a single word data type.

```
open_system('fxpdemo_mulhelpers_example1');
set_param('fxpdemo_mulhelpers_example1', 'SimulationCommand', 'Update');
```

Multiplication on a 32-bit target



Copyright 2012 The MathWorks, Inc.

Generate code for the model and review it:

```
evalc('slbuild(''fxpdemo_mulhelpers_example1'');'); % Suppress output
fid = fopen('fxpdemo_mulhelpers_example1_grt_rtw/fxpdemo_mulhelpers_example1.c'); ctext = fread
```

Two helper functions were generated: `mul_wide_u32()` and `mul_u32_loSR()`.

The step function calls the outer helper function `mul_u32_loSR()`:

```
match = regexp(ctext, 'void fxpdemo_mulhelpers_example1_step.*?\n\n}', 'match'); disp(match{1});
void fxpdemo_mulhelpers_example1_step(void)
{
  /* Product: '<Root>/MulHelper1' incorporates:
   * Inport: '<Root>/In1'
   * Inport: '<Root>/In2'
   */
  Y = mul_u32_loSR(U1, U2, 5U);
}
```

The outer helper function calls an inner helper function `mul_wide_u32()`:

```
match = regexp(ctext, 'uint32_T mul_u32_loSR.*?\n\n}', 'match'); disp(match{1});
```

```
uint32_T mul_u32_loSR(uint32_T a, uint32_T b, uint32_T aShift)
{
    uint32_T result;
    uint32_T u32_chi;
    mul_wide_u32(a, b, &u32_chi, &result);
    return u32_chi << (32U - aShift) | result >> aShift;
}
```

The inner helper function `mul_wide_u32()` is also generated in the file:

```
match = regexp(ctext, 'void mul_wide_u32.*?\n\}', 'match'); disp(match{1});
void mul_wide_u32(uint32_T in0, uint32_T in1, uint32_T *ptrOutBitsHi, uint32_T
                *ptrOutBitsLo)
{
    uint32_T in0Hi;
    uint32_T in0Lo;
    uint32_T in1Hi;
    uint32_T in1Lo;
    uint32_T outBitsLo;
    uint32_T productHiLo;
    uint32_T productLoHi;
    in0Hi = in0 >> 16U;
    in0Lo = in0 & 65535U;
    in1Hi = in1 >> 16U;
    in1Lo = in1 & 65535U;
    productHiLo = in0Hi * in1Lo;
    productLoHi = in0Lo * in1Hi;
    in0Lo *= in1Lo;
    in1Lo = 0U;
    outBitsLo = (productLoHi << 16U) + in0Lo;
    if (outBitsLo < in0Lo) {
        in1Lo = 1U;
    }

    in0Lo = outBitsLo;
    outBitsLo += productHiLo << 16U;
    if (outBitsLo < in0Lo) {
        in1Lo++;
    }

    *ptrOutBitsHi = (((productLoHi >> 16U) + (productHiLo >> 16U)) + in0Hi * in1Hi)
        + in1Lo;
    *ptrOutBitsLo = outBitsLo;
}
```

Why Was A Helper Function Generated For Example 1?

Multiplying two 32-bit numbers yields a 64-bit ideal result. In `example1`, the type of the ideal result is `uint64`.

The actual output type of `ufix32_E5` uses only bits 5 to 36 from these 64 bits.

```
bits = 0:63;
ideal_product_bits = ones(1,64);
most_significant_word = [ones(1,32) NaN*ones(1,32)];
least_significant_word = [NaN*ones(1,32) ones(1,32)];
output_bits = [NaN*ones(1, 27) ones(1,32) NaN*ones(1,5)];
```

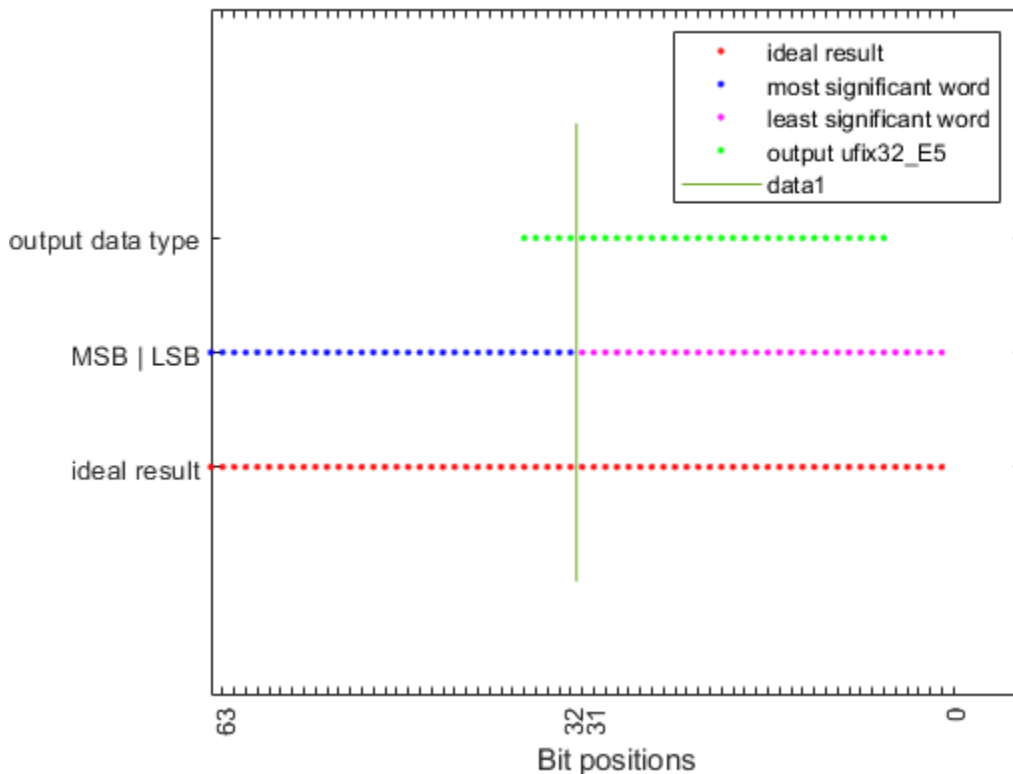
```

plot(bits, ideal_product_bits, '.r');
set(gca, 'YLim', [-1 5]);
set(gca, 'YTick', 1:3);
set(gca, 'YTickLabel', {'ideal result', 'MSB | LSB', 'output data type'});
set(gca, 'XTick', 1:64);

% "Move" interesting tick labels so they're readable and show
set(gca, 'XTickLabel', {'63', '31', '32', '0'});

set(get(gca, 'XLabel'), 'String', 'Bit positions')
hold on
plot(bits, most_significant_word*2, '.b');
plot(bits, least_significant_word*2, '.m');
plot(bits, output_bits*3, '.g');
legend('ideal result', 'most significant word', 'least significant word', 'output ufix32_E5') %
plot([31.5 31.5], [0 4]); % MSB-LSB boundary
clear bits ideal_product_bits most_significant_word least_significant_word output_bits

```



The output data type includes bits from both most and least significant words of the ideal product type. Even though the input and output data types are single word, it takes two words to construct the ideal result, so it can be converted to the output type. Helper functions are generated to manipulate this wide ideal product.

Why Are There Two Helper Functions?

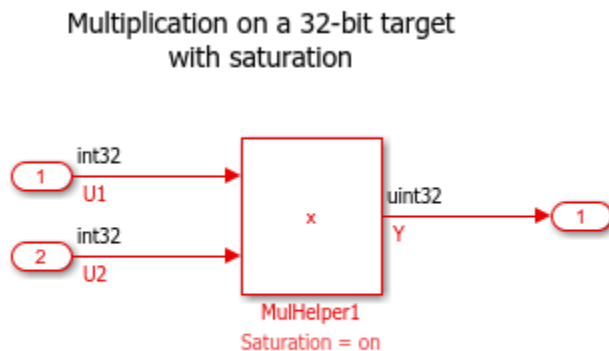
The combination of the double-word multiplication and the cast is encoded in the two helper functions. The inner function, `mul_wide_u32()` performs long multiplication of two 32-bit operands into a 64-bit result. The outer function `mul_u32_loSR()` performs data type conversion from the raw `uint64` result to the desired output data type.

```
close_system('fxpdemo_mulhelpers_example1', 0);
close all;
```

Example 2: Helper Functions Due To Saturation

The following model creates helper functions for multiplication, in order to implement saturation.

```
open_system('fxpdemo_mulhelpers_example2');
set_param('fxpdemo_mulhelpers_example2', 'SimulationCommand', 'Update');
```



Copyright 2012 The MathWorks, Inc.

```
evalc('slbuild(''fxpdemo_mulhelpers_example2'');'); % Suppress output
fid = fopen('fxpdemo_mulhelpers_example2_grt_rtw/fxpdemo_mulhelpers_example2.c'); ctext = fread
match = regexp(ctext, 'void fxpdemo_mulhelpers_example2_step.*?\n}', 'match'); disp(match{1});
```

```
void fxpdemo_mulhelpers_example2_step(void)
{
  /* Product: '<Root>/MulHelper1' incorporates:
   * Inport: '<Root>/In1'
   * Inport: '<Root>/In2'
   */
  Y = mul_us32_sat(U1, U2);
}
```

This model does not have the range problem that the model in example 1 had. The output type fits within the least significant word of the ideal product.

In this model, the ideal type can be negative but the output is unsigned. To saturate negative values to 0, the sign bit needs to be computed. This sign bit resides in the most significant word of the ideal result. Coder again generates two helpers, one to compute the raw 64-bit result, and one to perform a saturating cast.

```
close_system('fxpdemo_mulhelpers_example2', 0);
```


Avoiding Helper Functions By Changing Numerical Properties

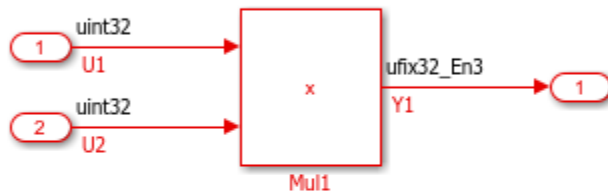
It is possible to avoid generating the multiplication helpers by:

- Using wrapping multiplications
- Using output types that don't exceed the range of the least significant word
- Multiplying smaller data types with a single word ideal product

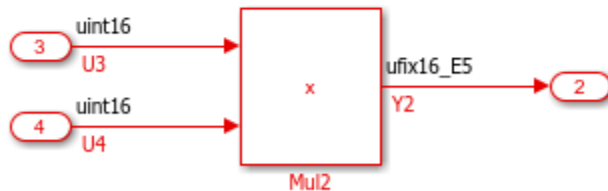
```
open_system('fxpdemo_mulhelpers_example3');
set_param('fxpdemo_mulhelpers_example3', 'SimulationCommand', 'Update');
```

Multiplication on a 32-bit target without helper functions

This multiplication does not generate a helper function
because the range of the output data type
which is $2^{28}-1$
is within the range of the ideal product's least significant word
which is $2^{31}-1$.



This multiplication does not generate a helper function
because the ideal product's data type fits in a single word.



Copyright 2012 The MathWorks, Inc.

```
evalc('slbuild(''fxpdemo_mulhelpers_example3'');'); % Suppress output
fid = fopen('fxpdemo_mulhelpers_example3_grt_rtw/fxpdemo_mulhelpers_example3.c'); ctext = fread
match = regexp(ctext, 'void fxpdemo_mulhelpers_example3_step.*?\n\}', 'match'); disp(match{1});
```

```
void fxpdemo_mulhelpers_example3_step(void)
{
  /* Product: '<Root>/Mul1' incorporates:
   * Inport: '<Root>/In1'
   * Inport: '<Root>/In2'
   */
  Y1 = (U1 * U2) << 3;
```

```

/* Product: '<Root>/Mul2' incorporates:
 * Inport: '<Root>/In3'
 * Inport: '<Root>/In4'
 */
Y2 = (uint16_T)((((uint32_T)U3 * U4) >> 5);
}

close_system('fxpdemo_mulhelpers_example3', 0);

```

Avoiding Helper Functions By Specifying Design Ranges

If large multiplicands are a necessity but the range of possible values on the operands is small enough, it may be possible to avoid the helper functions by specifying design ranges.

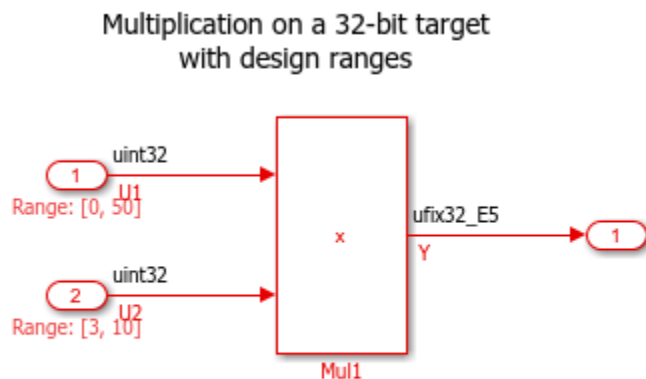
Note: This feature requires an Embedded Coder license.

Let's review a model similar to the one in example 1:

```

open_system('fxpdemo_mulhelpers_example4');
set_param('fxpdemo_mulhelpers_example4', 'SimulationCommand', 'Update');

```



Copyright 2012 The MathWorks, Inc.

Note that we selected the 'Optimize using the specified minimum and maximum values' checkbox on the optimization pane of the model configuration dialog.

```
get_param(bdroot, 'UseSpecifiedMinMax')
```

```
ans =
    'on'
```

Examine the step function in the generated code:

```

evalc('slbuild(''fxpdemo_mulhelpers_example4'');'); % Suppress output
fid = fopen('fxpdemo_mulhelpers_example4_ert_rtw/fxpdemo_mulhelpers_example4.c'); ctext = fread
match = regexp(ctext, 'void fxpdemo_mulhelpers_example4_step.*?\n}', 'match'); disp(match{1});

```

```
void fxpdemo_mulhelpers_example4_step(void)
{
  /* Product: '<Root>/Mul1' incorporates:
   * Inport: '<Root>/In1'
   * Inport: '<Root>/In2'
   */
  Y = (U1 * U2) >> 5;
}
```

The range of the possible values of the ideal product is from 0 to 500. A single word data type is sufficient to express this product.

```
close_system('fxpdemo_mulhelpers_example4', 0);
```

Caution: If, while running the generated code, values on the signals exceed the design ranges, incorrect results may be computed.

Avoiding Helper Functions Using CRL Operator Replacement

You can customize code generation by replacing the cumbersome multiplication operators with your own C macros or custom functions. You may want to do this if, for example, your CPU has resources to implement wide multiplications more efficiently.

References:

- “Define Code Replacement Library Optimizations” (Embedded Coder)
- “Register Code Replacement Library” (Embedded Coder)

```
clear ctext fid match
clear ans
```

Fixed-Point Optimizations Using Specified Minimum and Maximum Values

This example shows how to optimize fixed-point operations in generated code using minimum and maximum values that you specify in a model.

Overview

Minimum and maximum values can represent environmental limits or mechanical limits, such as the output ranges of sensors. The code generation software can use these values to create more efficient code by eliminating unreachable code branches and unnecessary utility functions.

Note: You must ensure that the specified minimum and maximum values are accurate and trustworthy. Otherwise, optimization might result in numerical mismatch with simulation.

The benefits of optimizing the generated code are:

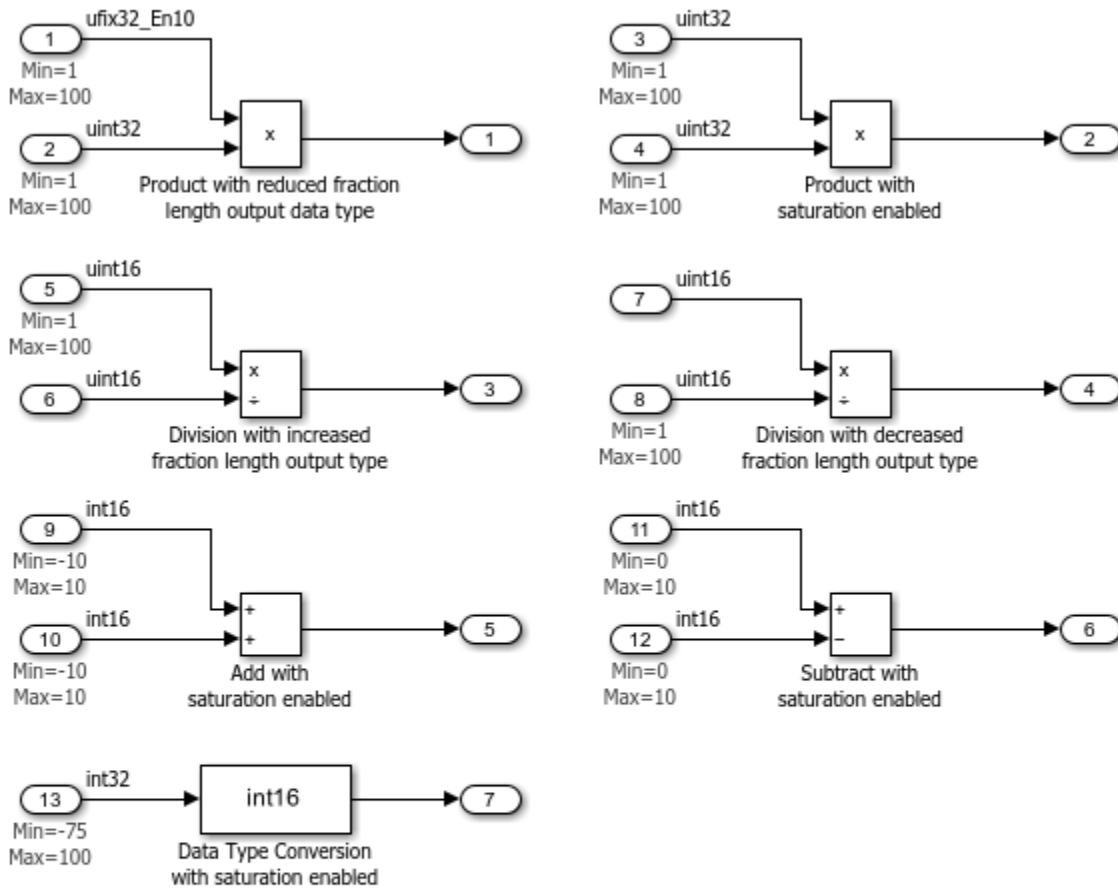
- Reducing the ROM and RAM consumption.
- Improving the execution speed.

Open Model

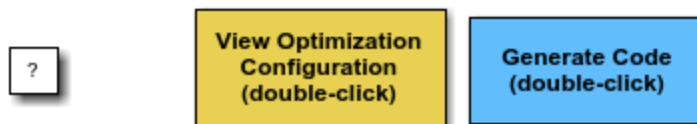
Open the model `fxpdemo_min_max_optimization`.

```
open_system('fxpdemo_min_max_optimization');
```

Fixed Point Optimizations Using Specified Min and Max Values



Optimizing various fixed-point operations using specified minimum and maximum values



Copyright 2010 The MathWorks, Inc.

Inspect Model

In this model, there are minimum and maximum values specified at the input ports upstream of the various fixed-point blocks. By utilizing these values, every fixed-point operation in the model is optimized in some way.

Generate Code without Optimization

First, without using the specified minimum and maximum values, generate code for this model. Double-click the blue button.

```
rtwbuild('fxpdemo_min_max_optimization');
### Starting build procedure for: fxpdemo_min_max_optimization
### Successful completion of code generation for: fxpdemo_min_max_optimization
```

Build Summary

Top model targets built:

| Model | Action | Rebuild Reason |
|------------------------------|-----------------|--|
| fxpdemo_min_max_optimization | Code generated. | Code generation information file does not exist. |

1 of 1 models built (0 models already up to date)
Build duration: 0h 1m 28.273s

Inspect Code without Optimization

An HTML Code Generation Report opens. Examine the code corresponding to the block "Product with reduced fraction length output data type". Right-click on this block and select **Code Generation > Navigate to code...**

```
rtwtrace('fxpdemo_min_max_optimization/Product with reduced fraction length output data type');
```

The generated code is:

```
rtY.Out1 = mul_u32_u32_u32_sr10(rtU.In1, rtU.In2);
```

To implement this fixed-point multiplication operation, the code generation software must generate a utility function `mul_u32_u32_u32_sr10`. Also, to implement `mul_u32_u32_u32_sr10`, it must generate a second utility function, `mul_wide_u32`. These functions consist of many lines of code and require several temporary variables.

Enable Optimization

- 1 Double-click the yellow button to open the Configuration Parameters dialog box.
- 2 In this dialog box, under **Code generation**, select **Optimize using specified minimum and maximum values**.

```
set_param('fxpdemo_min_max_optimization', 'UseSpecifiedMinMax', 'on');
```

Generate Code with Optimization

Now regenerate the code using the specified minimum and maximum values.

Double-click the blue button.

```
rtwbuild('fxpdemo_min_max_optimization');
### Starting build procedure for: fxpdemo_min_max_optimization
### Successful completion of code generation for: fxpdemo_min_max_optimization
```

Build Summary

Top model targets built:

| Model | Action | Rebuild Reason |
|-------|--------|----------------|
| ===== | | |

```
fxpdemo_min_max_optimization Code generated. Generated code was out of date.
```

```
1 of 1 models built (0 models already up to date)  
Build duration: 0h 0m 30.866s
```

Examine Optimized Code

Right-click on block "Product with reduced fraction length output data type" and select **Code Generation > Navigate to code...**

```
rtwtrace('fxpdemo_min_max_optimization/Product with reduced fraction length output data type');
```

The generated code is:

```
rtY.Out1 = rtU.In1 * rtU.In2 >> 10;
```

Using the specified minimum and maximum values, the code generation software determines that it can safely implement the reduced fraction length at the output with a right shift, and does not generate utility functions.

Examine Other Operations

Examine other operations in the generated code to see how the code generation software uses the specified minimum and maximum values. The code generation software now implements each fixed-point operation with simple C operations and eliminates unnecessary helper functions and code branches.

Fixed-Point Function Approximation

When a fixed-point library function is not available, fixed-point applications require an approximation of the function. Often, an interpolated look up table is used to store an approximation of the function over a specified range.

This example shows how to approximate the function $y = \sin(2\pi x)$ over a specified input range using a lookup table.

See also: `FunctionApproximation.Problem`, `FunctionApproximation.Options`

In this example, the input uses an unsigned 16-bit data type, `fixdt(0,16,16)`, and the output uses a signed 16-bit data type, `fixdt(1,16,14)`.

The goal of this example is to create an approximation that is accurate to 8 bits to the right of the binary point. This means that the worst case error should be less than 2^{-8} .

There are many sets of lookup table data points that would meet this goal. Different implementations can be chosen depending on goals such as memory usage and speed of computation. This example finds a solution that meets this accuracy goal with the minimal number of data points.

Approximate Function

Use the `FunctionApproximation.Options` object to specify accuracy and word length constraints.

```
options = FunctionApproximation.Options();
options.AbsTol = 2^-8;
options.RelTol = 0;
options.WordLengths = [8 16 32];
options.MemoryUnits = 'bytes';
options.OnCurveTableValues = true;
```

Specify the function to approximate and the input ranges and data types in the `FunctionApproximation.Problem` object.

```
functionToApproximate = @(x) sin(2*pi*x);

problem = FunctionApproximation.Problem(functionToApproximate, 'Options', options);
problem.InputTypes = numerictype(0,16,16);
problem.InputLowerBounds = 0;
problem.InputUpperBounds = 0.25;
problem.OutputType = numerictype(1,16,14);

% Create a LUT solution
solution = solve(problem);

% Change breakpoint specification to EvenPow2Spacing and create a LUT
% solution again
problem.Options.BreakpointSpecification = 'EvenPow2Spacing';
bestEvenPow2SpacingSolution = solve(problem);
```

Searching for fixed-point solutions.

| ID | Memory (bytes) | Feasible | Table Size | Breakpoints W/Ls | TableData WL | BreakpointSpec |
|----|----------------|----------|------------|------------------|--------------|----------------|
|----|----------------|----------|------------|------------------|--------------|----------------|

| | | | | | | |
|----|------------|---|----|----|----|----------|
| 0 | 4.0000e+00 | 0 | 2 | 8 | 8 | EvenPow |
| 1 | 1.9000e+01 | 0 | 17 | 8 | 8 | EvenPow |
| 2 | 3.5000e+01 | 1 | 33 | 8 | 8 | EvenPow |
| 3 | 2.8000e+01 | 1 | 26 | 8 | 8 | EvenPow |
| 4 | 2.4000e+01 | 1 | 22 | 8 | 8 | EvenPow |
| 5 | 2.1000e+01 | 1 | 19 | 8 | 8 | EvenPow |
| 6 | 1.5000e+01 | 0 | 13 | 8 | 8 | EvenPow |
| 7 | 1.7000e+01 | 1 | 15 | 8 | 8 | EvenPow |
| 8 | 1.2000e+01 | 0 | 10 | 8 | 8 | EvenPow |
| 9 | 1.1000e+01 | 0 | 9 | 8 | 8 | EvenPow |
| 10 | 1.4000e+01 | 1 | 12 | 8 | 8 | EvenPow |
| 11 | 1.3000e+01 | 0 | 11 | 8 | 8 | EvenPow |
| 12 | 9.0000e+00 | 0 | 7 | 8 | 8 | EvenPow |
| 13 | 6.0000e+00 | 0 | 2 | 16 | 8 | EvenPow |
| 14 | 1.3000e+01 | 0 | 9 | 16 | 8 | EvenPow |
| 15 | 6.0000e+00 | 0 | 2 | 8 | 16 | EvenPow |
| 16 | 4.0000e+00 | 0 | 2 | 8 | 8 | EvenPow |
| 17 | 1.1000e+01 | 0 | 9 | 8 | 8 | EvenPow |
| 18 | 6.0000e+00 | 0 | 2 | 16 | 8 | EvenPow |
| 19 | 1.3000e+01 | 0 | 9 | 16 | 8 | EvenPow |
| 20 | 6.0000e+00 | 0 | 2 | 8 | 16 | EvenPow |
| 21 | 8.0000e+00 | 0 | 2 | 16 | 16 | EvenPow |
| 22 | 1.4000e+01 | 1 | 7 | 8 | 8 | Explicit |
| 23 | 3.5000e+01 | 1 | 33 | 8 | 8 | EvenPow |

Best Solution

| ID | Memory (bytes) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|----------------|----------|------------|-----------------|--------------|----------------|
| 10 | 1.4000e+01 | 1 | 12 | 8 | 8 | EvenPow |

Searching for fixed-point solutions.

| ID | Memory (bytes) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|----------------|----------|------------|-----------------|--------------|----------------|
| 0 | 4.0000e+00 | 0 | 2 | 8 | 8 | EvenPow |
| 1 | 1.9000e+01 | 0 | 17 | 8 | 8 | EvenPow |
| 2 | 3.5000e+01 | 1 | 33 | 8 | 8 | EvenPow |
| 3 | 1.1000e+01 | 0 | 9 | 8 | 8 | EvenPow |
| 4 | 6.0000e+00 | 0 | 2 | 16 | 8 | EvenPow |
| 5 | 2.1000e+01 | 0 | 17 | 16 | 8 | EvenPow |
| 6 | 1.3000e+01 | 0 | 9 | 16 | 8 | EvenPow |
| 7 | 6.0000e+00 | 0 | 2 | 8 | 16 | EvenPow |
| 8 | 2.0000e+01 | 0 | 9 | 8 | 16 | EvenPow |

Best Solution

| ID | Memory (bytes) | Feasible | Table Size | Breakpoints WLS | TableData WL | BreakpointSpec |
|----|----------------|----------|------------|-----------------|--------------|----------------|
| 2 | 3.5000e+01 | 1 | 33 | 8 | 8 | EvenPow |

Explore Solutions

```
%The software returns several implementations that meet the requirements
%specified in the |FunctionApproximation.Problem| and
%|FunctionApproximation.Options| objects. You can explore these different
%implementations.
```

```
feasibleSolutions = solution.FeasibleSolutions;
tableDataVec = [feasibleSolutions.TableData];
evenSpacingSolutions = find([tableDataVec.IsEvenSpacing]);
unevenSpacingSolutions = find(~[tableDataVec.IsEvenSpacing]);
```

```

evenSolutionsMemoryUsage = arrayfun(@(x) x.totalMemoryUsage(), feasibleSolutions(evenSpacingSo
unevenSolutionsMemoryUsage = arrayfun(@(x) x.totalMemoryUsage(), feasibleSolutions(unevenSpacingSo

bestEvenSpacingSolution = feasibleSolutions(evenSpacingSolutions(evenSolutionsMemoryUsage == mi
bestUnevenSpacingSolution = feasibleSolutions(unevenSpacingSolutions(unevenSolutionsMemoryUsage =

xeven = bestEvenSpacingSolution.TableData.BreakpointValues{1};
yeven = bestEvenSpacingSolution.TableData.TableValues;

xuneven = bestUnevenSpacingSolution.TableData.BreakpointValues{1};
yuneven = bestUnevenSpacingSolution.TableData.TableValues;

xpow2 = bestEvenPow2SpacingSolution.TableData.BreakpointValues{1};
ypow2 = bestEvenPow2SpacingSolution.TableData.TableValues;

```

Compare Memory Usage

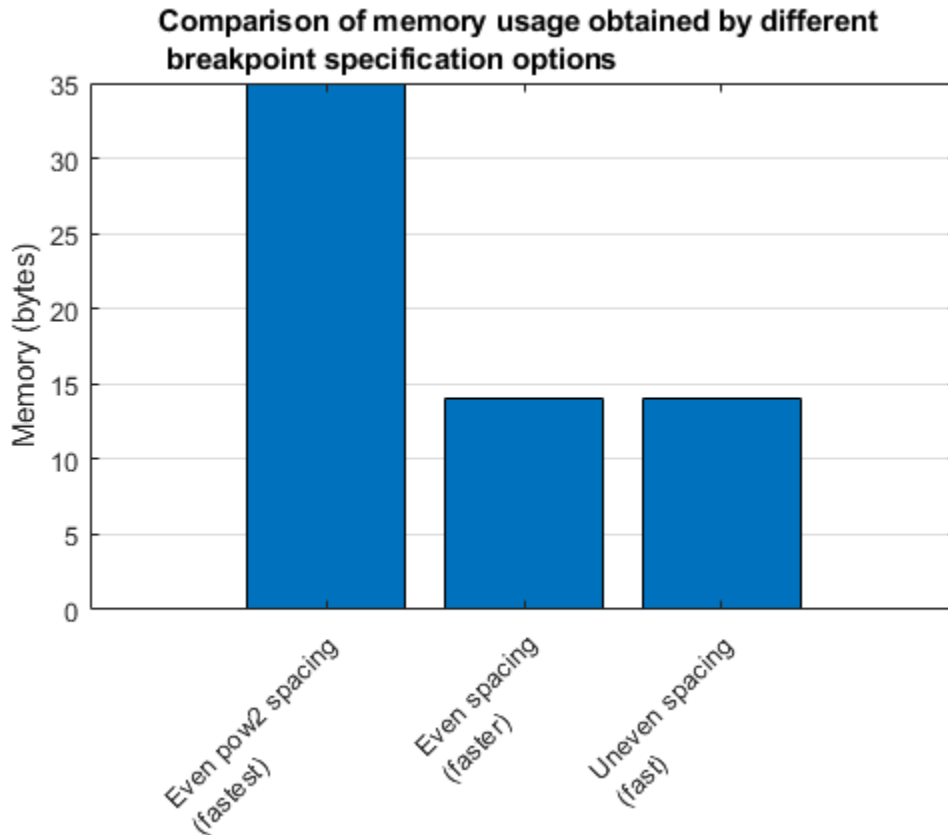
Compare the memory used by the lookup tables.

```

memoryValues = [...
    totalMemoryUsage(bestEvenPow2SpacingSolution), ...
    totalMemoryUsage(bestEvenSpacingSolution), ...
    totalMemoryUsage(bestUnevenSpacingSolution)];

figure();
xTickLabels = {'Even pow2 spacing \newline(fastest)', 'Even spacing \newline(faster)', 'Uneven spacing \newline(fastest)'};
hMemory = bar(memoryValues);
title('Comparison of memory usage obtained by different \newline breakpoint specification options');
hMemory.Parent.XTickLabel = xTickLabels;
hMemory.Parent.XTickLabelRotation = 45;
hMemory.Parent.YLabel.String = 'Memory (bytes)';
hMemory.Parent.Box = 'on';
hMemory.Parent.YGrid = 'on';

```



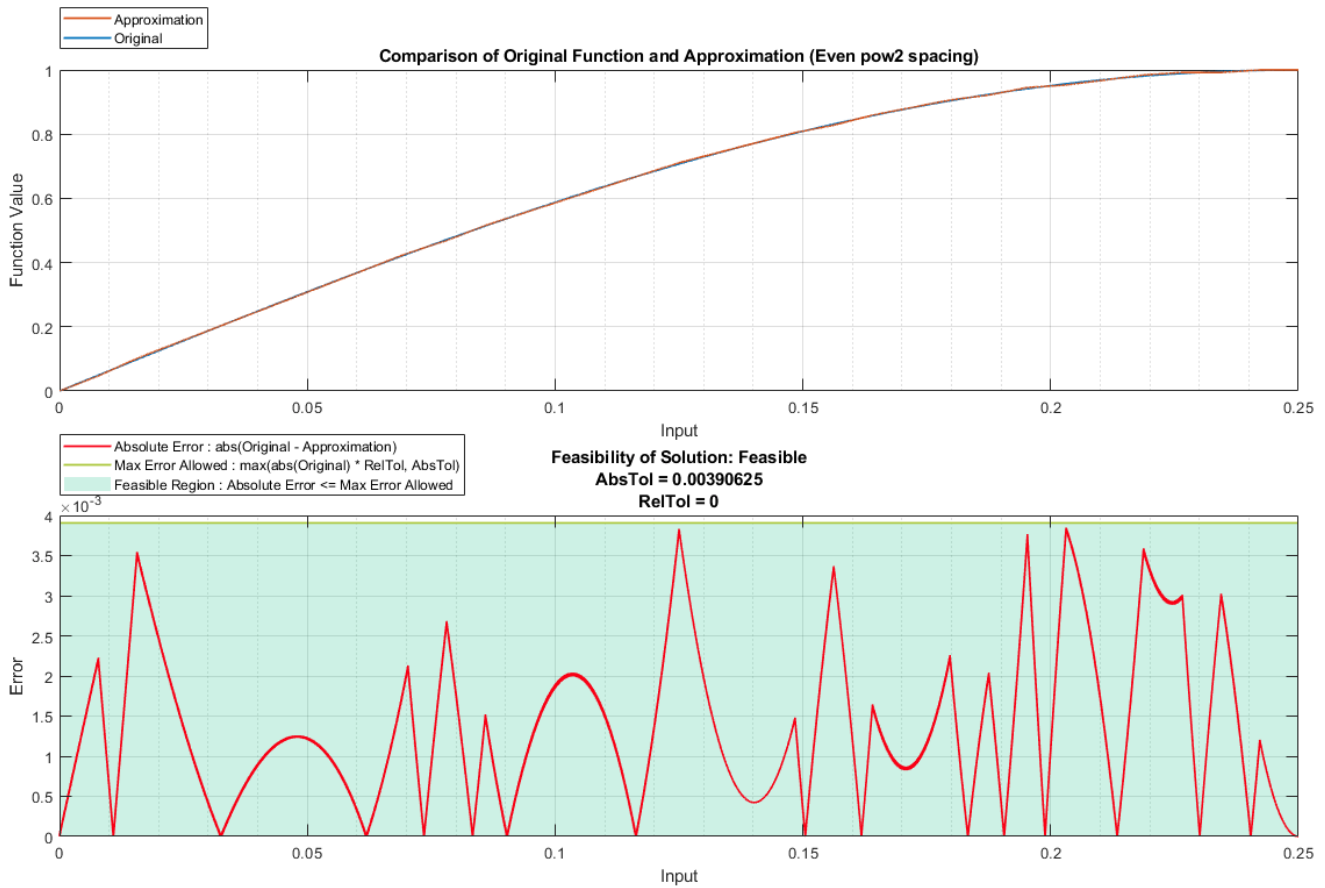
The amount of memory used by the tables using even spacing and uneven spacing are the same, but the number of points are different. This is because when storing the breakpoints for tables using even spacing, only the first point and spacing are stored. In contrast, all breakpoints are stored for tables using uneven spacing.

Even pow2 spacing stores more than double the points stored for even spacing. From a memory usage perspective, even pow2 spacing is the least optimal for this function. However, computations for even pow2 spacing are performed using arithmetic shifts instead of multiplication, which can lead to faster execution times.

Compare Solutions to Original Function

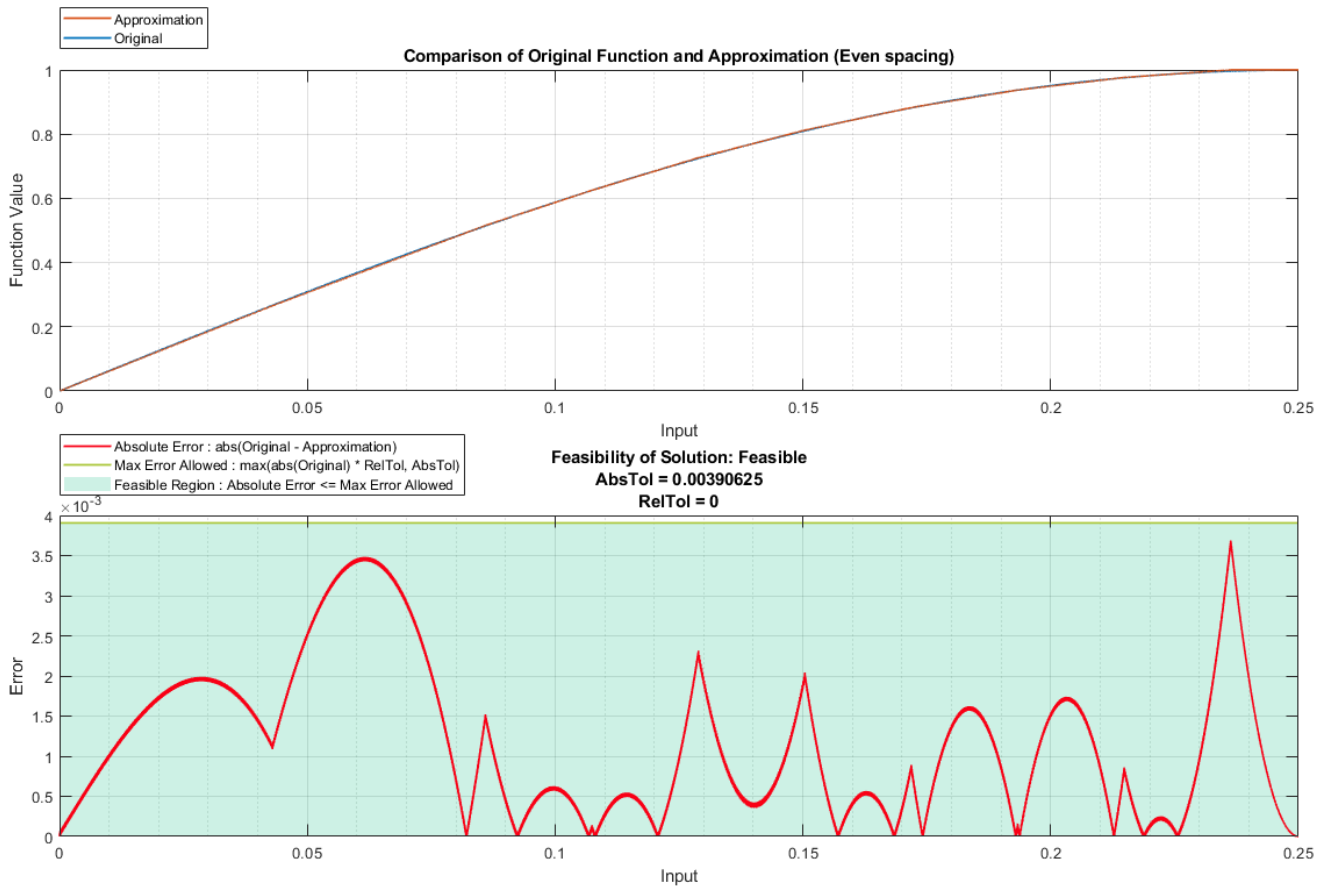
Compare the solution using even pow2 spacing to the original function.

```
[~, hEvenPow2Spacing] = compare(bestEvenPow2SpacingSolution);
hEvenPow2Spacing.Children(4).Title.String = [hEvenPow2Spacing.Children(4).Title.String ' (Even p
```



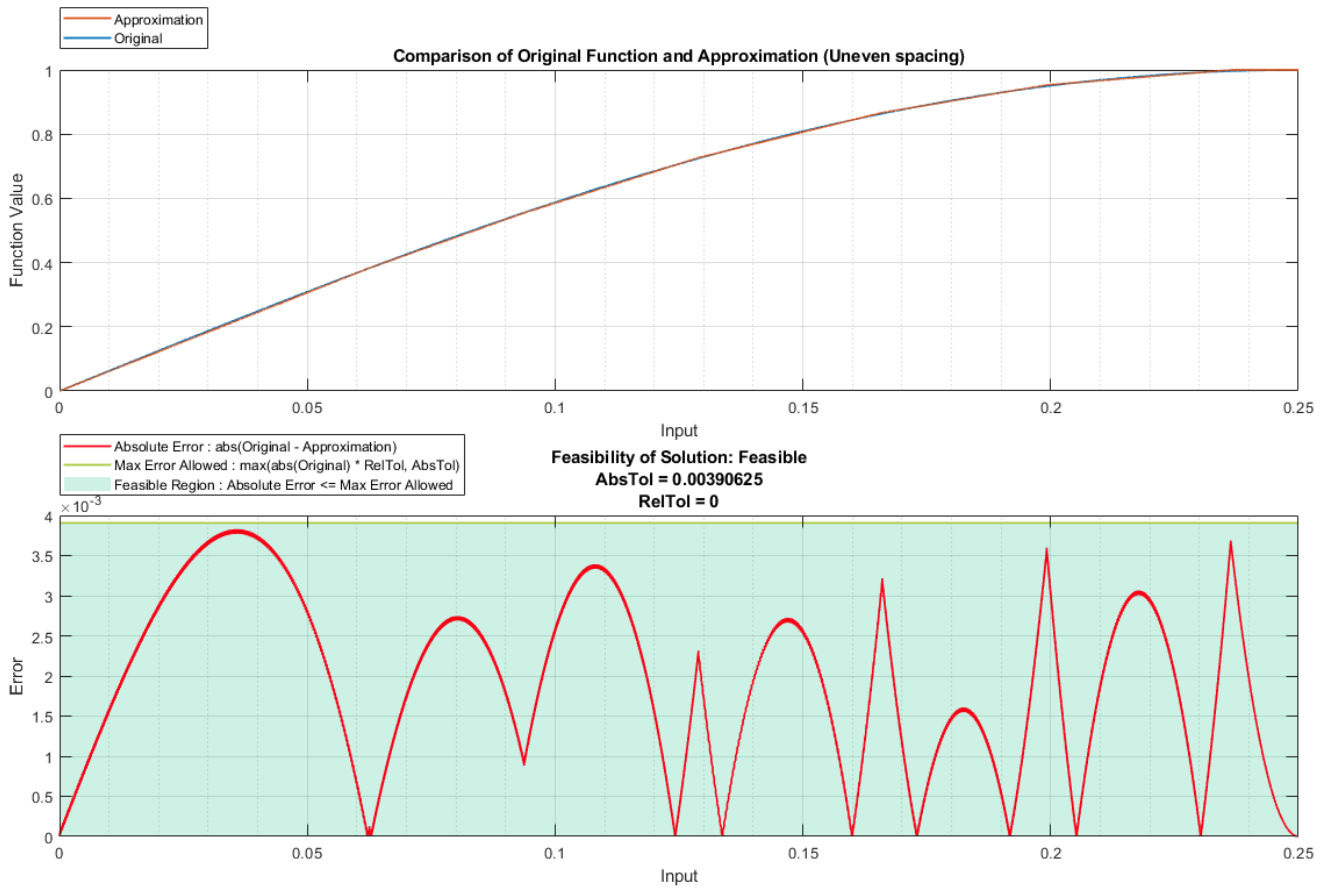
Compare the solution using even spacing to the original function.

```
[~, hEvenSpacing] = compare(bestEvenSpacingSolution);
hEvenSpacing.Children(4).Title.String = [hEvenSpacing.Children(4).Title.String ' (Even spacing)'];
```



Compare the solution using uneven spacing to the original function.

```
[~, hUnevenSpacing] = compare(bestUnevenSpacingSolution);
hUnevenSpacing.Children(4).Title.String = [hUnevenSpacing.Children(4).Title.String ' (Uneven spacing)']
```



Use the Approximation in a Simulink® Model

You can use this approximation directly in a Simulink® Lookup Table (n-D) block.

```
modelName = 'fxpdemo_approx';
open_system(modelName)

modelWorkspace = get_param(modelName, 'ModelWorkspace');

modelWorkspace.assignin('xevenFirstPoint' , xeven(1) );
modelWorkspace.assignin('xevenSpacing' , diff(xeven(1:2)) );
modelWorkspace.assignin('yeven' , yeven );
modelWorkspace.assignin('TableDTeven' , bestEvenSpacingSolution.TableData.TableDataType
modelWorkspace.assignin('BreakpointDTeven' , bestEvenSpacingSolution.TableData.BreakpointData

modelWorkspace.assignin('xuneven' , xuneven);
modelWorkspace.assignin('yuneven' , yuneven);
modelWorkspace.assignin('TableDTuneven' , bestUnevenSpacingSolution.TableData.TableDataType
modelWorkspace.assignin('BreakpointDTuneven' , bestUnevenSpacingSolution.TableData.BreakpointData

modelWorkspace.assignin('xpow2FirstPoint' , xpow2(1) );
modelWorkspace.assignin('xpow2Spacing' , diff(xpow2(1:2)) );
modelWorkspace.assignin('ypow2' , ypow2 );
modelWorkspace.assignin('TableDtpow2' , bestEvenPow2SpacingSolution.TableData.TableDataType
```

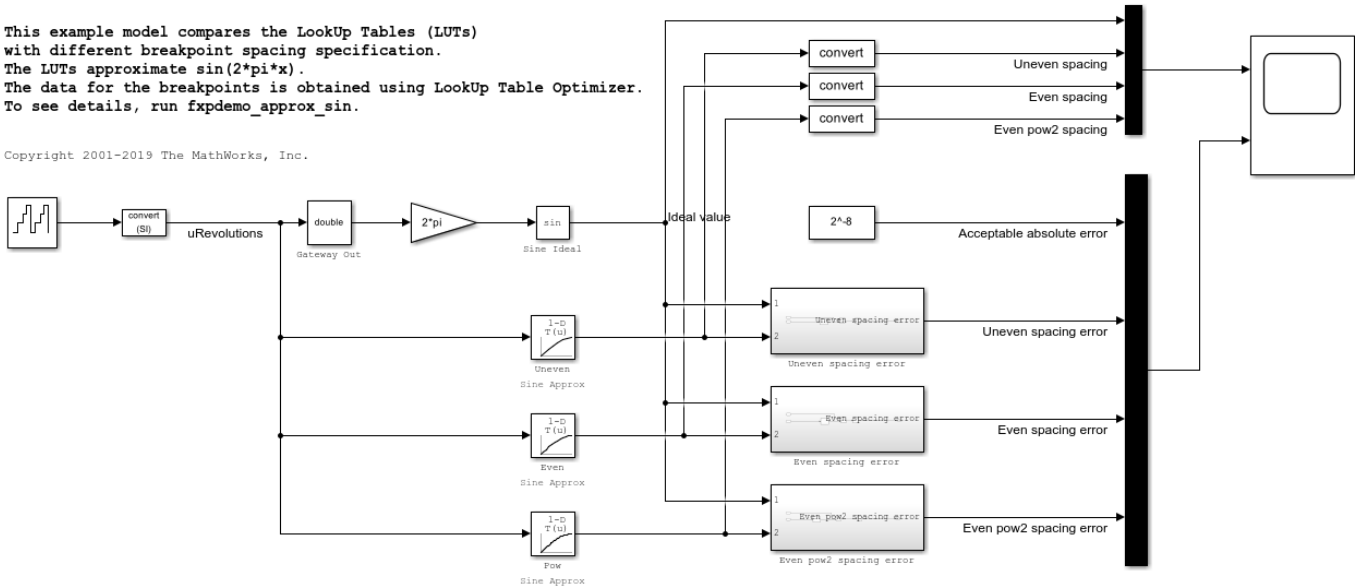
```

modelWorkspace.assignin('BreakpointDTpow2' , bestEvenPow2SpacingSolution.TableData.BreakpointT
set_param(modelName, 'Dirty', 'off');

```

This example model compares the LookUp Tables (LUTs) with different breakpoint spacing specification. The LUTs approximate $\sin(2\pi*x)$. The data for the breakpoints is obtained using LookUp Table Optimizer. To see details, run `fxpdemo_approx_sin`.

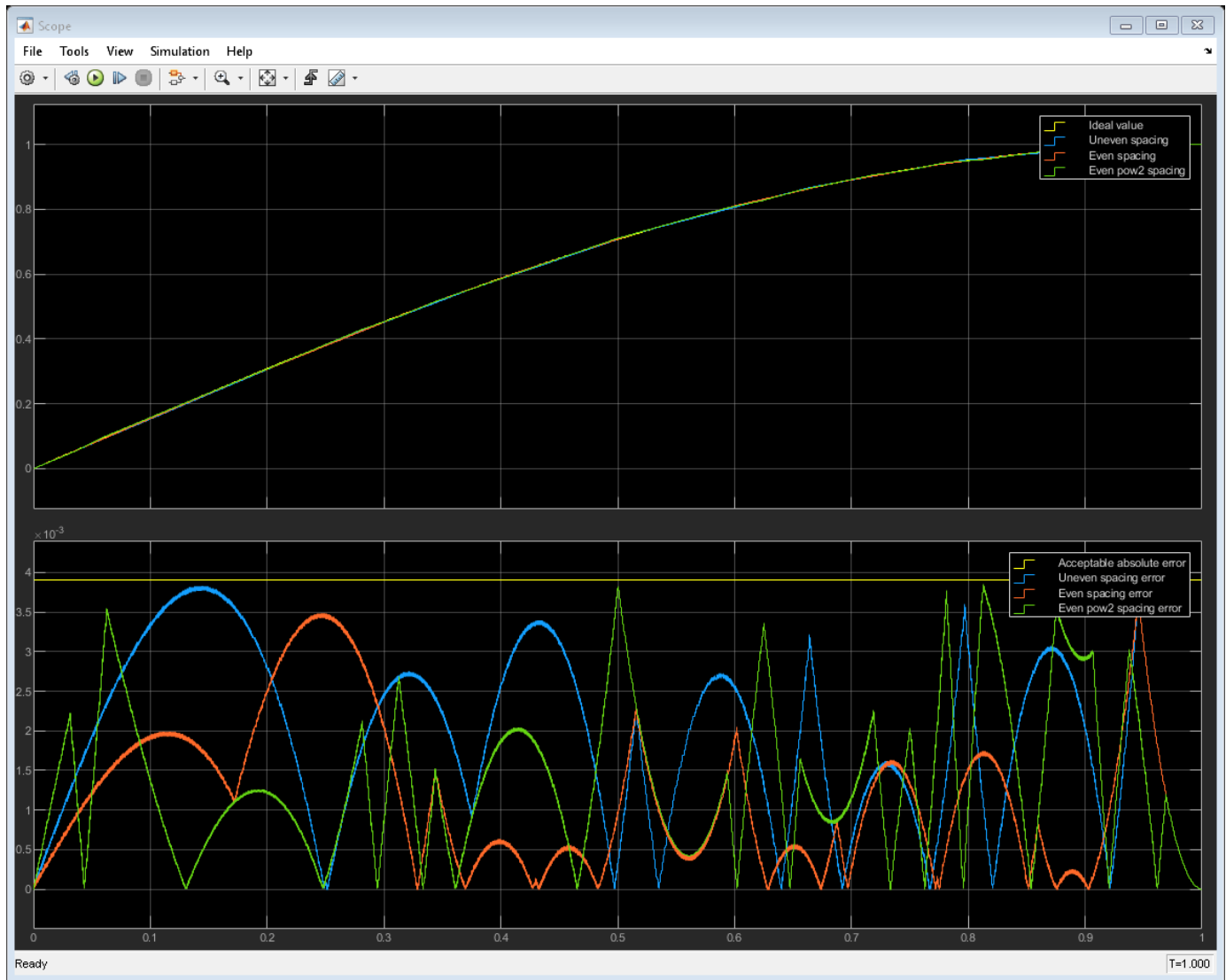
Copyright 2001-2019 The MathWorks, Inc.



Summary

The ideal function and the three approximations are used in the model `fxpdemo_approx`. If you have Simulink® Coder™ installed, you can generate code for the model. If inline parameters is ON, the generated code will show the large efficiency differences in the implementation of unevenly spaced, evenly spaced, and power of 2 spacing.

```
sim(modelName)
```



Clean up

```
close_system(modelName, 0);
```


Fixed-Point Conversion Using Fixed-Point Tool and Derived Range Analysis

This example shows how to use derived range analysis to collect ranges that then can be used by the Fixed-Point Tool to propose fixed-point scaling.

Overview

The Fixed-Point Tool can collect range data either by simulating the model and logging range information or by running derived range analysis. Range information obtained from simulation is completely dependent on the input to the simulation. Accurate range information can be obtained only if the inputs exercise the system over its full operating range. Unlike range collection using simulation logging, derived range analysis does not depend on sample inputs. Instead it uses formal methods to generate range information.

The benefits of derived range analysis over simulation logging are:

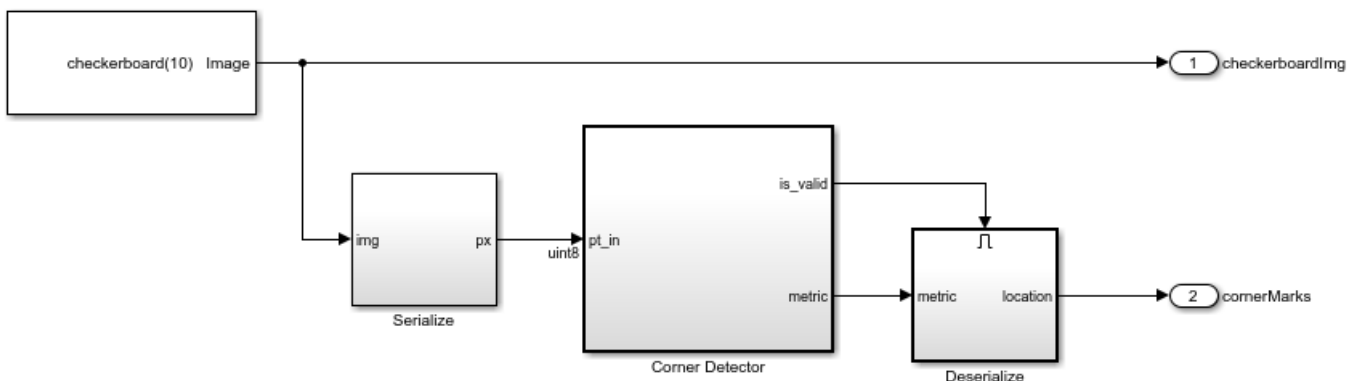
- Accurate results regardless of quality of input data
- Faster analysis times

Open and Inspect Model

Open the `ex_fxpdemo_corner_detection` model.

This model implements a corner detection algorithm based on the Harris corner detection method. The top level of the model includes blocks required to run the simulation. Note that the input to the simulation is a checkerboard image and the output is a 2 by n workspace matrix variable, `cornerMarks`, that contains locations of all found corners. The Corner Detector subsystem implements the algorithm. Within the Corner Detector subsystem, the Sobel Edge block applies Sobel operator to the input data, and the Corner Metric subsystem calculates the Harris corner metric. This example focuses on analyzing ranges of the Corner Metric subsystem using the Fixed-Point Tool.

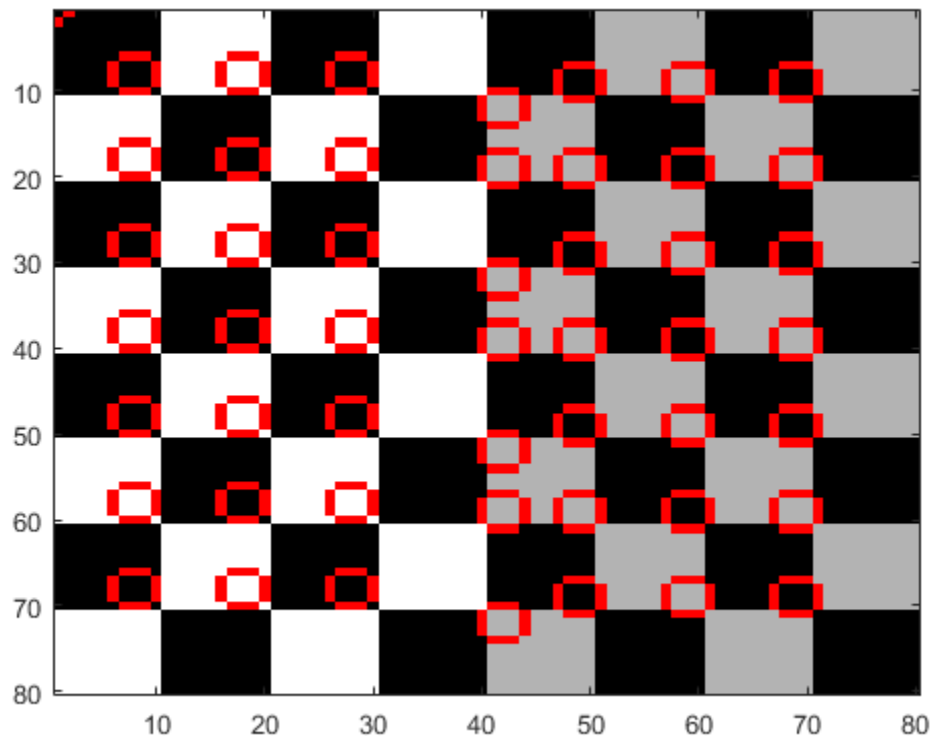
Note: Design ranges on inputs to the Corner Metric subsystem have been explicitly specified by the user prior to running the analysis. Those design ranges can be obtained either from running derived range analysis, simulating the model with logging on the system, or inferred from the model.



Simulate Model

Some blocks inside the Corner Metric subsystem have fixed-point types manually specified. However, some of the types were chosen poorly, and the model fails to detect corners correctly. Simulate the model by clicking the **Simulate** button. Observe that numerous overflows occur and that most of the corners are not marked.

In the Fixed-Point Tool, in the **Visualization of Simulation Data** pane, overflow markers indicate the blocks in the model that overflowed during simulation. The Fixed-Point Tool shows overflows occurred in the Gaussian Filter Accumulators.



Prepare Model For Conversion

To open the Fixed-Point Tool, right click on the Corner Metric subsystem and select Fixed-Point Tool .

Run Derived Range Analysis

The Fixed-Point Tool derives ranges by analyzing the generated code. Fixed-point data types generate more code and can make it harder for the analysis to derive accurate ranges. To improve the accuracy of the results, derive ranges of models using double-precision data types.

In the Fixed-Point Tool, in the **Collect Ranges** section of the tool-strip, click the **Derived Ranges** button to specify the range collection method. To begin the range analysis click the **Collect Ranges** button. This action overrides the data types in the model with double-precision types before performing the analysis to improve the accuracy of the results. After the analysis finishes, the Fixed-Point Tool displays collected range information.

Propose Fixed-Point Data Types

Range information obtained from derived range analysis can be used by the Fixed-Point Tool to propose fixed-point data types for blocks in the model. This can be done by clicking **Propose Data Types** button in the Fixed-Point Tool.

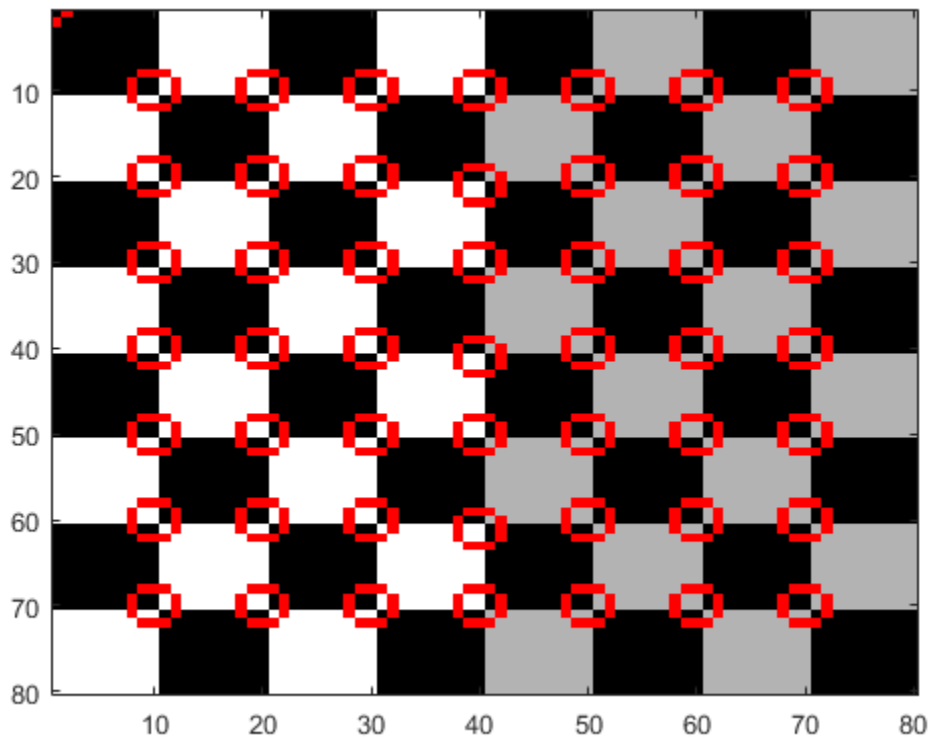
Apply Proposed Data Types

To apply the proposed data types, click the **Apply Data Types** button. By default, the Fixed-Point Tool applies all of the proposed data types. To apply a subset of the proposals, use the **Accept** check box to specify the proposals that you want to apply.

Verify Proposed Data Types

Proposed types should handle all possible inputs correctly. Set the model to use the newly applied types, simulate the model, and observe that all of the corners are now detected.

Note: Applying proposed data types updates the data type visualization and removes the corresponding overflow indicators

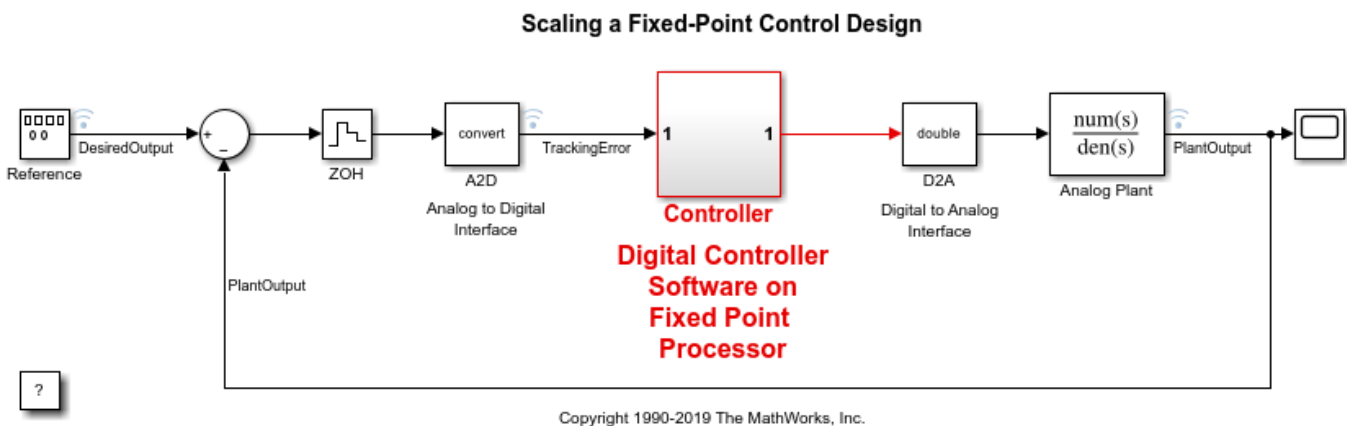


Fixed-Point Tool

This example uses the `fxpdemo_feedback` model to demonstrate how to use the Iterative Fixed-Point Conversion workflow in the Fixed-Point Tool. The iterative workflow in the tool automates common tasks of collecting min-max range data during simulations for use in automatically selecting fixed-point data types for blocks. You can configure any fixed-point capable block in Simulink® to accept the fixed-point data types proposed by the tool. You can manually specify data types for key blocks such as input ports and selectively disable automated scaling in a model on a block-by-block basis. This gives the Fixed-Point Tool more information to work with and results in automatically selected fixed-point scalings that are compatible with the key manually selected scalings.

Open the `fxpdemo_feedback` Model

The `fxpdemo_feedback` model contains a digital controller with initial guesses for fixed-point data types already specified in the Controller subsystem. In this example, you specify a known tolerance for the plant output, propose new fixed-point data types for the Controller subsystem, then apply the proposed data types to the model.



Open the Fixed-Point Tool

On the Simulink® **Apps** tab, under **Code Generation**, click the app icon.

Start the iterative Fixed-Point Conversion workflow. In the Fixed-Point Tool, click **New > Iterative Fixed-Point Conversion**.

Set Up the Model for Conversion to Fixed-Point

- 1 Select the subsystem that you want to convert. Under **System Under Design (SUD)**, select the Controller subsystem.
- 2 Choose the range collection method to use. Under **Range Collection Mode**, select **Simulation ranges**.
- 3 Specify **Simulation Inputs**. For this example, use the default model inputs for simulation.
- 4 Specify signal tolerances for logged signals. Set the **Absolute Tolerance** of the Plant Output signal to 0.1.

ITERATIVE FIXED-POINT CONVERSION

Settings | MATLAB Functions | Propose Data Types | Apply Data Types | Simulate with Embedded Types | Run to compare in SDI | Compare Results | Restore Original Model

WORKFLOW | PREPARE | COLLECT

Workflow Browser

Setup

Model Hierarchy

- Simulink Root
 - Data Objects
 - fxpdemo_feedback
 - Controller

System Under Design (SUD)

Select the system to analyze or convert.

Selected system under design: `fxpdemo_feedback/Controller`

- Simulink Root
 - fxpdemo_feedback
 - Controller

Range Collection Mode

Select whether to collect ranges through simulation or through static analysis that derives the ranges.

Simulation ranges
 Derived ranges
 Simulation with derived ranges

Simulation Inputs

Specify inputs for simulations. You can choose to use the current model inputs, or select a Simulink.SimulationInput object from the base workspace.

Simulation inputs:

Signal Tolerances

Specify tolerances for signals in your model that have signal logging enabled. After simulating with embedded types, the Workflow Browser displays whether the embedded run meets the specified signal tolerances.

Filter signal list:

| Signal Name | Absolute Tolerance | Relative Tolerance | Time Tolerance (seconds) |
|-----------------|--------------------|--------------------|--------------------------|
| TrackingError | | | |
| PlantOutput | 0.1 | | |
| In1:1 | | | |
| Combine Terms:1 | | | |
| PlantInput | | | |
| Prev Out:1 | | | |
| Up Cast:1 | | | |
| DesiredOutput | | | |

Prepare for Conversion to Fixed-Point

To prepare the model for fixed-point conversion, click **Prepare**. The Fixed-Point Tool creates a backup version of the model and checks the model for compatibility with the conversion process. For more information about preparation checks, see “Use the Fixed-Point Tool to Prepare a System for Conversion” on page 38-2.

Workflow Browser: Setup, Preparation Results

Model Hierarchy: Simulink Root, Data Objects, fxpdemo_feedback, Controller

Selected system under design: fxpdemo_feedback/Controller

Select a result below for more information

| Selection | Check | Status |
|----------------------------------|-------------------------------------|--------|
| <input checked="" type="radio"/> | Create Restore Point | ✓ |
| <input type="radio"/> | Hardware Implementation Consistency | ✓ |
| <input type="radio"/> | Diagnostic Settings | ✓ |
| <input type="radio"/> | Unsupported Constructs | ✓ |
| <input type="radio"/> | System Under Design Boundary | ✓ |

Preparation is complete for the selected system under design

Progress: 100%

Preparation Details

Check Details

To ensure your original design is saved before making fixed-point data type changes, create a restore point for the model.

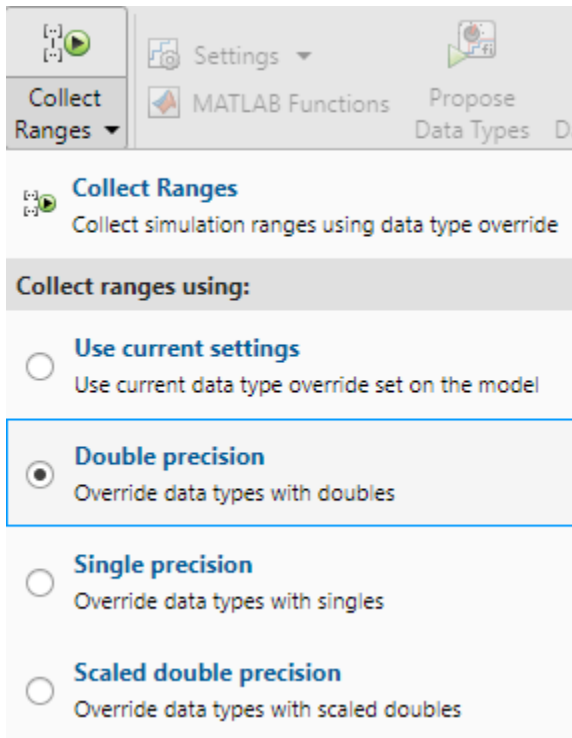
Check Status

A restore point was created for the model. To restore the model to this state, click the **Restore Original Model** button.

- fxpdemo_feedback

Collect Ranges

To collect ranges, click the **Collect Ranges** button arrow and select **Double precision**. Click **Collect Ranges** to start the range collection simulation.



When you select **Double precision** as the range collection mode, the tool simulates the `fxpdemo_feedback` model with data type override enabled. Data type override performs a global override of the fixed-point data types in the model, thereby avoiding quantization effects. This enables you to establish an ideal floating-point baseline for the behavior of your model.

The results of range collection are stored in `BaselineRun`. The **Results** spreadsheet displays a summary of the statistics collected during the range collection simulation, including the currently specified data types on the model (**SpecifiedDT**), simulation minimum, and simulation maximum values. The compiled data type (**CompiledDT**) column displays `double` for all objects in the `Controller` subsystem, indicating that data type override was applied during the range collection simulation.

You can click on any result to view additional details in the **Results Details** pane. The **Visualization of Simulation Data** pane displays a summary of histograms of the bits used by each object in your model. You can customize the information displayed in the **Results** spreadsheet, or use the **Explore** tab to sort and filter these results based on additional criteria. For more information, see “Control Views in the Fixed-Point Tool” on page 39-13.

ITERATIVE FIXED-POINT CONVERSION EXPLORE

Settings MATLAB Functions Propose Data Types Apply Data Types Simulate with Embedded Types Run to compare in SDI Compare Results Restore Original Model

Workflow Browser: Setup, Preparation Results, BaselineRun

Model Hierarchy: Simulink Root, Data Objects, fxdemo_feedback, Controller

| Name | CompiledDT | SpecifiedDT | SimMin | SimMax |
|------------------------------------|------------|---------------------------|--------------------|--------------------|
| Combine Terms : Accumulator | double | Inherit: Inherit via i... | -6.475416336873... | 4.32700180757929 |
| Combine Terms : Output | double | fixdt(1,32,12) | -2.413500903789... | 4.32700180757929 |
| Denominator Terms : Accumulator | double | fixdt(1,32,12) | -8.516638478410... | 5.39648751512156 |
| Denominator Terms : Output | double | fixdt(1,32,12) | -6.475416336873... | 3.4877081684371... |
| Denominator Terms : Product output | double | fixdt(1,32,12) | -8.516638478410... | 5.39648751512156 |
| Down Cast | double | fixdt(1,16,5) | -2.413500903789... | 4.32700180757929 |
| Numerator Terms : Accumulator | double | fixdt(1,32,12) | -5.677304459288... | 5.700524518426912 |
| Numerator Terms : Output | double | fixdt(1,32,12) | -3.367372640959... | 3.5439615259925... |
| Numerator Terms : Product output | double | fixdt(1,32,12) | -5.677304459288... | 5.700524518426912 |
| Up Cast | double | fixdt(1,16,14) | -2 | 3.9999999999711... |

Result Details: fxdemo_feedback/Controller/Up Cast

Needs Attention: There are overflows associated with this result.

| Property | Specified Data Type |
|-----------|---------------------|
| Data Type | fixdt(1,16,14) |
| Minimum | -2 |
| Maximum | 1.99993896484375 |
| Precision | 6.103515625e-05 |

| Property | Minimum | Maximum |
|------------|---------|-----------------|
| Simulation | -2 | 3.9999999999... |

Simulation Data Overview using fixdt(1,16,14)

| Values | Potential Overflows | In-Range | Potential Underflows |
|----------|---------------------|----------|----------------------|
| Positive | 9 | 51 | 139 |
| Negative | 0 | 65 | 135 |
| Zero | 0 | 0 | 0 |

Histograms of all results in the model

Propose Fixed-Point Data Types

To configure the settings to use for data type proposals, expand the **Settings** button arrow. For this example, use the default settings.

Settings

PROPOSE

Propose: Fraction Length

Propose signedness: Yes

Safety margin for simulation min/max (%): 2

CONVERT TO FIXED POINT

Convert double/single/half types: Yes

Convert inherited types: Yes

Default word length: 16

Default fraction length: 4

| <u>Original Data Type</u> | <u>Word Length</u> | <u>Fraction Length</u> |
|---------------------------|--------------------|------------------------|
| Double/Single/Half | → 16 | Will propose |
| Inherited | → 16 | Will propose |
| Fixed point | → No change | Will propose |

To propose data types based on the ranges collected and the data type proposal settings specified, click **Propose Data Types**. The tool uses all available range data to calculate data type proposals which can include design minimum or maximum values, simulation minimum or maximum values, and derived minimum or maximum values. Data types are proposed for all objects in the system under design whose **Lock output data type setting against changes by the fixed-point tools** parameter is cleared.

The **Results** spreadsheet updates to show the proposed data types in the **ProposedDT** column. The Fixed-Point Tool allows you to selectively apply data type proposals to objects in your model. In the spreadsheet, use the **Accept** check boxes to specify the proposed data types that you want to apply to your model. By default, the app accepts all data type proposals which differ from the currently specified data types. For this example, use the default.

The screenshot displays the Fixed-Point Designer interface during the 'Propose Data Types' workflow. The 'Results' table is as follows:

| Name | CompiledDT | SpecifiedDT | ProposedDT | Accept | SimMin | SimMax |
|-------------------------------|------------|----------------------|----------------|--------|-----------------|-----------------|
| Combine Terms : Accumulator | double | Inherit: Inherit ... | n/a | | -6.475416336... | 4.3270018075... |
| Combine Terms : Output | double | fixdt(1,32,12) | fixdt(1,32,28) | ✓ | -2.413500903... | 4.3270018075... |
| Denominator Terms : Accu... | double | fixdt(1,32,12) | fixdt(1,32,27) | ✓ | -8.516638478... | 5.3964875151... |
| Denominator Terms : Output | double | fixdt(1,32,12) | fixdt(1,32,28) | ✓ | -6.475416336... | 3.4877081684... |
| Denominator Terms : Produ... | double | fixdt(1,32,12) | fixdt(1,32,27) | ✓ | -8.516638478... | 5.3964875151... |
| Down Cast | double | fixdt(1,16,5) | fixdt(1,16,12) | ✓ | -2.413500903... | 4.3270018075... |
| In1 | | Inherit: auto | fixdt(1,8,4) | ✓ | | |
| Numerator Terms : Accumul... | double | fixdt(1,32,12) | fixdt(1,32,28) | ✓ | -5.677304459... | 5.7005245184... |
| Numerator Terms : Output | double | fixdt(1,32,12) | fixdt(1,32,29) | ✓ | -3.367372640... | 3.5439615259... |
| Numerator Terms : Product ... | double | fixdt(1,32,12) | fixdt(1,32,28) | ✓ | -5.677304459... | 5.7005245184... |
| Out1 | | Inherit: auto | n/a | | | |
| Up Cast | double | fixdt(1,16,14) | fixdt(1,16,12) | ✓ | -2 | 3.999999999... |

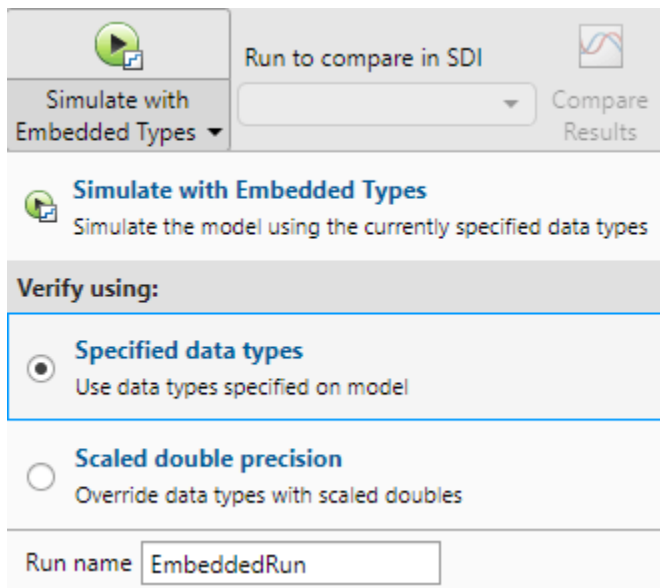
The 'Visualization of Simulation Data' section shows histograms for all results. The legend indicates the following counts for the 'Up Cast' variable:

| Values | Potential Overflows | In-Range | Potential Underflows |
|----------|---------------------|----------|----------------------|
| Positive | 0 | 48 | 151 |
| Negative | 0 | 49 | 151 |
| Zero | 0 | 0 | 0 |

Apply Fixed-Point Data Types to the Model and Verify New Settings

To write the proposed data types to the model, click **Apply Data Types**. The tool updates the **SpecifiedDT** column to show that the data types have been applied to the model.

Simulate the model using the applied fixed-point data types. Expand the **Simulate with Embedded Types** button arrow and select **Specified data types**. Then click **Simulate with Embedded Types**.



Next, simulate the model using the fixed-point data types currently specified on the model. Expand the **Simulate with Embedded Types** button arrow and select **Specified data types**, then click **Simulate with Embedded Types**.

The Fixed-Point Tool simulates the model using the new fixed-point data types and logs minimum and maximum values and overflow data for all objects in the system under design. This information is stored in a new run named **EmbeddedRun**. The icon next to **EmbeddedRun** displays a pass status, indicating that all signals in the system under design meet the specified tolerances. The **Visualization of Simulation Data** pane updates to display the new **EmbeddedRun** data.

The screenshot displays the MATLAB Fixed-Point Designer interface during a simulation comparison. The top toolbar shows the 'Run to compare in SDI' dropdown menu set to 'EmbeddedRun'. The 'Results' table lists various simulation components and their data types. The 'Visualization of Simulation Data' section shows a histogram of results for the 'Up Cast' signal, with a legend indicating 'Overflows', 'Representable', 'In-Range', and 'Underflows'. The 'Result Details' panel for 'fxpdemo_feedback/Controller/Up Cast' provides specific property values for the specified data type.

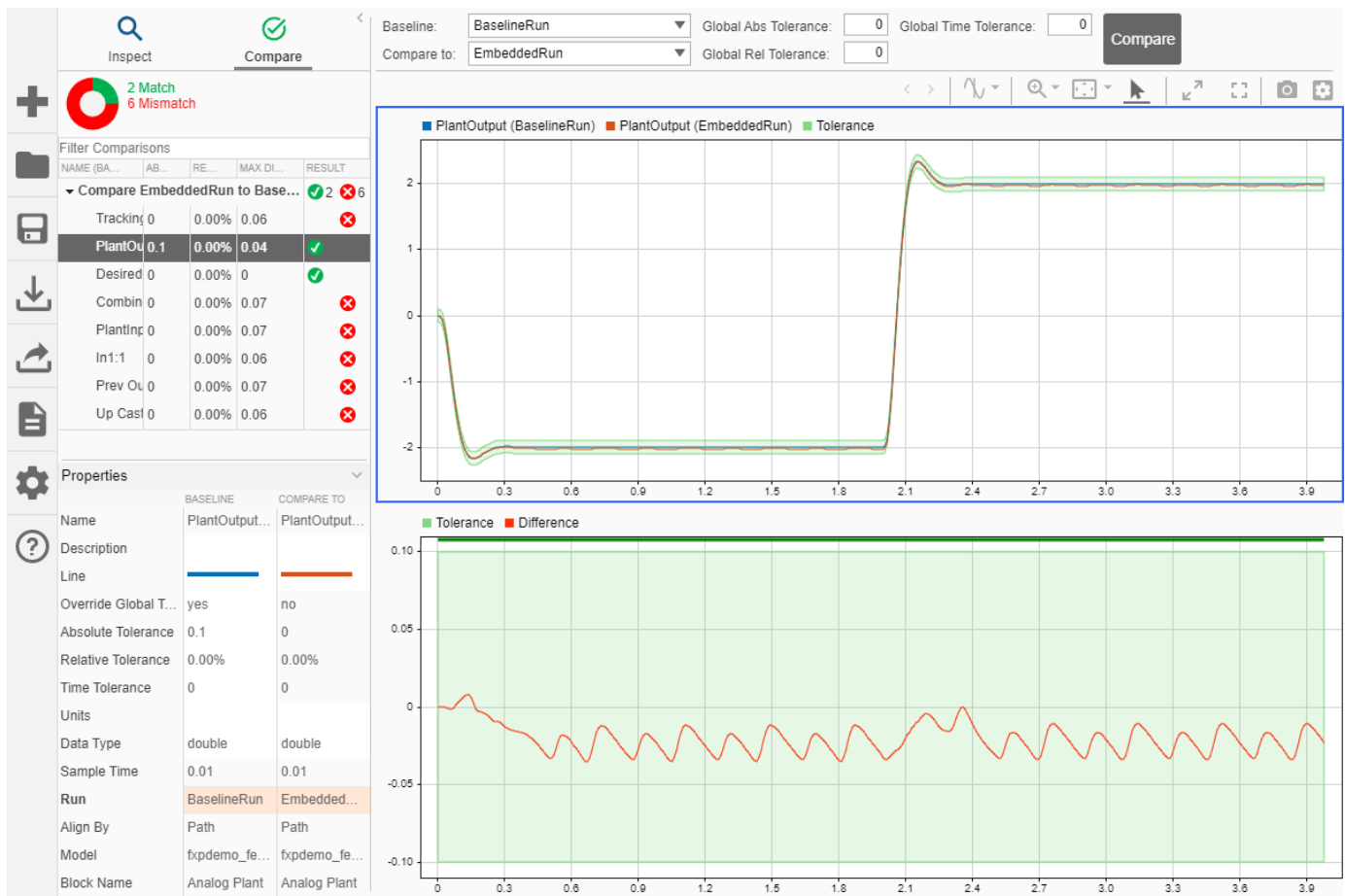
| Property | Specified Data Type |
|-----------|---------------------|
| Data Type | fixdt(1,16,12) |
| Minimum | -8 |
| Maximum | 7.999755859375 |
| Precision | 0.000244140625 |

| Property | Minimum | Maximum |
|------------|---------|---------|
| Simulation | -2 | 4.0625 |

| Values | Potential Overflows | In-Range | Potential Underflows |
|----------|---------------------|----------|----------------------|
| Positive | 0 | 85 | 0 |
| Negative | 0 | 25 | 0 |
| Zero | 0 | 289 | 0 |

To compare the ideal results stored in **BaselineRun** with the newly applied fixed-point data types, select **EmbeddedRun** from the **Run to compare in SDI** drop down menu. Then click **Compare Results** to open the Simulation Data Inspector.

In the Simulation Data Inspector, select **PlantOutput** as the signal to compare.



The plot of the plant output signal for EmbeddedRun is within the specified tolerance band.

If the behavior of the converted system does not meet your requirements or if you wish to explore the effect of additional data type selections, you can propose new data types after applying new proposal settings. Continue iterating until you find settings for which the fixed-point behavior of the system is acceptable.

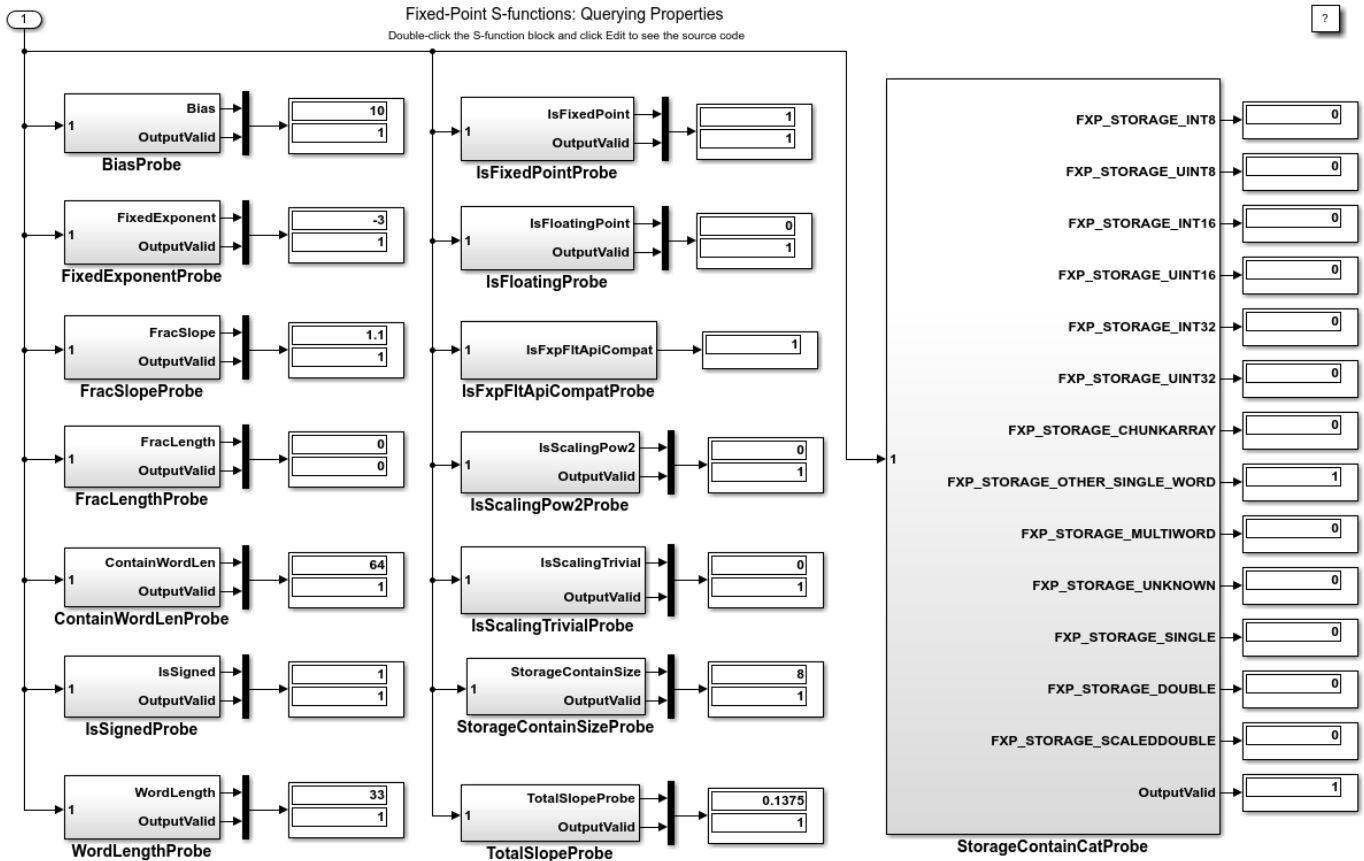
After the conversion process, if you want to restore your model to its state at the start of the conversion process, click **Restore Original Model**. Any changes made to your model after the preparation stage of conversion are removed.

See Also

Fixed-Point Tool

Fixed-Point S-Functions: Querying Properties

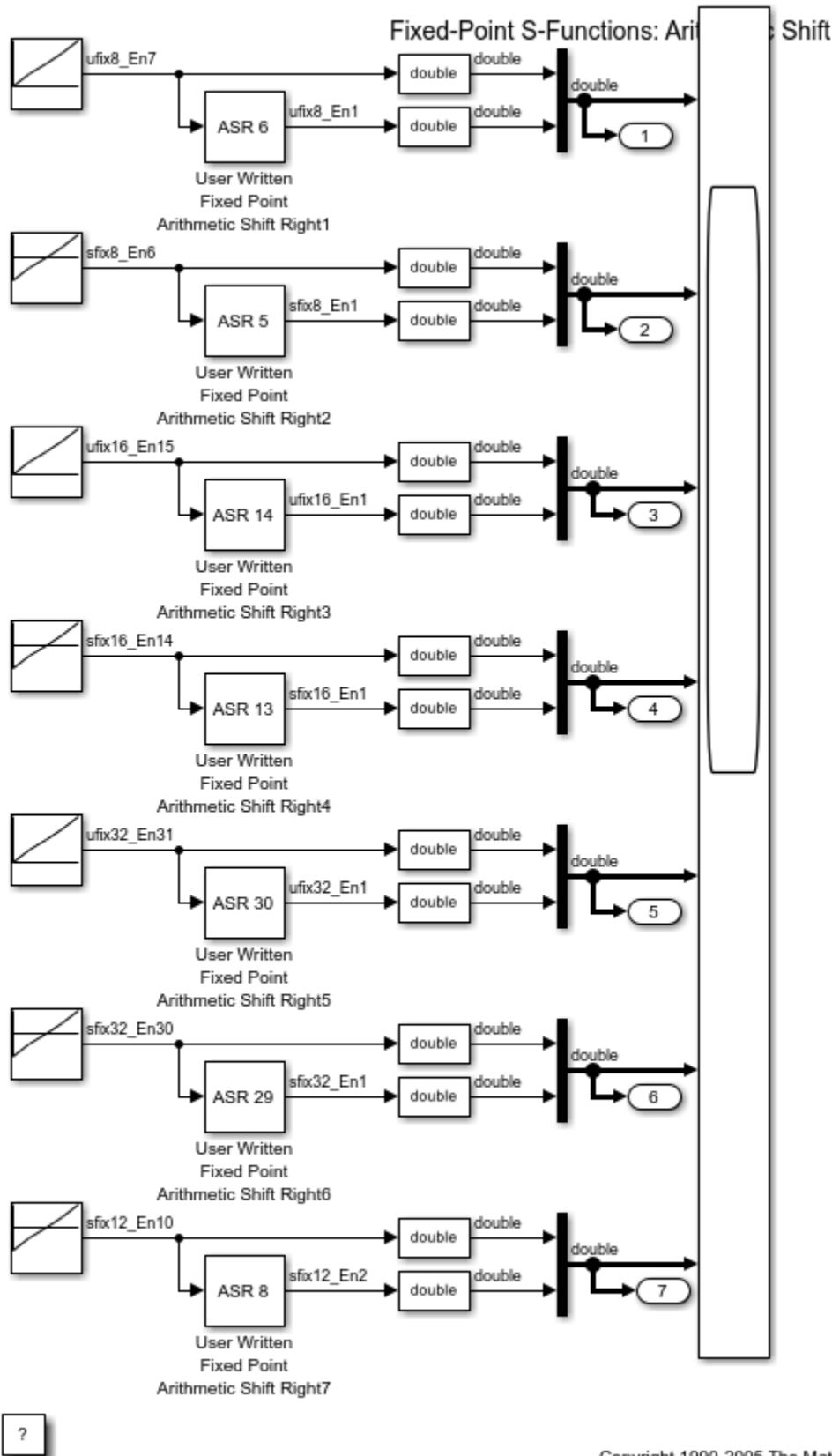
This example shows several ways to use S-functions to probe signal properties. Use this model for simulation only; it does not support code generation.

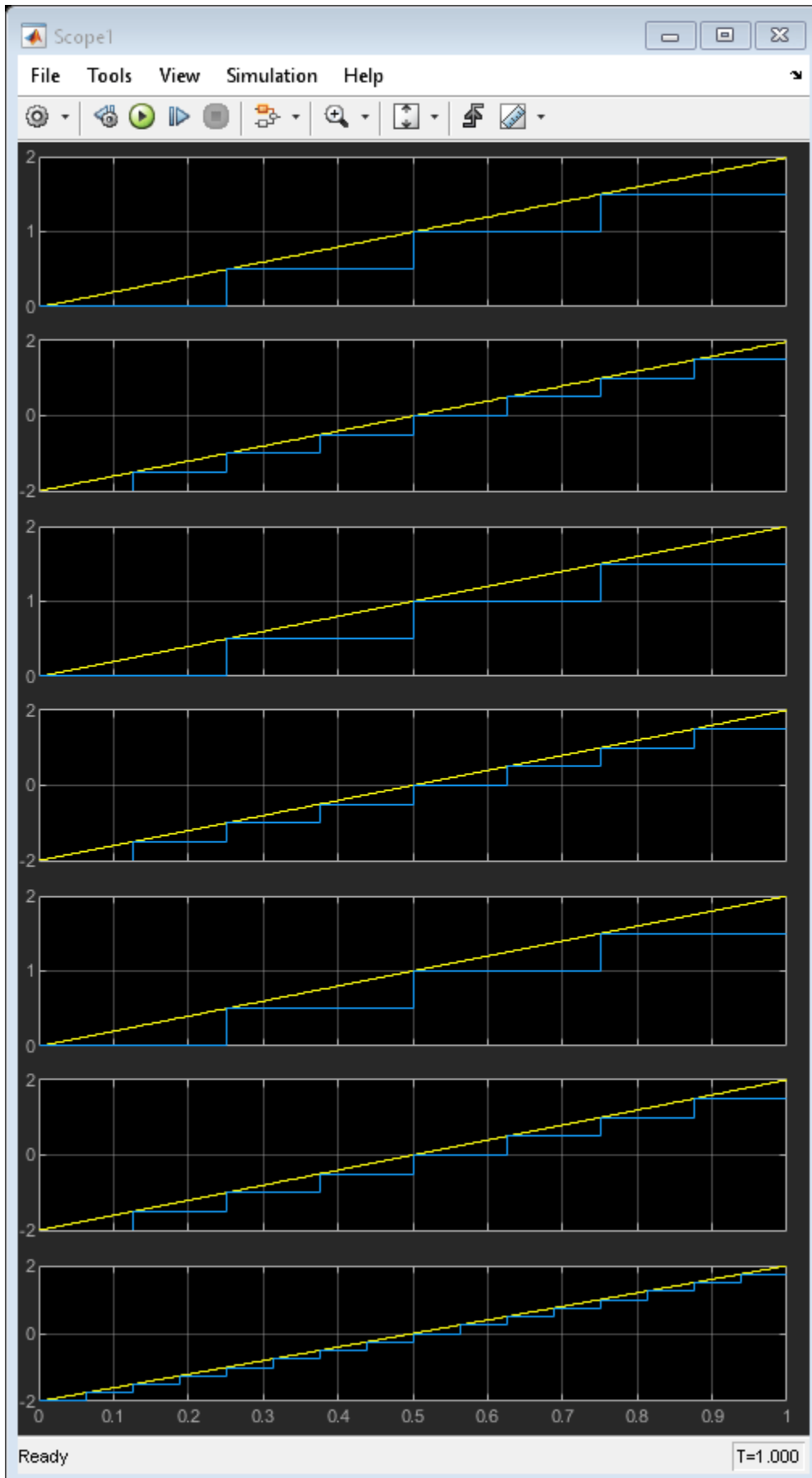


Fixed-Point S-Functions: Arithmetic Shift

This example shows a custom C language S-function written to perform an arithmetic shift. This operation is available in Simulink® with the "Shift Arithmetic" block, which can be used for comparison with this S-function example.

To see the source code for the S-function, use the right-click context menu to select "Look under mask". When the dialog box appears, press the Edit button.

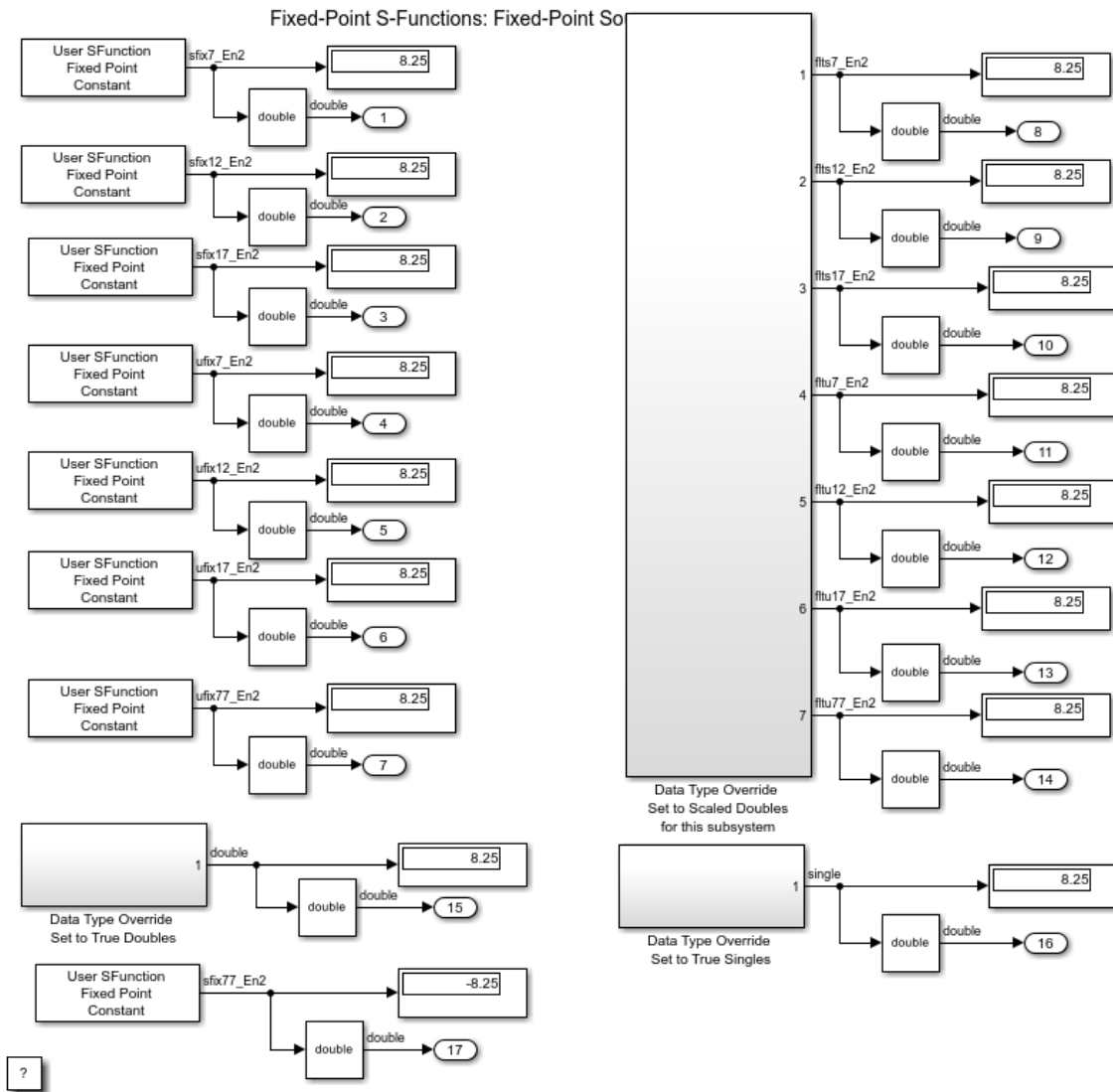




Fixed-Point S-Functions: Fixed-Point Source

This example shows a custom C language S-function written to generate a constant value. This operation is available in Simulink® with the "Constant" block, which can be used for comparison with this S-function example.

To see the source code for the S-function, use the right-click context menu to select "Look Under Mask". When the dialog box appears, press the Edit button.

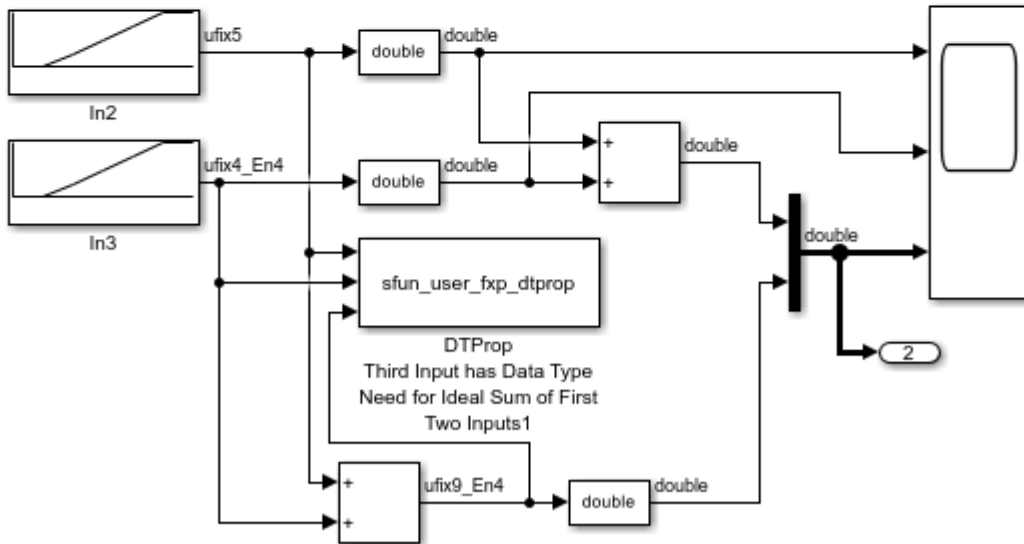
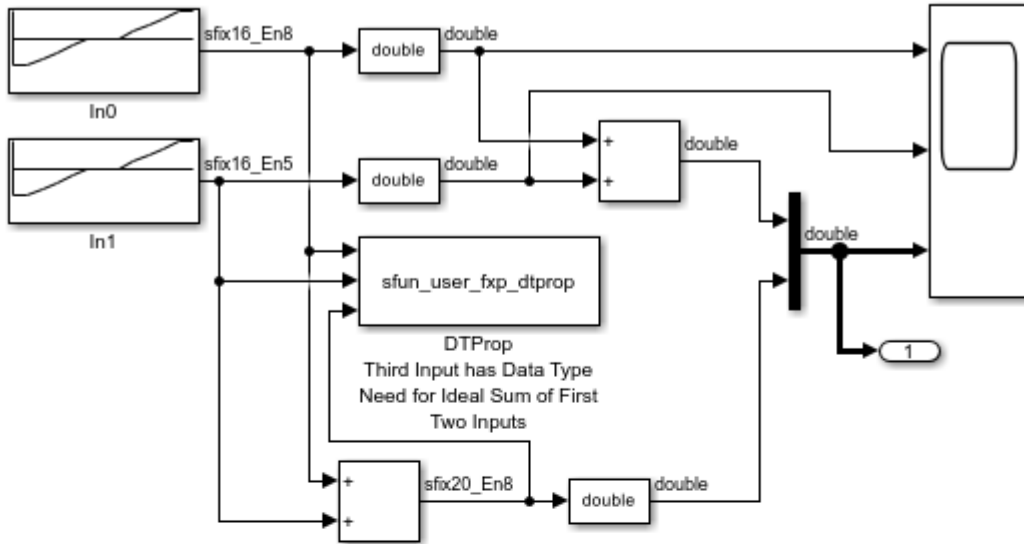


Fixed-Point S-Functions: Data Type Propagation

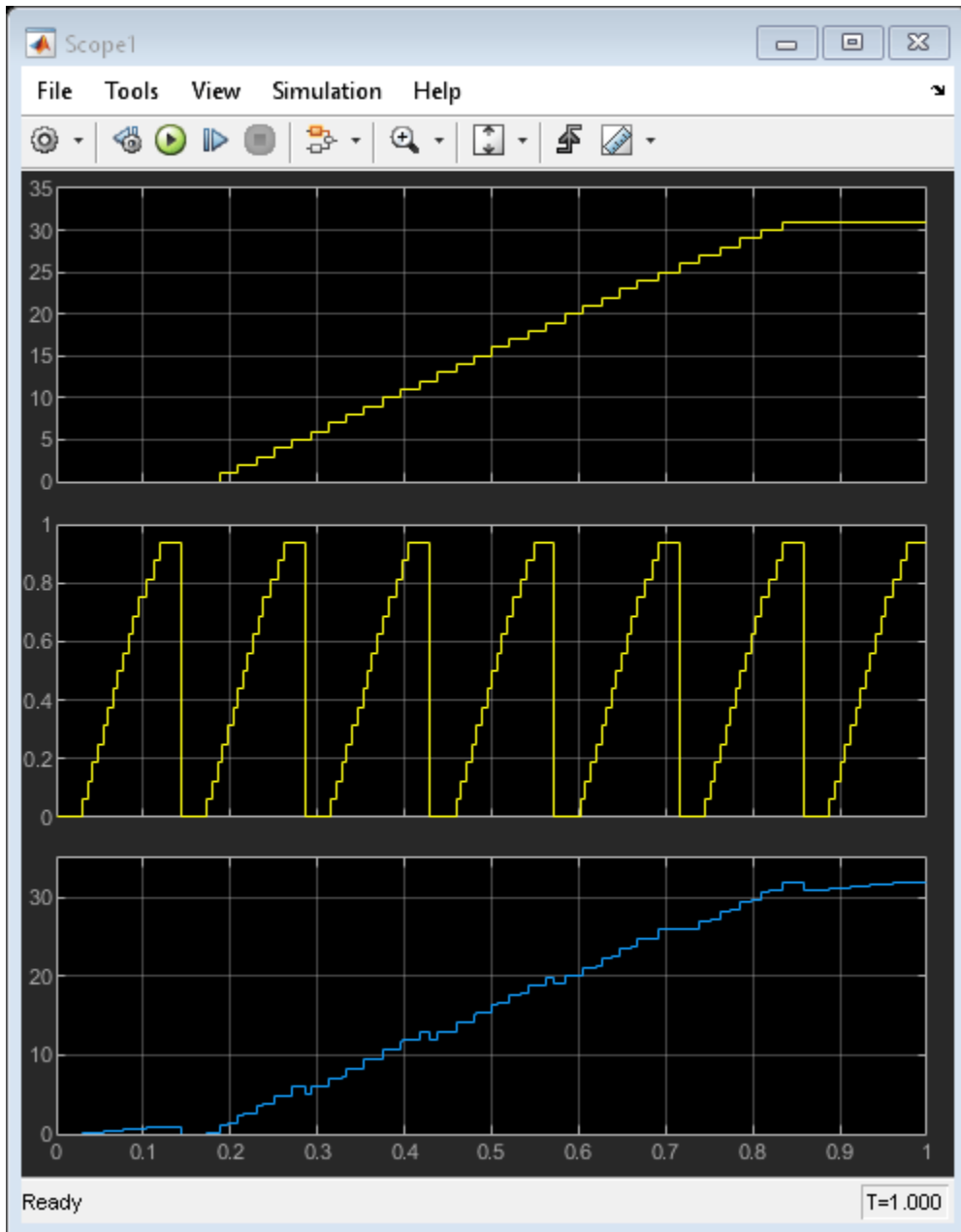
This model shows how to propagate fixed-point data types in fixed-point S-Functions. It exercises a custom C language S-function written to enforce data types across multiple signals. This operation is available in Simulink® with the "Data Type Propagation" block, which can be used for comparison with this S-function example.

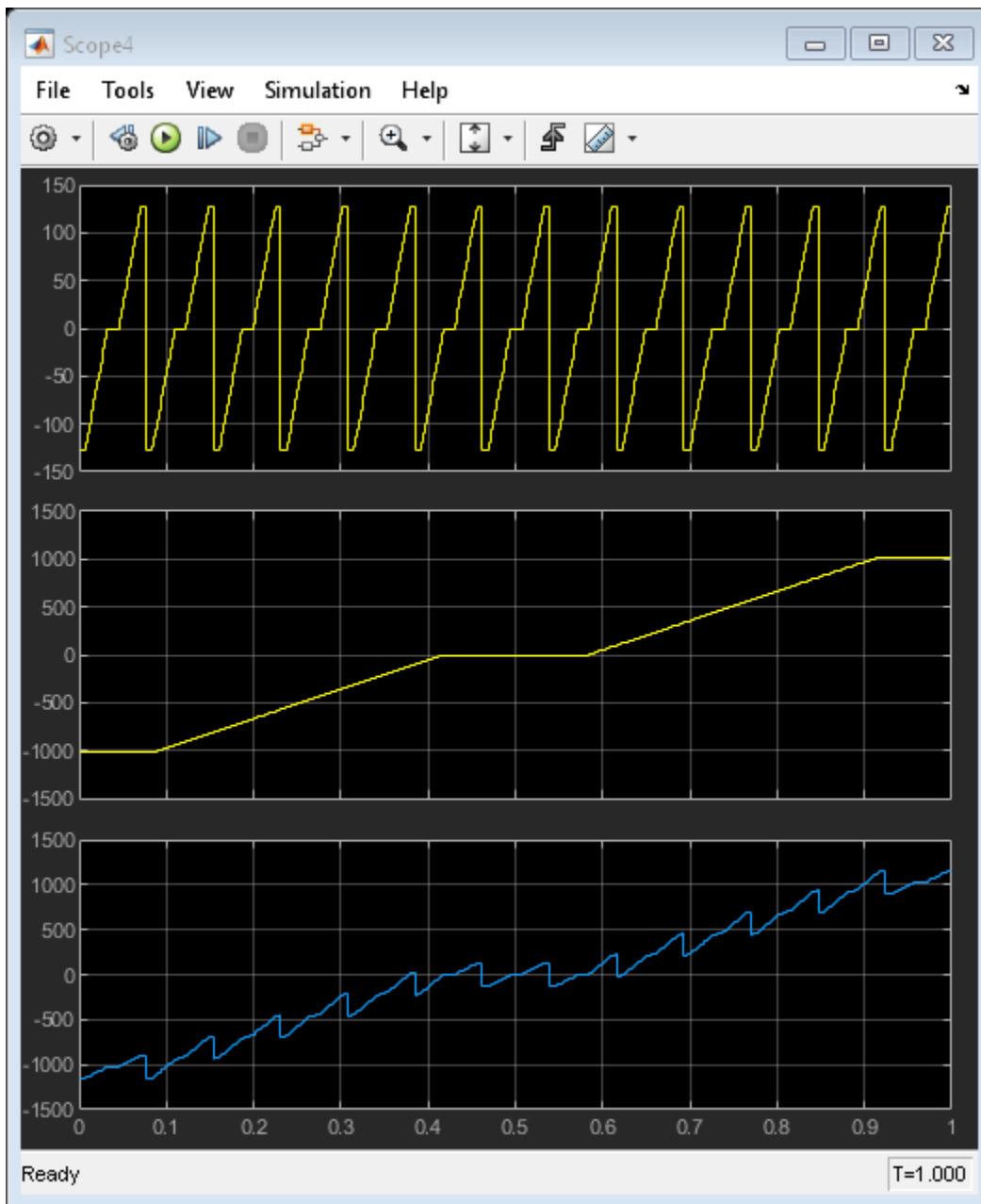
To see the source code for the S-function, use the right-click context menu to select "Block Parameters". When the dialog box appears, press the Edit button.

Fixed-Point S-Functions: Data Type Propagation



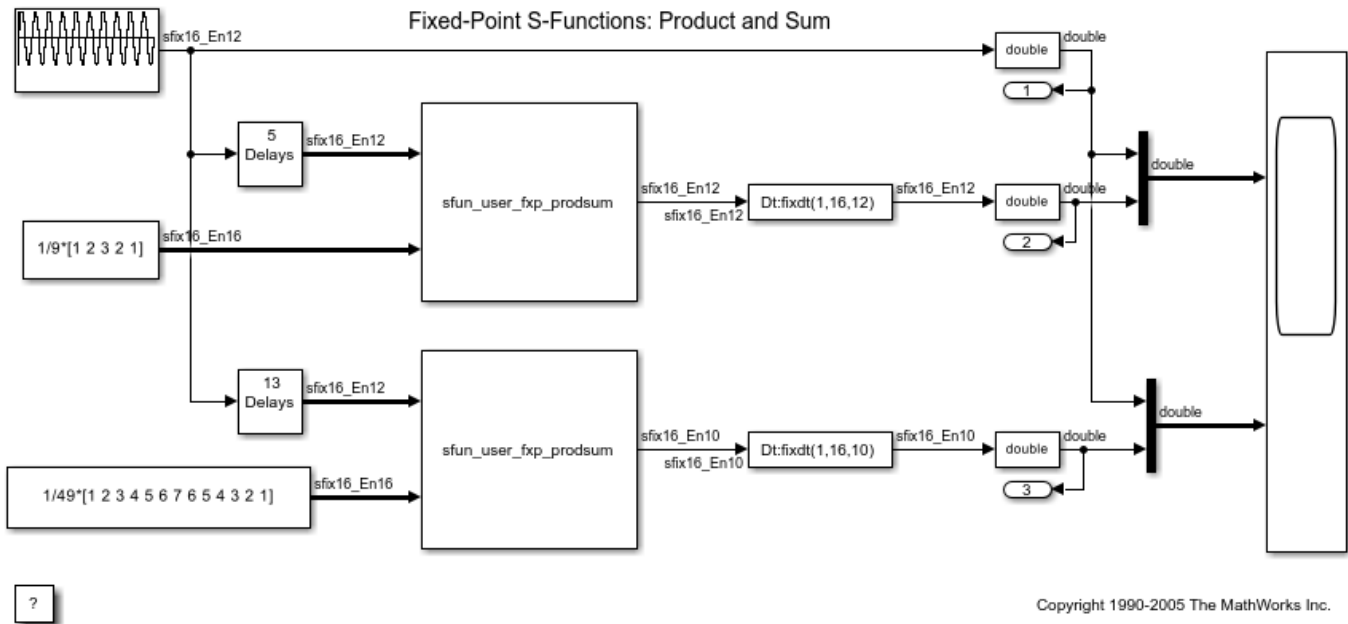
?



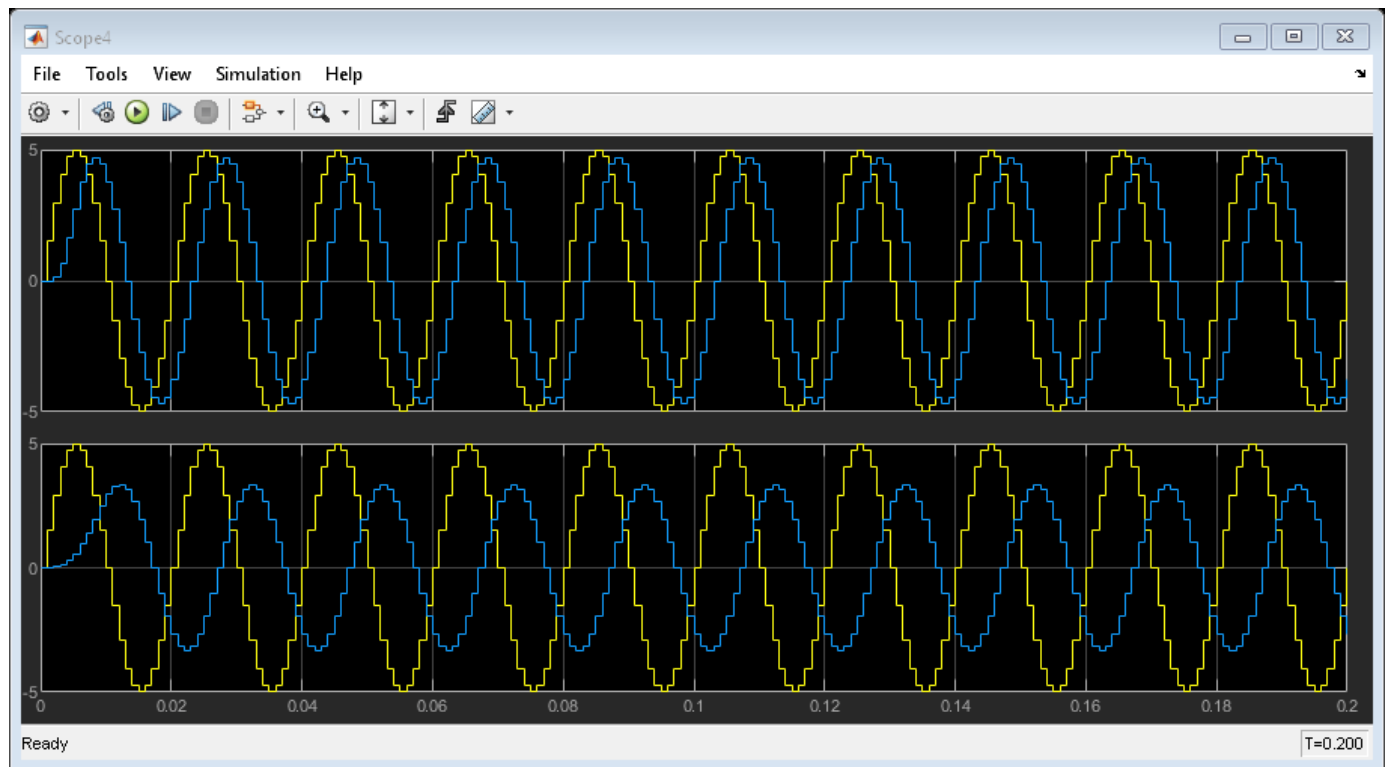


Fixed-Point S-Functions: Product and Sum

This model shows how to exercise a custom C language S-function written to compute a fixed-point "product and sum" operation. To see the source code for the S-function, use the right-click context menu to select "Block Parameters". When the dialog box appears, press the Edit button.



Copyright 1990-2005 The MathWorks Inc.



How to Use HDL Optimized Normalized Reciprocal

This example shows how and when to use the `normalizedReciprocal` function and the Normalized Reciprocal HDL Optimized block to compute the normalized reciprocal of an input.

The reciprocal of a value can have a large range, and fixed-point types have limited range compared to floating-point types. If an input value u is small, then $1/u$ is large, and if u is large, then $1/u$ is small. Therefore, a fixed-point type for $1/u$ must have high precision and large range which requires a large word length.

In applications where the range of a product of a reciprocal and another variable is known, then it is efficient to compute a normalized reciprocal, multiply it by the other variable, and then apply a shift to the final output. For example, this is the case in least-squares applications where a division by a pivot element is required in back-substitution.

Calculating the Normalized Reciprocal

Given an input u , Normalized Reciprocal computes the normalized reciprocal, y , and exponent, e , such that

$$(2.^e) .* y = 1./u,$$

and

$$0.5 < |y| \leq 1.$$

If $u = 0$ and u is fixed-point or scaled-double, then $y = 2^{-\text{eps}(y)}$.

If $u = 0$ and u is a floating-point type, then $y = \text{inf}$.

If $u \sim 0$, this function returns the equivalent of

$$\begin{aligned} [y, e] &= \text{log2}(1./\text{abs}(\text{double}(u))) \\ y(u < 0) &= -y(u < 0) \end{aligned}$$

except that it is computed using only shifts and adds.

Choose the MATLAB Function or the Simulink Block

For C code generation and system design, use the MATLAB function `normalizedReciprocal`. This function does not compute with latency. For simulation, compile the function into a MEX file for speed using `fiaccel`, `buildInstrumentedMex`, or `codegen`.

To generate optimized HDL code, use the Normalized Reciprocal HDL Optimized block. This block is optimized for high throughput and small area in HDL, and simulates with the same latency present in the generated HDL code.

The block and function produce identical numerical outputs.

Compute the Normalized Reciprocal Using MATLAB

Compute the normalized reciprocal of a fixed-point input, u , then compare this value to the actual value of the reciprocal.

```
u = fi([-pi, 0.01, pi])
```

```
u =  
    -3.1416    0.0100    3.1416  
        
      DataTypeMode: Fixed-point: binary point scaling  
      Signedness: Signed  
      WordLength: 16  
      FractionLength: 13  
  
[y,e] = normalizedReciprocal(u)  
  
y =  
    -0.6367    0.7806    0.6367  
        
      DataTypeMode: Fixed-point: binary point scaling  
      Signedness: Signed  
      WordLength: 16  
      FractionLength: 14  
  
e =  
  
    1x3 int32 row vector  
  
    -1     7    -1  
  
computed_reciprocal = 2.^double(e).* double(y)  
  
computed_reciprocal =  
    -0.3183    99.9141    0.3183  
  
actual_reciprocal = 1./double(u)  
  
actual_reciprocal =  
    -0.3183    99.9024    0.3183
```

You can see that the normalized reciprocal and the actual reciprocal are close in value.

Define Inputs for Simulink Model

Define a fixed-point input, *u*, to take the normalized reciprocal of using the Normalized Reciprocal HDL Optimized block.

```
x = linspace(0.001,100,100);  
x = [fliplr(-x),x];  
u = fi(x,1,18);
```

Latency of the Normalized Reciprocal HDL Optimized block

The Normalized Reciprocal HDL Optimized block works by normalizing the input using a binary search, which has a latency of approximately \log_2 of the word length of the input, followed by a

CORDIC reciprocal kernel, which has a latency approximately the same as the word length of the input.

The Normalized Reciprocal HDL Optimized block is always ready to accept data. After the initial latency, valid samples are output every sample. The latency in samples for a fixed-point input u is

$$D = \text{ceil}(\log_2(u.\text{WordLength})) + u.\text{WordLength} + 5$$

You can use the function `normalizedReciprocalLatency`, included in this example, to compute the latency for inputs with fixed point, double, or single numeric types.

To align the input samples with the output of the Normalized Reciprocal HDL Optimized block, use the latency D in a delay.

```
D = normalizedReciprocalLatency(u)
```

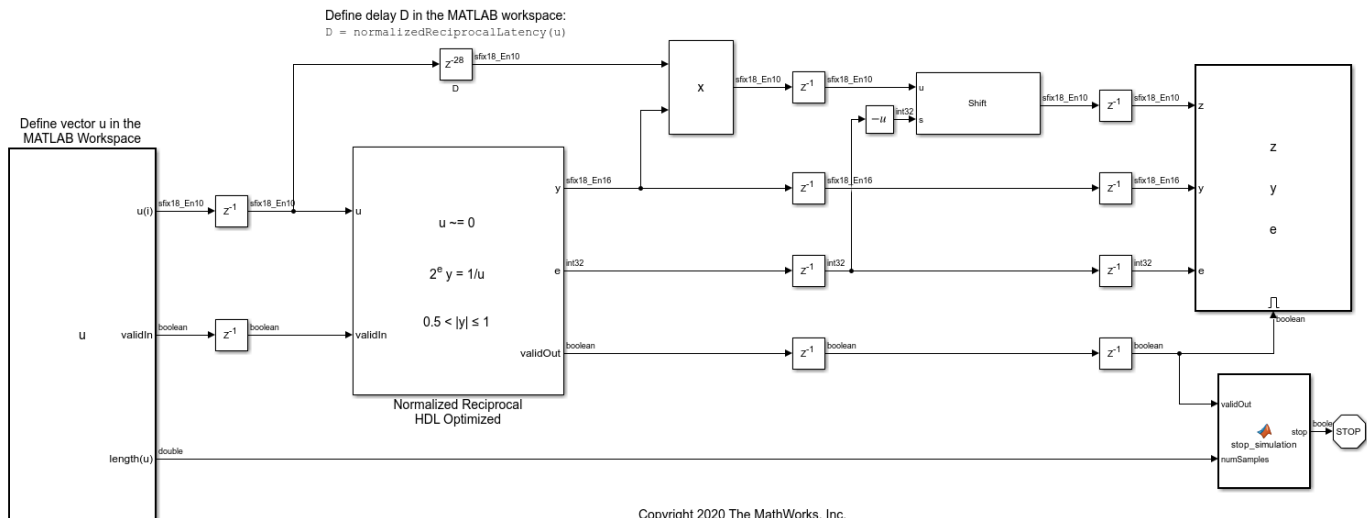
```
D =
    28
```

Run the Normalized Reciprocal Simulink Model

As a simple example where the range of a product of a reciprocal and another variable is known, compute the normalized reciprocal of a value and multiply it by itself, so the final product should be equal to one. Even though the input value and reciprocal have large ranges, the product of the value and its reciprocal have a known range.

Note that the product $u*y$ can be computed in the same type as u because $0.5 < \text{abs}(y) \leq 1$, and so the product does not grow larger in magnitude.

```
model = 'NormalizedReciprocalModel';
open_system(model)
out = sim(model);
```



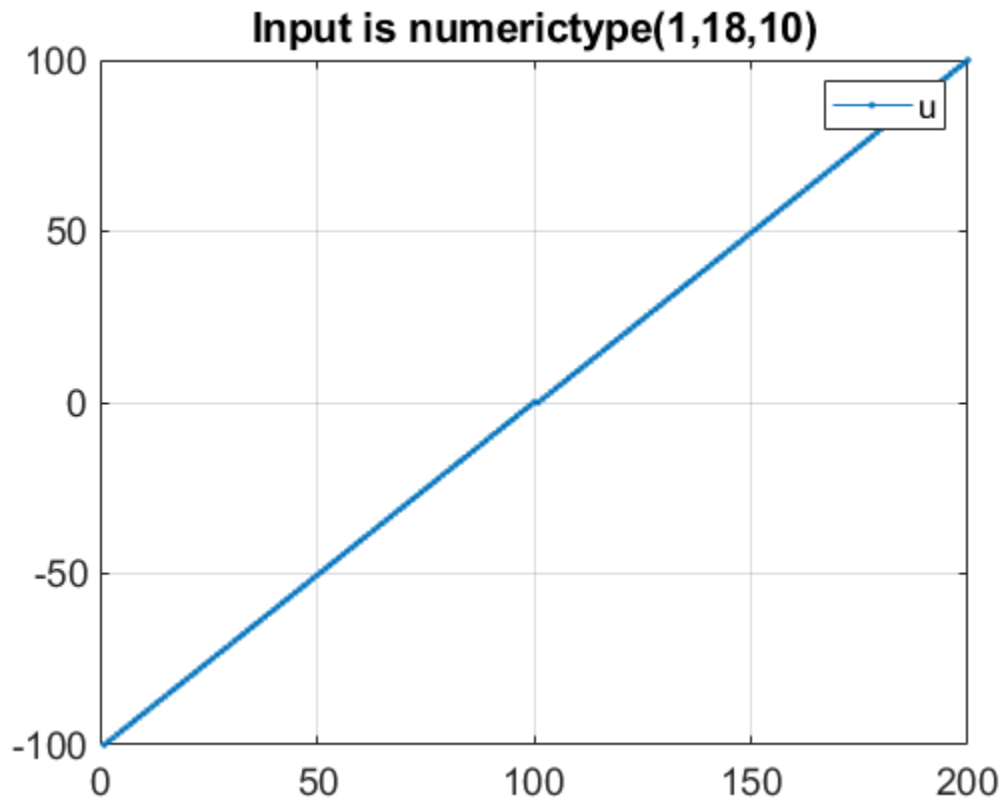
Copyright 2020 The MathWorks, Inc.

Analyze the Results of the Normalized Reciprocal Simulink Model

In the following plots, you can see that the input u and the reciprocal $1/u$ have large ranges, but the normalized reciprocal y is in the range -1 to 1.

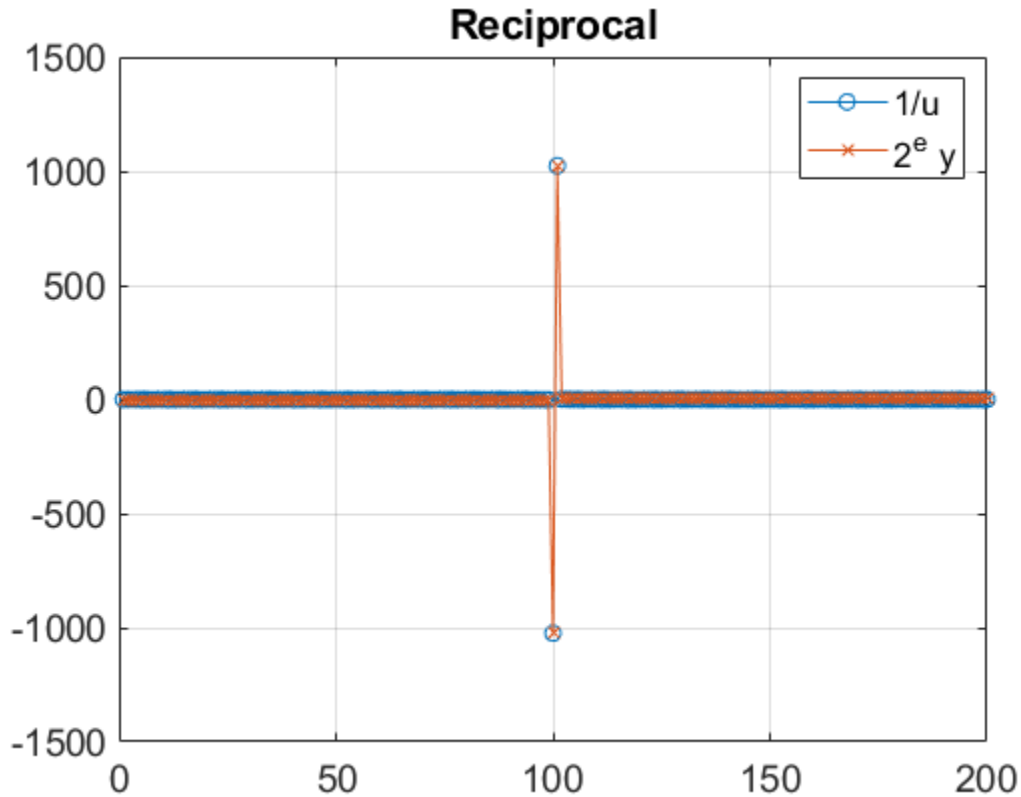
The input, u , is a vector of fi values which span the range of -100 to 100.

```
normalizedReciprocalPlot(1,u,out.y,out.e,out.z);
```



The actual value of the reciprocal, $1/u$, is nearly identical to the normalized reciprocal, $(2.^e) .* y$. The reciprocal $1/u$ has a large range, which means that to compute a straight reciprocal in fixed point would require a data type with a high dynamic range (i.e. large word length and large fraction length).

```
normalizedReciprocalPlot(2,u,out.y,out.e,out.z);
```

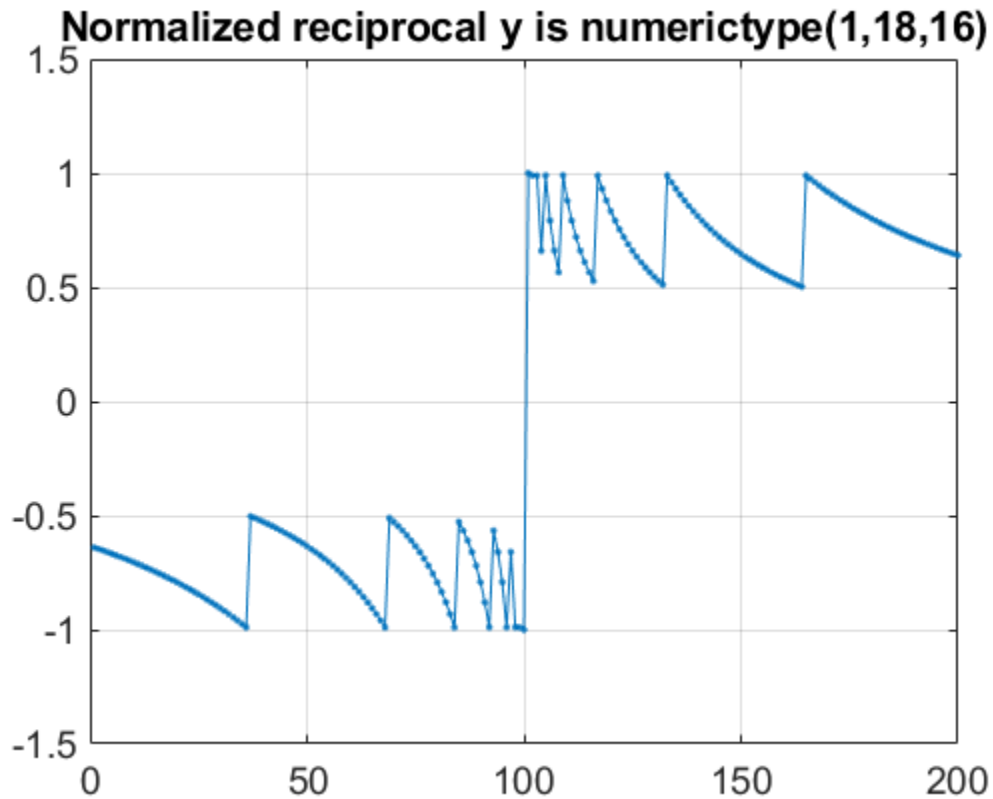


The normalized reciprocal, y , is in the range

$$0.5 < |y| \leq 1$$

and can be efficiently and accurately stored in a data type with the same word length as u . The numeric type of y has a fraction length equal to two less than the word length of u . This additional integer bit ensures that y can be positive or negative, and can also exactly represent the values -1 and +1.

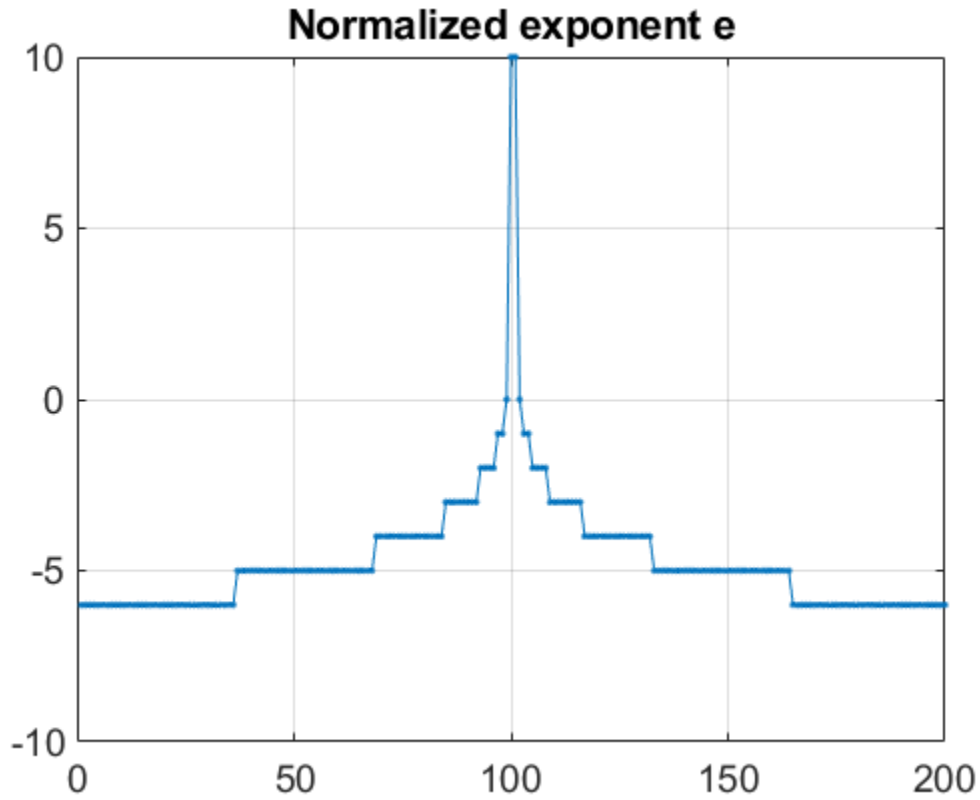
```
normalizedReciprocalPlot(3,u,out.y,out.e,out.z);
```



The normalized exponent, e , and the magnitude of the input, u , have an inverse relationship. When u is large in magnitude, then e is a large negative value. When u is close to zero, then e is a large positive value. The relationship is

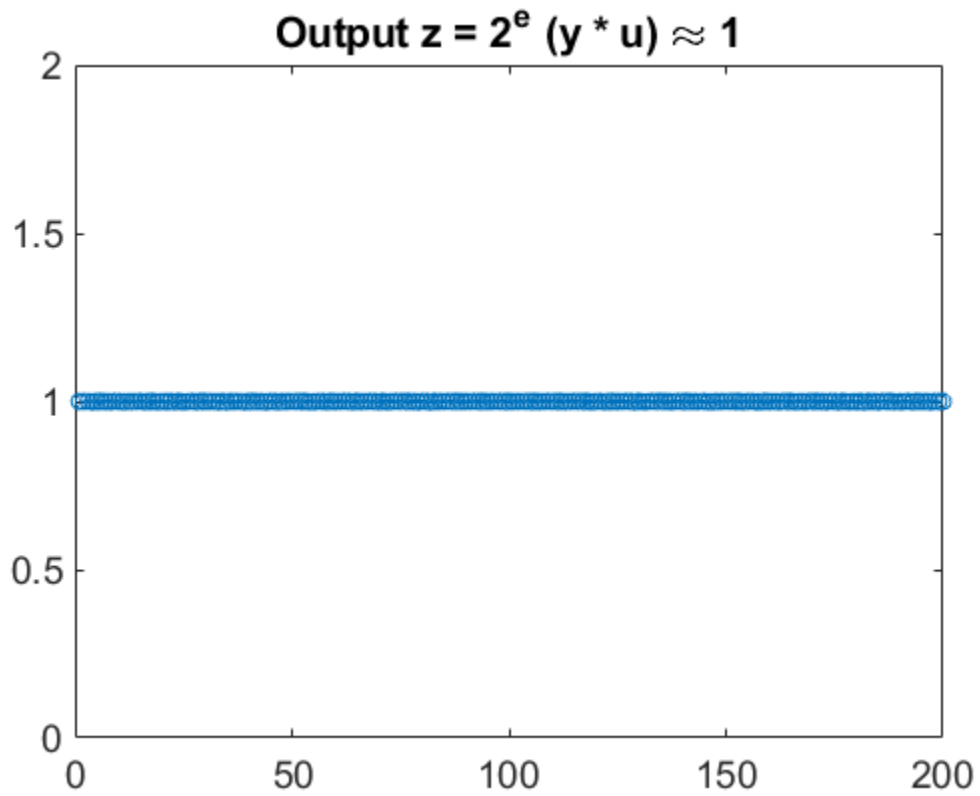
$$(2.^e) .* y = 1./u,$$

```
normalizedReciprocalPlot(4,u,out.y,out.e,out.z);
```



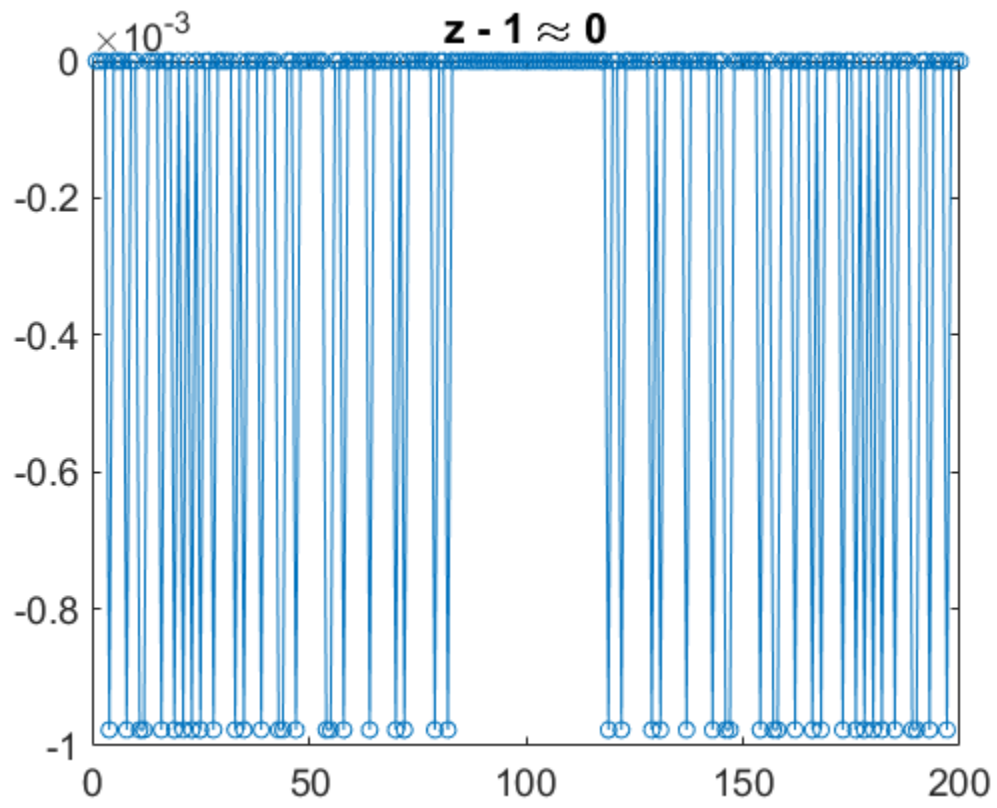
The model multiplies the normalized reciprocal y by the input u and then scales the result by shifting by e . Since y has a magnitude less than 1, the product can be done in the fixed-point type of u . The output, $z = 2^e (y * u)$, is approximately equal to 1.

```
normalizedReciprocalPlot(5,u,out.y,out.e,out.z);
```



The output, z , has a small roundoff error as a result of computing the normalized reciprocal with fixed-point data types.

```
normalizedReciprocalPlot(6,u,out.y,out.e,out.z);
```

Implement Hardware-Efficient Real Divide HDL Optimized

This example demonstrates how to perform the division of real numbers using hardware-efficient MATLAB® code embedded in Simulink® models. The model used in this example is suitable for HDL code generation for fixed-point inputs. The algorithm employs a fully pipelined architecture, which is suitable for FPGA or ASIC devices where throughput is of concern. This implementation also uses available on-chip resources judiciously, making it suitable for resource-conscious designs as well.

Division of Real Numbers

The division operation for two real numbers a and b , where $b \neq 0$, is defined as $\frac{a}{b} = c$, such that $a = b \cdot c$.

The CORDIC Algorithm

CORDIC is an acronym for COordinate Rotation DIGital Computer, and can be used to efficiently compute many trigonometric, hyperbolic, and arithmetic functions. For a detailed explanation of the CORDIC algorithm and its application in the calculation of a trigonometric function, see [Compute Sine and Cosine Using CORDIC Rotation Kernel](#).

Fully Pipelined Fixed-Point Computations

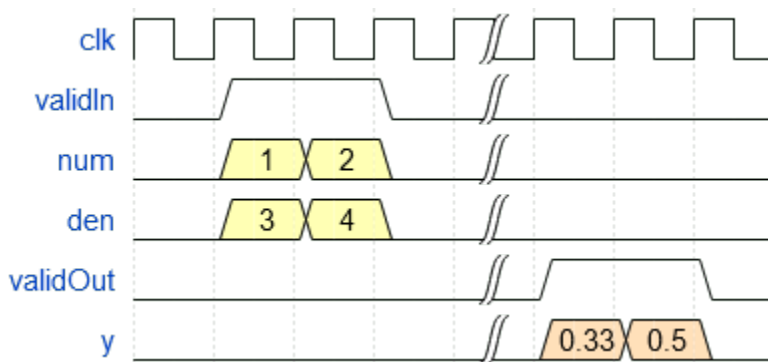
The Real Divide HDL Optimized block supports HDL code generation for fixed-point data with binary-point scaling. It is designed with this application in mind, and employs hardware specific semantics and optimizations. One of these optimizations is pipelining its entire internal circuitry to maintain a very high throughput.

When deploying intricate algorithms to FPGA or ASIC devices, there is often a trade-off between resource usage and total throughput for a given computation. Resource-sharing often reduces the resources consumed by a design, but also reduces the throughput in the process. Simple arithmetic and trigonometric computations, which typically form parts of bigger computations, require high throughput to drive circuits further in the design. Thus, fully pipelined implementations consume more on-chip resources but are beneficial in large designs.

All of the key computational units in the Real Divide HDL Optimized block are fully pipelined internally. This includes not only the CORDIC circuitry used to perform the Givens rotations, but also the adders and shifters used elsewhere in the design, thus ensuring maximum throughput.

How to Interface with the Real Divide HDL Optimized Block

Because of its fully pipelined nature, the Real Divide HDL Optimized block is able to accept input data on any cycle, including consecutive cycles. To send input data to the block, the `validIn` signal must be `true`. When the block has finished the computation and is ready to send the output, it will change `validOut` to `true` for one clock cycle. For inputs sent on consecutive cycles, `validOut` will also be set to `true` on consecutive cycles. Both the numerator and the denominator must be sent together on the same cycle.



Open the Model and Define Input Data

To open the example model, at the command line, enter:

```
mdl = 'fxpdemo_realDivide';
open_system(mdl)
```

The model contains the Real Divide HDL Optimized block connected to a data source which takes in arrays of inputs (numerators and denominators) and passes an input value from each array to the block on consecutive cycles. The output computed for each value is stored in a workspace variable. The simulation terminates when all inputs have been processed.

Define arrays of inputs `realDivideNumerators` and `realDivideDenominators`. For this example, the inputs are doubles. Note that both the numerator and the denominator should have the same datatype.

```
rng('default');
realDivideNumerators = 9*rand(1000,1) + 1;
realDivideDenominators = 9*rand(1000,1) + 1;
```

Define the output datatype to be used in the model. For this example, the outputs are also doubles. Note that fixed-point type outputs can only be used with fixed-point type inputs.

```
OutputType = 'double';
```

Simulate the Model and Examine the Output

Simulate the model.

```
sim(mdl);
```

When the simulation is complete, a new workspace variable, `realDivideOutputs`, is created to hold the computed value for each pair of inputs.

To examine the error of the calculation, compare the output of the Real Divide HDL Optimized block to that of the built-in MATLAB® divide function.

```
expectedOutput = realDivideNumerators./realDivideDenominators;
actualOutput = realDivideOutputs;
maxError = max(abs(expectedOutput - actualOutput))
```

```
maxError = 0
```

Implement Hardware-Efficient Complex Divide HDL Optimized

This example demonstrates how to perform the division of complex numbers using hardware-efficient MATLAB® code embedded in Simulink® models. The model used in this example is suitable for HDL code generation for fixed-point inputs. The algorithm employs a fully pipelined architecture, which is suitable for FPGA or ASIC devices where throughput is of concern. This implementation also uses available on-chip resources judiciously, making it suitable for resource-conscious designs as well.

Division of Complex Numbers

The division operation for two complex numbers $a + bi$ and $c + di$, where $b \neq 0$, is defined as

$$z = \frac{a + bi}{c + di}$$

After multiplying the denominator by its complex conjugate, this can be re-written as

$$z = \frac{(ac + bd) + (bc - ad)i}{c^2 + d^2}$$

The CORDIC Algorithm

CORDIC is an acronym for COordinate Rotation DIGital Computer, and can be used to efficiently compute many trigonometric, hyperbolic, and arithmetic functions. For a detailed explanation of the CORDIC algorithm and its application in the calculation of a trigonometric function, see [Compute Sine and Cosine Using CORDIC Rotation Kernel](#).

Fully Pipelined Fixed-Point Computations

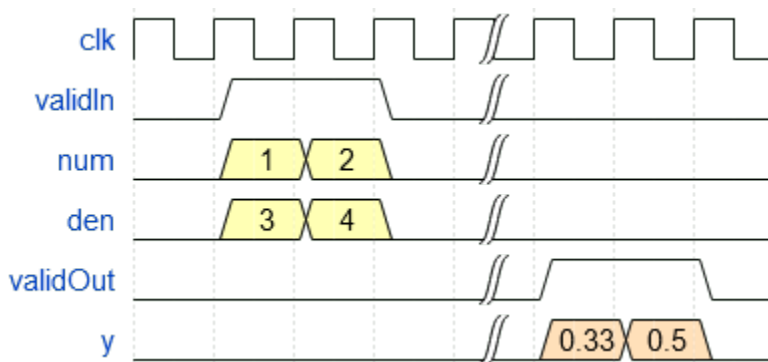
The Complex Divide HDL Optimized block supports HDL code generation for fixed-point data with binary-point scaling. It is designed with this application in mind, and employs hardware specific semantics and optimizations. One of these optimizations is pipelining its entire internal circuitry to maintain a very high throughput.

When deploying intricate algorithms to FPGA or ASIC devices, there is often a trade-off between resource usage and total throughput for a given computation. Resource-sharing often reduces the resources consumed by a design, but also reduces the throughput in the process. Simple arithmetic and trigonometric computations, which typically form parts of bigger computations, require high throughput to drive circuits further in the design. Thus, fully pipelined implementations consume more on-chip resources but are beneficial in large designs.

All of the key computational units in the Complex Divide HDL Optimized block are fully pipelined internally. This includes not only the CORDIC circuitry used to perform the Givens rotations, but also the adders and shifters used elsewhere in the design, thus ensuring maximum throughput.

Interfacing with the Complex Divide HDL Optimized Block

Because of its fully pipelined nature, the Complex Divide HDL Optimized block is able to accept input data on any cycle, including consecutive cycles. To send input data to the block, the `validIn` signal must be set to `true`. When the block has finished the computation and is ready to send the output, it will set `validOut` to `true` for one clock cycle. For inputs sent on consecutive cycles, `validOut` will also be set to `true` on consecutive cycles. Both the numerator and the denominator must be sent together on the same cycle.



Open the Model and Define Input Data

To open the example model, at the command line, enter:

```
mdl = 'fxpdemo_complexDivide';
open_system(mdl)
```

The model contains the Complex Divide HDL Optimized block connected to a data source which takes in arrays of inputs (numerators and denominators) and passes an input value from each array to the block on consecutive cycles. The output computed for each value is stored in a workspace variable. The simulation terminates when all inputs have been processed.

Define arrays of inputs `complexDivideNumerators` and `complexDivideDenominators`. For this example, the inputs are doubles. Note that both the numerator and the denominator should have the same datatype.

```
rng('default');
complexDivideNumerators = (9*rand(1000,1) + 1) + (9*rand(1000,1) + 1)*1i;
complexDivideDenominators = (9*rand(1000,1) + 1) + (9*rand(1000,1) + 1)*1i;
```

Define the output datatype to be used in the model. For this example, the outputs are also doubles. Note that fixed-point type outputs can only be used with fixed-point type inputs.

```
OutputType = 'double';
```

Simulate the Model and Examine the Output

Simulate the model.

```
sim(mdl);
```

When the simulation is complete, a new workspace variable, `complexDivideOutputs`, is created to hold the computed value for each pair of inputs.

Examine the error of the calculation by comparing the output of the Complex Divide HDL Optimized block to that of the built-in MATLAB® divide function.

```
expectedOutput = complexDivideNumerators./complexDivideDenominators;
actualOutput = complexDivideOutputs;
maxError = max(abs(expectedOutput - actualOutput))
```

```
maxError = 3.5958e-15
```

Implement Hardware-Efficient Hyperbolic Tangent

This example demonstrates how to compute the hyperbolic tangent of a given real-valued set of data using hardware-efficient MATLAB® code embedded in Simulink® models. The model used in this example is suitable for HDL code generation for fixed-point inputs. The algorithm employs an architecture that shares computational and memory units across different steps, which is beneficial when deploying to FPGA or ASIC devices with constrained resources. This implementation thus has a smaller throughput than a fully pipelined implementation, but it also has a smaller on-chip footprint, making it suitable for resource-conscious designs.

The Hyperbolic Tangent

The hyperbolic tangent function is the hyperbolic analogue of the circular tan function, and is defined as the ratio of the hyperbolic sine and hyperbolic cosine functions for a given angle α .

$$\tanh(\alpha) = \frac{\sinh(\alpha)}{\cosh(\alpha)}$$

The CORDIC Algorithm

CORDIC is an acronym for COordinate Rotation DIGital Computer, and can be used to efficiently compute many trigonometric and hyperbolic functions. For a detailed explanation of the CORDIC algorithm and its application in the calculation of a trigonometric function, see Compute Sine and Cosine Using CORDIC Rotation Kernel.

Hardware Efficient Fixed-Point Computations

The Hyperbolic Tangent HDL Optimized block supports HDL code generation for fixed-point data with binary-point scaling. It is designed with this application in mind, and employs hardware specific semantics and optimizations. One of these optimizations is resource sharing.

When deploying intricate algorithms to FPGA or ASIC devices, there is often a trade-off between resource usage and total throughput for a given computation. Fully pipelined and parallelized algorithms have the greatest throughput, but they are often too resource intensive to deploy on real devices. By implementing scheduling logic around one or several core computational circuits, it is possible to reuse resources throughout a computation. The result is an implementation with a much smaller footprint, at the cost of a reduced total throughput. This is often an acceptable trade-off, as resource shared designs can still meet overall latency requirements.

All of the key computational units in the Hyperbolic Tangent HDL Optimized block are reused throughout the computation life cycle. This includes not only the CORDIC circuitry used to perform the Givens rotations, but also the adders and multipliers used for updating the angles. This saves both DSP and fabric resources when deploying to FPGA or ASIC devices.

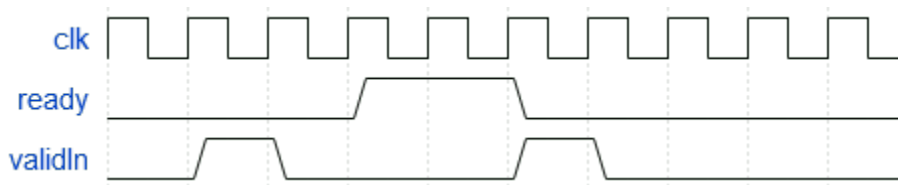
Supported Datatypes

Single, double, binary-point scaled fixed-point, and binary-point scaled-double data types are supported for simulation. However, only binary-point scaled fixed-point data types are supported for HDL code generation.

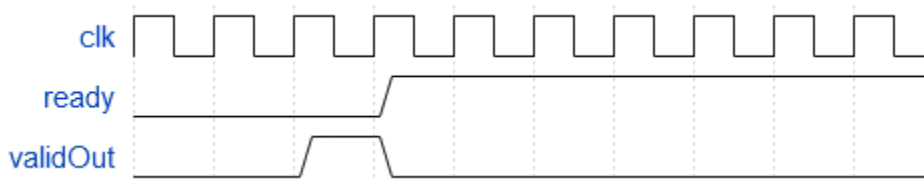
Interfacing with the Hyperbolic Tangent HDL Optimized Block

The Hyperbolic Tangent HDL Optimized block accepts data when the `ready` output is high, indicating that the block is ready to begin a new computation. To send input data to the block, the `validIn`

signal must be asserted. If the block successfully registers the input value it will de-assert the `ready` signal, and the user must then wait until the signal is asserted again to send a new input. This protocol is summarized in the following wave diagram. Note how the first valid input to the block is discarded because the block was not ready to accept input data.



When the block has finished the computation and is ready to send the output, it will assert `validOut` for one clock cycle. Then `ready` will be asserted, indicating that the block is ready to accept a new input value.



Simulate the Example Model

Open the example model by entering at the command line:

```
mdl = 'fxpdemo_tanh';
open_system(mdl)
```

The model contains the Hyperbolic Tangent HDL Optimized block connected to a data source which takes in an array of inputs and passes an input value from the array to the Hyperbolic Tangent HDL Optimized block when it is ready to accept a new input. The output computed for each value is stored in a workspace variable. The simulation terminates when all inputs have been processed.

Define an array of inputs, `tanhInput`. For this example, inputs are doubles.

```
tanhInput = -10:0.05:10;
```

Simulate the model.

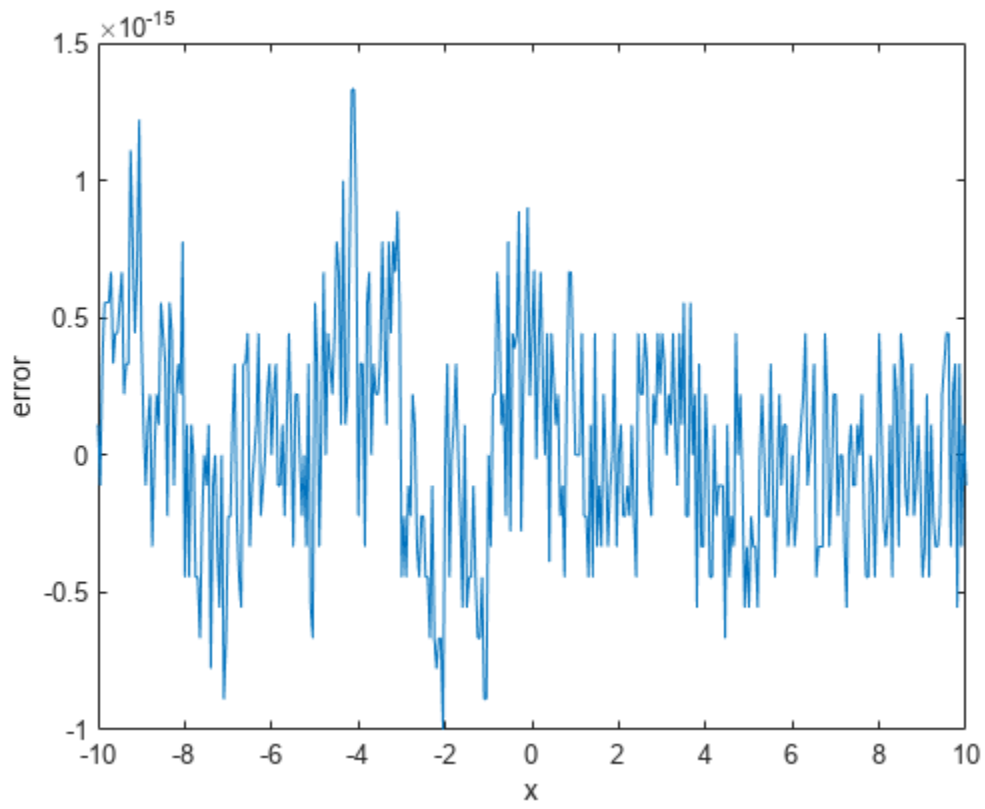
```
sim(mdl);
```

When the simulation is complete, a new workspace variable, `tanhOutput`, is created to hold the computed value for each input.

Plot the Output

Plot the error of the calculation by comparing the output of the Hyperbolic Tangent HDL Optimized block to that of the MATLAB® `tanh` function.

```
figure(1);
plot(tanhInput, tanhOutput - tanh(tanhInput));
xlabel('x');
ylabel('error');
```



Implement HDL Optimized Modulo By Constant

This example shows how to use the Modulo by Constant HDL Optimized block.

The modulo operation,

$$Y = X \bmod D = X - \lfloor \frac{X}{D} \rfloor \times D$$

is an important building block for many mathematical algorithms. However, this formula for $X \bmod D$ is computationally inefficient for fixed-point and integer inputs. Many embedded processors lack instructions for integer division. Those that do have them require many clock cycles to compute the answer. Division is also inefficient in commercially-available FPGAs, whose arithmetic circuits are designed for efficient multiplication, addition, and subtraction. Finally, for fixed-point modulo operations, it is difficult to optimize the word length of internal data types used for the calculation because the division operation is unbounded, even for small-wordlength inputs.

The denominator in the modulo problem is a compile-time constant, so the block can compute the floored division by using a multiplication followed by a cast. Rewriting the division operation as

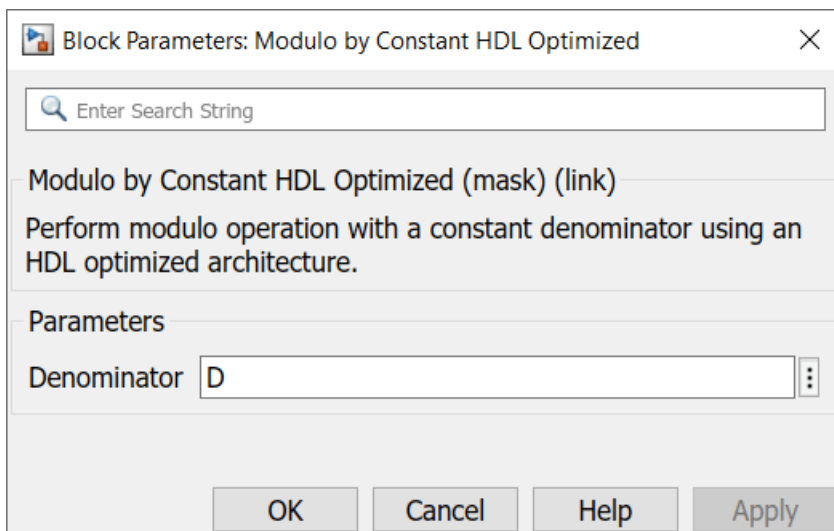
$$\frac{X}{D} = X \times \frac{1}{D}$$

shows this. The constant $1/D$ is calculated to the precision necessary to maintain both accuracy and computational efficiency. The cast that follows discards any fractional bits, which is an efficient operation on both microprocessors and FPGAs.

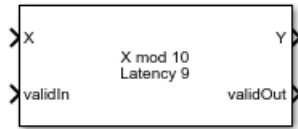
The following example shows how to use the Modulo by Constant HDL Optimized block to perform this operation and provides sample resource usage and performance statistics.

How to Use the Modulo by Constant HDL Optimized Block

The Modulo by Constant HDL Optimized block computes the modulo operation using the general strategy described above. The block requires you to specify the **Denominator** parameter as shown below.



The block is shown below when **Denominator** is set to 10. The block icon displays both the mathematical expression for the modulo operation and the latency of the block.



From the value of **Denominator** and the datatype of X , the block can compute all necessary constants and datatypes. Since it is designed for FPGA deployment, it uses the control signals `validIn` and `validOut` to indicate when X and Y are valid. Additionally, it simulates with the same latency as the generated HDL code.

To use the block, first create fixed-point input data. The format shown below is consumable by the From Workspace block.

```
>> X.time = (0:1:200).';
>> X.signals.values = fi(0:0.125:25,0,18,2).';
>> X.signals.dimensions = 1;
```

Using the same format, create a boolean `validIn` signal that toggles from false to true repeatedly.

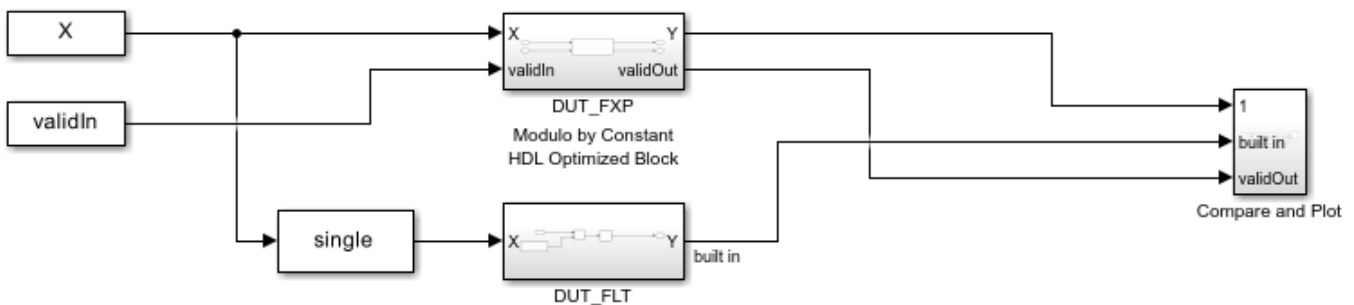
```
>> validIn.time = (0:1:200).';
>> validIn.signals.values = [false; repmat([true false]', 100, 1)];
>> validIn.signals.dimensions = 1;
```

To finish setting up the data for the problem, set D equal to the constant denominator to use for the modulo operation.

```
>> D = 10;
```

Open the model.

```
>> open_system('modulo_by_constant_block_example')
```

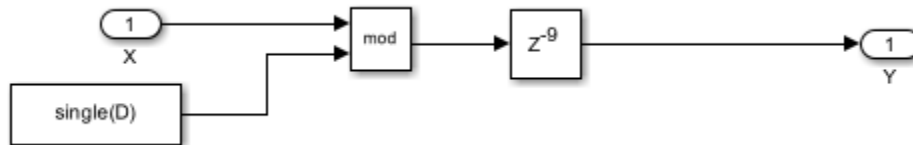


Copyright 2021 The MathWorks, Inc.

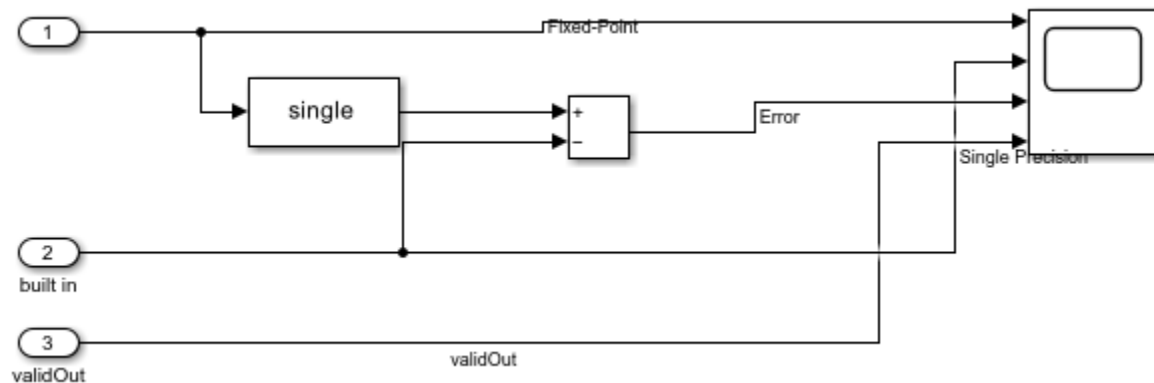
The `DUT_FXP` subsystem computes $X \bmod D$ for fixed-point inputs using the Modulo by Constant HDL Optimized block.



The DUT_FLT subsystem computes $X \bmod D$ using a Math Function block with the **Function** parameter set to `mod`. Because the Simulink `mod` operation only supports floating-point and integer inputs, the input data is cast to single precision before being input to the Math Function block. It additionally adds in a delay to match the latency of the DUT_FXP subsystem.



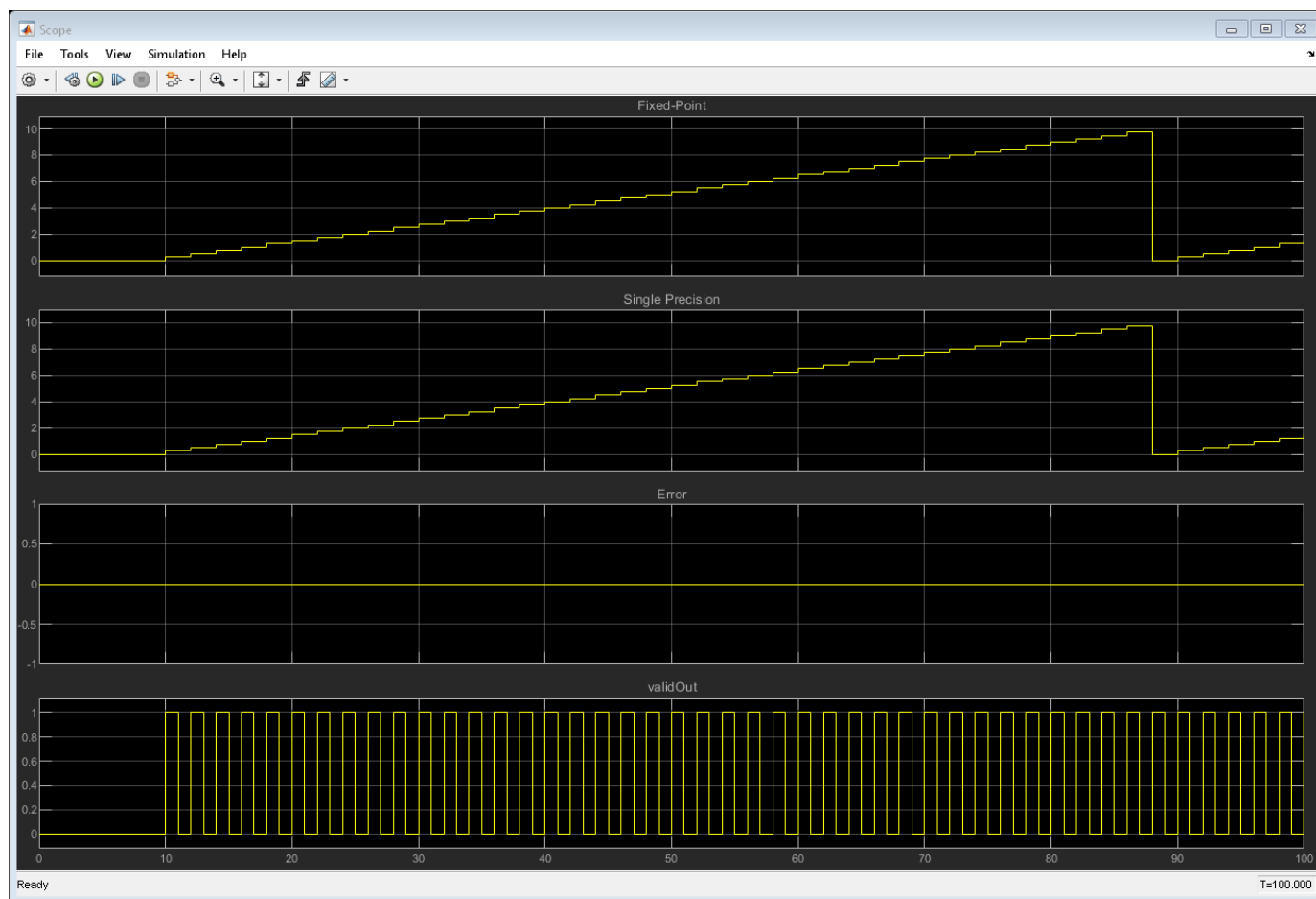
The Compare and Plot subsystem plots all outputs and computes the difference between the fixed-point computation and the floating-point ideal.



Simulate the model and examine the scope to compare the fixed-point and floating-point results.

```
>> sim('modulo_by_constant_block_example')
```

The results from the Modulo by Constant HDL Optimized and Math Function blocks agree exactly, as the plot below shows. Note that this plot displays the latency in the system, as there is a delay between the start of the simulation and the first time `validOut` goes high.



Generate HDL Code

If you have an HDL Coder license, you can generate and deploy HDL Code for the DUT_FXP subsystem as shown below.

```
>> makehdl('modulo_by_constant_block_example/DUT_FXP');
```

Implemented HDL Statistics

Sample statistics for resource usage on a Xilinx® Virtex®-7 XC7VX485 FFG1157-1 device are shown below. The implemented design is able to run at greater than 500MHz on this device.

| Resources | Usage |
|-----------------|-------|
| LUT | 33 |
| LUTRAM | 8 |
| Slice Registers | 57 |
| DSP48 | 1 |

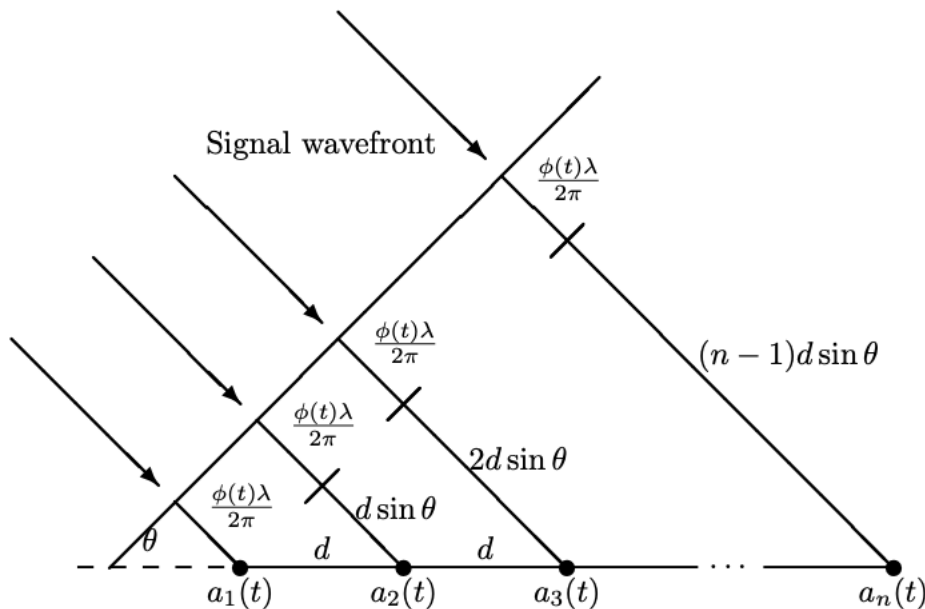
Fixed-Point HDL-Optimized Minimum-Variance Distortionless-Response (MVDR) Beamformer

This example shows how to implement a fixed-point HDL-optimized minimum-variance distortionless-response (MVDR) beamformer, also known as Capon's method [5] and robust adaptive beamforming (RAB) [6]. For more information on beamformers, see “Conventional and Adaptive Beamformers” (Phased Array System Toolbox).

MVDR Objective

The MVDR beamformer preserves the gain in the direction of arrival of a desired signal and attenuates interference from other directions [1], [2].

Given readings from a sensor array, such as the uniform linear array (ULA) in the following diagram, form data matrix A from samples of the array, where $a(t)$ is an n -by-1 column vector of readings from the array sampled at time t , and $a(t)^H$ is one row of matrix A . Many more samples are taken than there are elements in the array. This results in the number of rows in A being much greater than the number of columns. An estimate of the covariance matrix is $A^H A$, where A^H is the Hermitian or complex-conjugate transpose of A .



Compute the MVDR beamformer response by solving the following equation for x , where b is a steering vector pointing in the direction of the desired signal.

$$(A^H A)x = b$$

The MVDR weight vector w is computed from x and b using the following equation, which normalizes x to preserve the gain in the direction of arrival of the desired signal.

$$w = \frac{x}{b^H x}$$

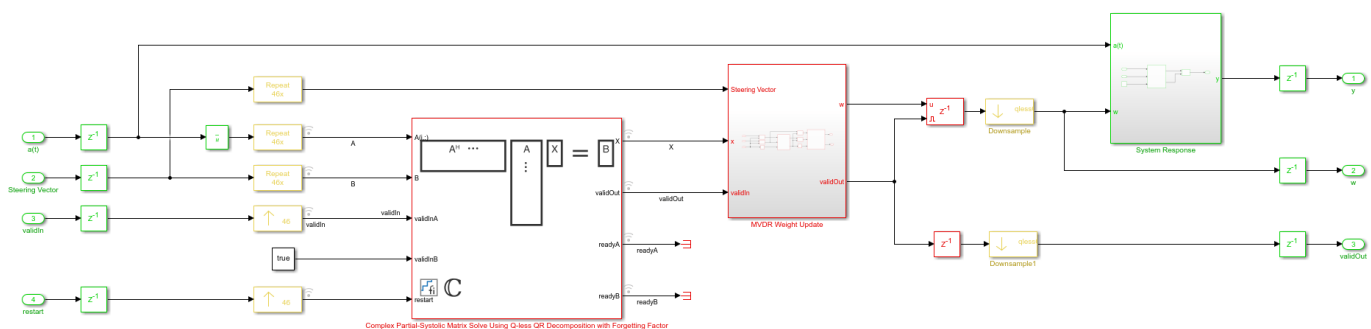
The MVDR system response is the inner product between the MVDR weight vector w and a current sample from the sensor array $a(t)$.

$$y = w^H a(t)$$

HDL-Optimized MVDR

The three equations in the previous section are implemented by the three main blocks in the following model. The rate changes give the matrix solve additional clock cycles to update before the next input sample. The number of clock cycles between a valid input and when the complex matrix solve block is ready is its input wordlength to allow time for CORDIC iterations, plus 9 cycles for internal delays.

```
load_system('MVDRBeamformerHDLOptimizedModel');
open_system('MVDRBeamformerHDLOptimizedModel/MVDR - HDL Optimized');
```



Instead of forming data matrix A and computing the Cholesky factorization of covariance matrix $A^H A$, the upper-triangular matrix of the QR decomposition of A is computed directly and updated as each data vector $a(t)$ streams in from the sensor array. Because the data is updated indefinitely, a forgetting factor is applied after each factorization. To integrate with an equivalent of a matrix of m rows, the forgetting factor α should be set to

$$\alpha = \exp(-1/(2m)).$$

This example simulates the equivalent of a matrix with $m = 300$ rows, so the forgetting factor is set to 0.9983.

The **Complex Partial-Systolic Matrix Solve Using Q-less QR Decomposition with Forgetting Factor** block is implemented using the method found in [3]. The upper-triangular matrix R from the QR decomposition of A is identical to the Cholesky factorization of $A^H A$ except the signs of values on the diagonal. Solving the matrix equation $(A^H A)x = b$ by computing the Cholesky factorization of $A^H A$ is not as efficient or as numerically sound as computing the QR decomposition of A directly [4].

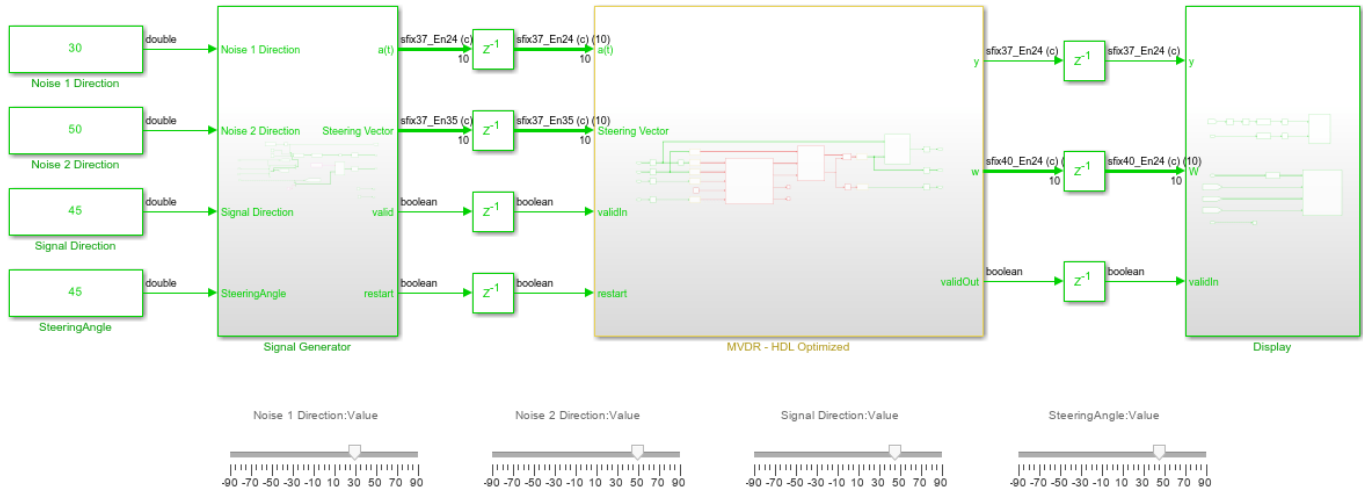
Run Model

Open and simulate the model.

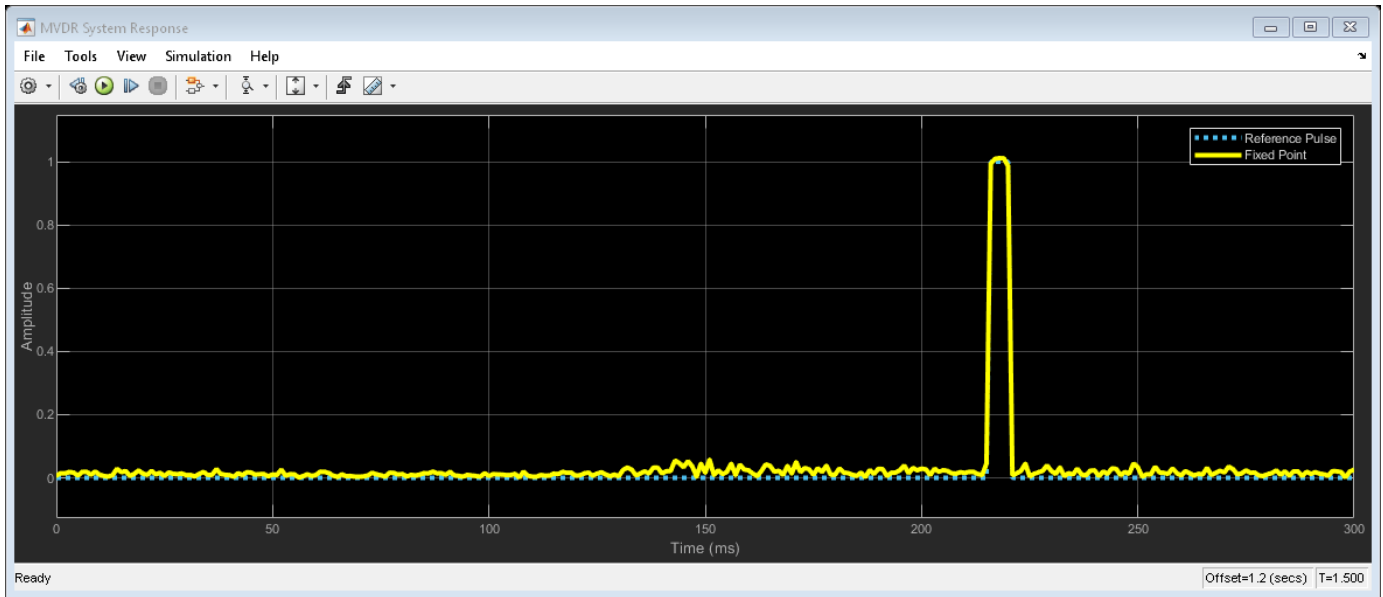
```
open_system('MVDRBeamformerHDLOptimizedModel');
sim('MVDRBeamformerHDLOptimizedModel');
```

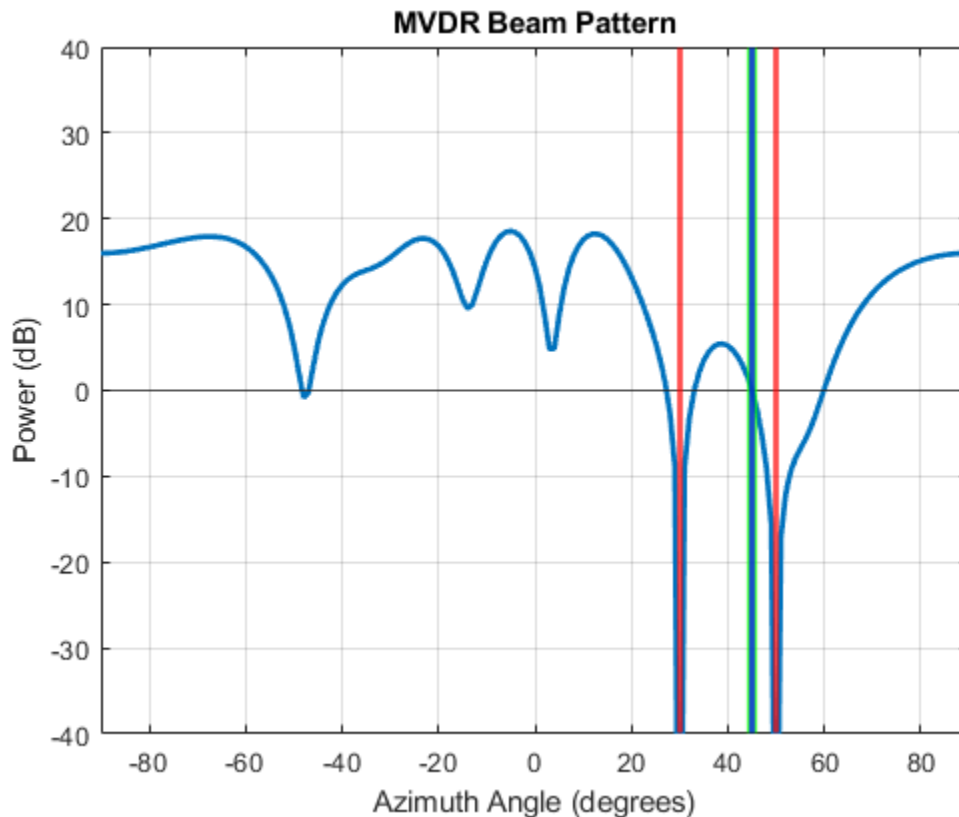
Minimum Variance Distortionless Response (MVDR) Adaptive Beamforming with Interference

Assumptions:
 -Antenna is a uniform linear array (ULA) with half wavelength between elements
 -The SNR at each antenna has a power of 50dB
 -The operating frequency of the system is 100 MHz



Copyright 2020-2021 The MathWorks, Inc.





The desired rectangular pulse appears when the noise sources are nulled. This example simulates with the same latency as the hardware, so you can see the signal settle over time as the simulation starts and when the directions are changed.

When the signal direction and steering angle are aligned as indicated by the blue and green lines, you can see that the beam pattern has a gain of 0 dB. The noise sources are nulled as indicated by the red lines.

As the model is simulating, you can adjust the signal direction, steering angle and noise directions by dragging the sliders, or by editing the constant values.

Set Parameters

The parameters for the beamformer are set in the model workspace. You can modify the parameters by editing and running the `setMVDRExampleModelWorkspace` function.

References

- [1] V. Behar et al. "Parameter Optimization of the adaptive MVDR QR-based beamformer for jamming and multipath suppression in GPS/GLONASS receivers". In: Proc. 16th Saint Petersburg International Conference on Integrated Navigation systems. Saint Petersburg, Russia, May 2009, pp. 325--334.
- [2] Jack Capon. "High-resolution frequency-wavenumber spectrum analysis". In: vol. 57. 1969, pp. 1408--1418.
- [3] C.M. Rader. "VLSI Systolic Arrays for Adaptive Nulling". In: IEEE Signal Processing Magazine (July 1996), pp. 29--49.

[4] Charles F. Van Loan. Introduction to Scientific Computing: A Matrix-Vector Approach Using Matlab. Second edition. Prentice-Hall, 2000. isbn: 0-13-949157-0.

[5] J. Capon, "High-resolution frequency-wavenumber spectrum analysis," in Proceedings of the IEEE, vol. 57, no. 8, pp. 1408-1418, Aug. 1969, doi: 10.1109/PROC.1969.7278.

[6] H. Cox, R. Zeskind and M. Owen, "Robust adaptive beamforming," in IEEE Transactions on Acoustics, Speech, and Signal Processing, vol. 35, no. 10, pp. 1365-1376, October 1987, doi: 10.1109/TASSP.1987.1165054.

Hardware-Efficient Euler Rotations Using CORDIC

This example shows how to implement Euler rotations using a CORDIC kernel. The resulting model is deployable to FPGA or ASIC devices.

Euler Rotations

Euler rotations are a common convention for describing arbitrary three dimensional rotations. You can modify them to describe either rotation of a vector in a fixed coordinate system, or rotation of a coordinate system with respect to a fixed vector. In this example, assume rotations of vectors with respect to fixed coordinate systems.

You can perform an Euler rotation of a vector as follows. First, rotate the vector by α about the z-axis. Then, rotate by β about the resultant x-axis. Finally, rotate by γ about the rotated z-axis to obtain the final coordinates of the vector.

It is common to construct Euler rotations by multiplying three individual rotation matrices. For example, you can rotate a vector \mathbf{X} using the matrix transformation below.

$$R_3(\gamma)R_1(\beta)R_3(\alpha)\mathbf{X}.$$

The matrix R_k represents a rotation about the axis k . For example, the matrix R_3 is given by

$$R_3(\alpha) = \begin{pmatrix} \cos(\alpha) & -\sin(\alpha) & 0 \\ \sin(\alpha) & \cos(\alpha) & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

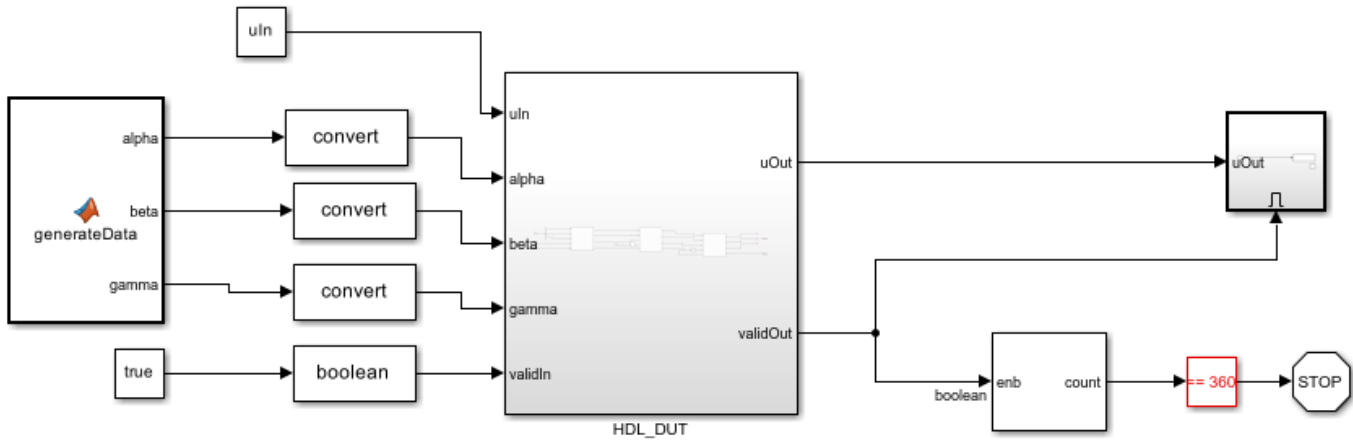
While this form is easily understood, it has several inefficiencies. If α , β , and γ are variables, you must recalculate **sin** and **cos** for each angle prior to forming up the matrix. You further need to multiply and add these results, which can result in unnecessary word length growth. Finally, the entire matrix multiplication must be properly pipelined for maximum efficiency. This can be a time consuming process.

Deploy Euler Rotations Using CORDIC Algorithm

Euler rotations operate by performing rotations in planes of intersecting coordinate axes. Therefore, efficient two dimensional rotations can be used to build up the full transformation. The Coordinate Rotation Digital Computer (CORDIC) algorithm performs these rotations using a series of shift and add operations, followed by a multiplication by a precomputed constant. It is extremely efficient and often used as a key component of high-frequency, high-throughput systems. It also eliminates the need to explicitly calculate any trigonometric functions at runtime, thus freeing further computational resources.

You can perform a full Euler rotation by using CORDIC to first rotate an input vector by α in the XY plane, followed by a CORDIC rotation by β in the resulting YZ plane. A final CORDIC rotation by γ in the resulting XY plane completes the transformation. This model shows a subsystem implementing this transformation.

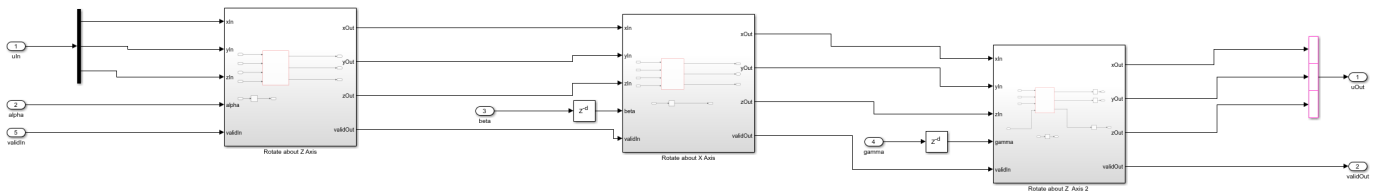
```
open_system("euler_rotations")
```



Copyright 2021 The MathWorks, Inc.

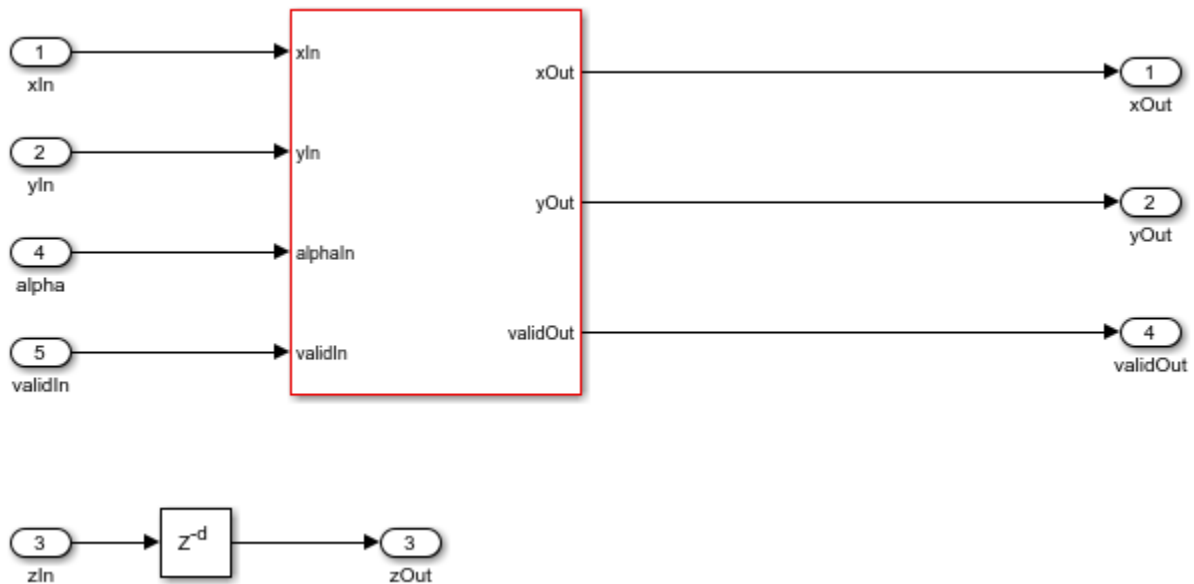
This subsystem shows the details of the actual Euler transformation. Note that the angles β and γ are delayed to align all data pipelines.

```
open_system("euler_rotations/HDL_DUT")
```



This subsystem illustrates a single planar rotation. The x and y components of the vector, x_{In} and y_{In} , are input to the CORDIC Rotation block, along with the angle α and $valid_{In}$. The z component of the vector, z_{In} , is passed along in a series of registers whose latency matches the latency of the CORDIC rotation.

```
open_system("euler_rotations/HDL_DUT/Rotate about Z Axis")
```



Define Data and Simulate

Define data and simulate the model.

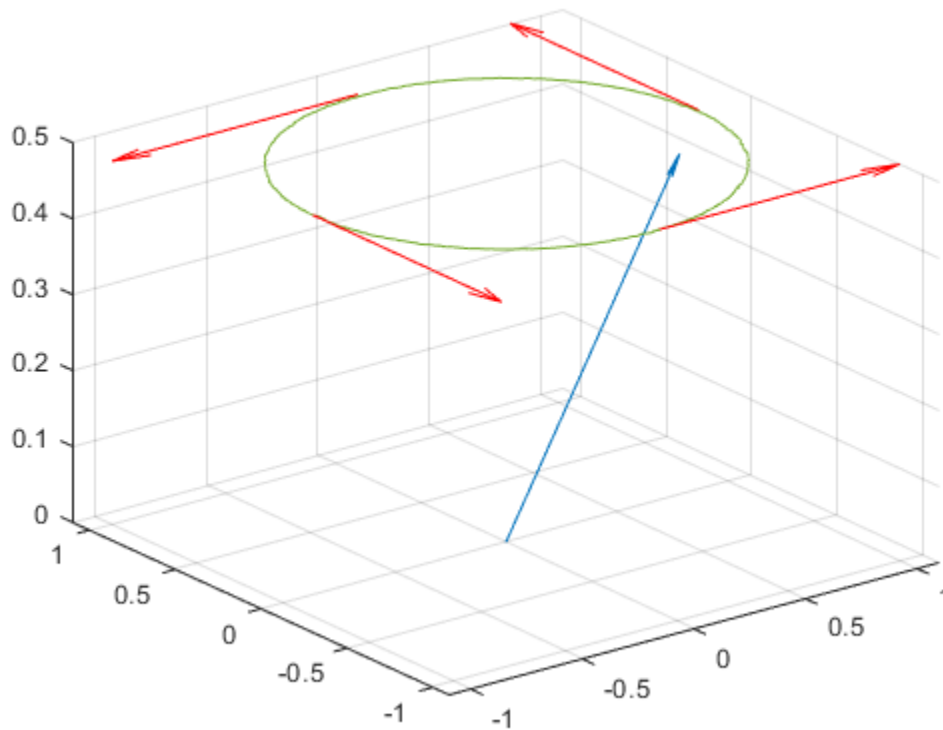
```
uIn = fi([sqrt(3)/2;0;1/2],1,16,8);
dt = numerictype(1,16,8);
nIters = 18;
out = sim("euler_rotations");
```

The model outputs the first 360 datapoints to the MATLAB® workspace. The model has been chosen to simulate one rotation about the z-axis.

Trajectory of Rotation

This figure shows the trajectory of the vector over the course of the simulation. The vector from the origin shows the initial vector, while the vectors on the raised line demonstrate the direction of rotation.

```
outData = out.simout.Data;
quiver3(0,0,0,sqrt(3)/2, 0, 0.5);
hold on;
plot3(outData(:,1), outData(:,2), 0.5*ones(length(outData)));
quiver3(sqrt(3)/2*[1;0;-1;0],sqrt(3)/2*[0;1;0;-1],[0.5;0.5;0.5;0.5],...
        0.5*[0;-1;0;1],0.5*[1;0;-1;0],[0;0;0;0], 'r-');
```



HDL Statistics

Generating HDL code for the system with the datatypes chosen gives excellent performance. The resource usage is shown below, and the device operates at approximately 373 MHz. All characterization was performed using Xilinx Vivado® using a ZC706 Evaluation Board.

```
Resource = ["LUT"; "LUTRAM"; "FF"];
Used = [3957; 99; 3673];
table(Resource, Used)
```

ans =

3x2 table

| Resource | Used |
|----------|------|
| "LUT" | 3957 |
| "LUTRAM" | 99 |
| "FF" | 3673 |

Modifying and Extending Euler Rotations

You can rearrange the constituent pieces used to develop this transformation to yield countless other transformations. The only requirement is that the full transformation you build be composed of

rotations along planes where principle axes intersect. This is one way to develop efficient, FPGA-ready solutions.

Simulation Data Inspector

- “View Data in the Simulation Data Inspector” on page 55-2
- “Import Data from a CSV File into the Simulation Data Inspector” on page 55-11
- “Microsoft Excel Import, Export, and Logging Format” on page 55-15
- “Configure the Simulation Data Inspector” on page 55-23
- “How the Simulation Data Inspector Compares Data” on page 55-31
- “Save and Share Simulation Data Inspector Data and Views” on page 55-36
- “Inspect and Compare Data Programmatically” on page 55-42
- “Limit the Size of Logged Data” on page 55-48

View Data in the Simulation Data Inspector

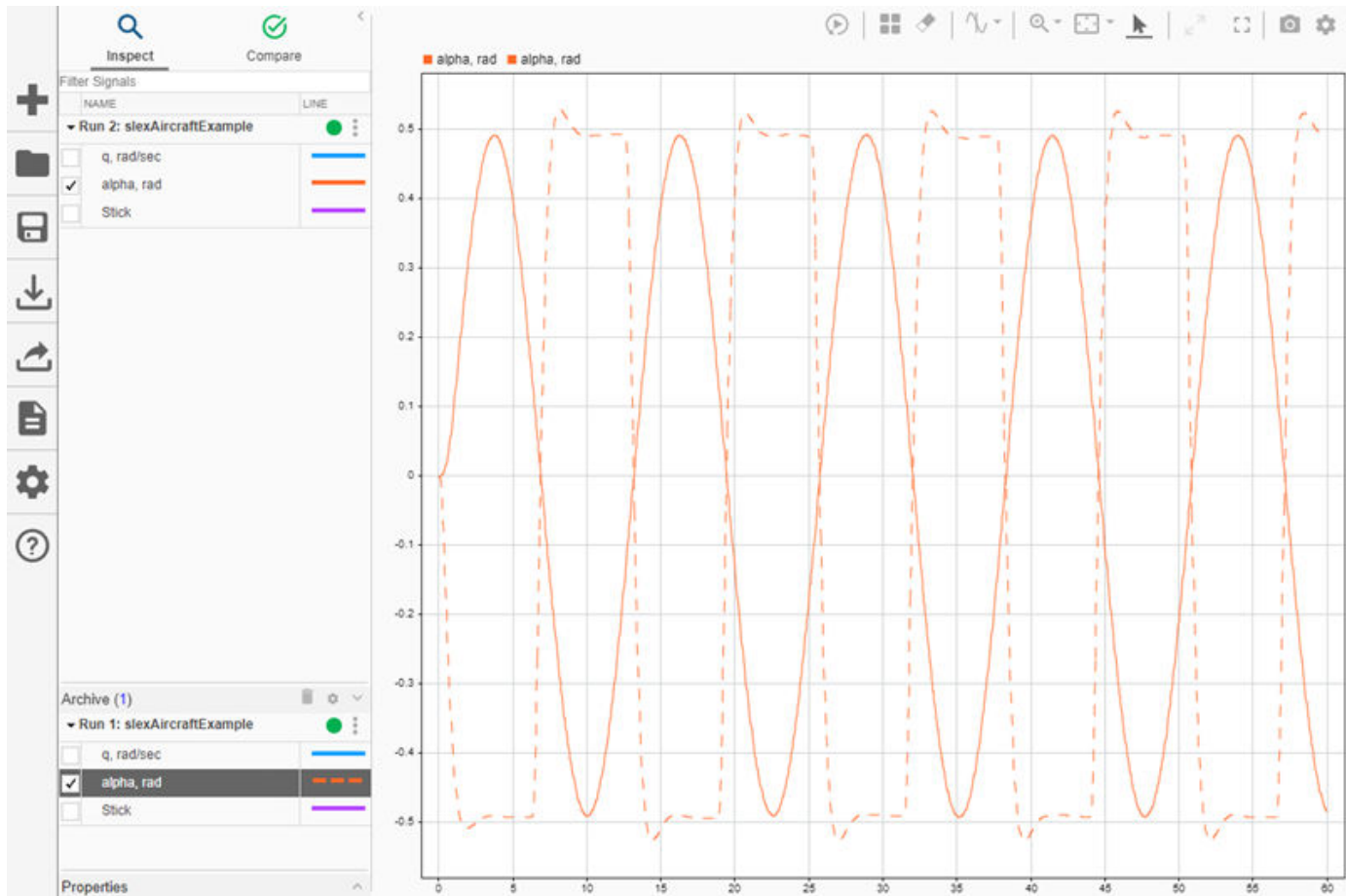
You can use the Simulation Data Inspector to visualize the data you generate throughout the design process. Simulation data that you log in a Simulink model logs to the Simulation Data Inspector. You can also import test data and other recorded data into the Simulation Data Inspector to inspect and analyze it alongside the logged simulation data. The Simulation Data Inspector offers several types of plots, which allow you to easily create complex visualizations of your data.

View Logged Data

Logged signals as well as outputs and states logged using the `Dataset` format automatically log to the Simulation Data Inspector when you simulate a model. You can also record other kinds of simulation data so the data appears in the Simulation Data Inspector at the end of the simulation. To see states and output data logged using a format other than `Dataset` in the Simulation Data Inspector, open the Configuration Parameters dialog box and, in the **Data Import/Export** pane, select the **Record logged workspace data in Simulation Data Inspector** parameter.

Note When you log states and outputs using the `Structure` or `Array` format, you must also log time for the data to record to the Simulation Data Inspector.

The Simulation Data Inspector displays available data in the table in the **Inspect** pane. To plot a signal, select the check box next to the signal. You can modify the layout and add different visualizations to analyze the simulation data. For more information, see “Create Plots Using the Simulation Data Inspector”.



The Simulation Data Inspector manages incoming simulation data using the archive. By default, the previous run moves to the archive when you start a new simulation. You can plot signals from the archive, or you can drag runs of interest back into the work area.

Import Data from the Workspace or a File

You can import data from the base workspace or from a file to view on its own or alongside simulation data. The Simulation Data Inspector supports all built-in data types and many data formats for importing data from the workspace. In general, whatever the format, sample values must be paired with sample times. The Simulation Data Inspector allows up to 8000 channels per signal in a run created from imported workspace data.

You can also import data from these types of files:

- MAT file
- CSV file — Format data as shown in “Import Data from a CSV File into the Simulation Data Inspector”.
- Microsoft® Excel® file — Format data as described in “Microsoft Excel Import, Export, and Logging Format”.

- MDF file — MDF file import is supported for Linux and Windows operating systems. The MDF file must have a `.mdf`, `.mf4`, `.mf3`, `.data`, or `.dat` file extension and contain data with only integer and floating data types.
- ULG file — Flight log data import requires a UAV Toolbox license.

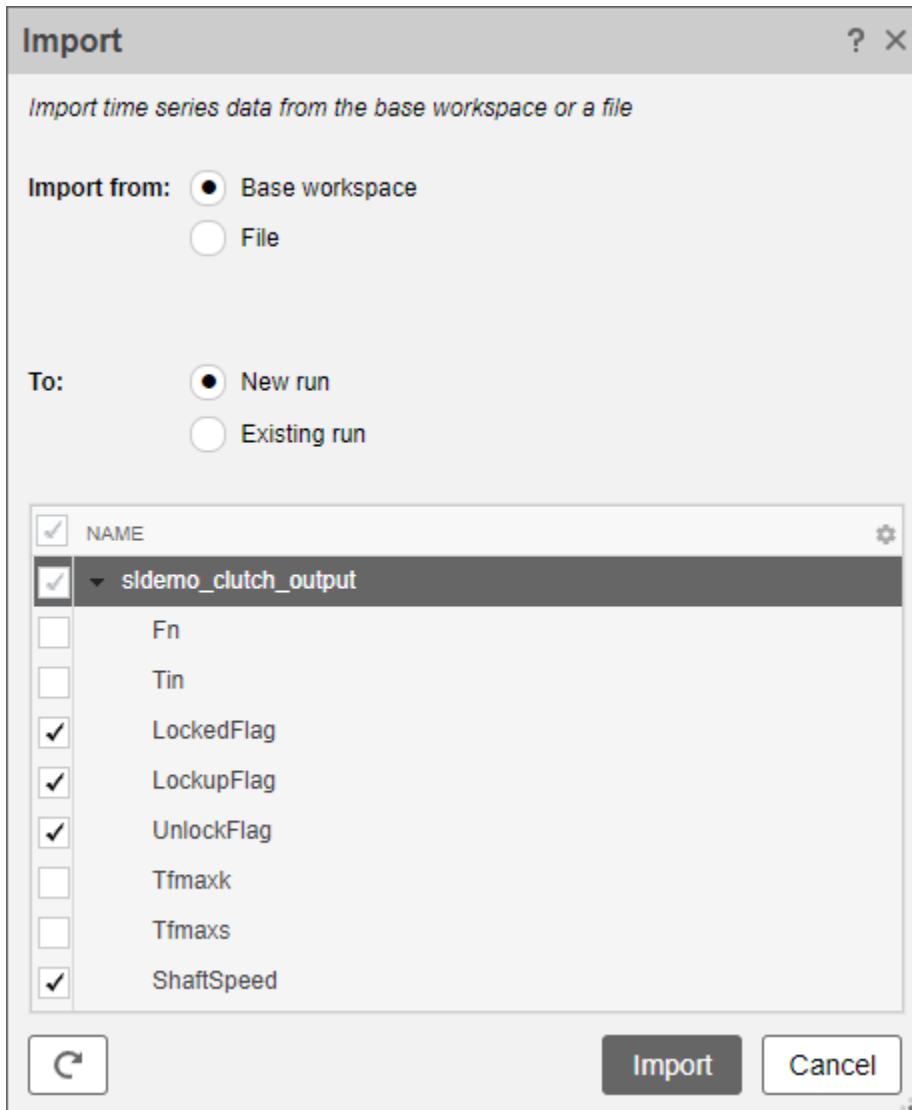
To import data from the workspace or from a file that is saved in a data or file format that the Simulation Data Inspector does not support, you can write your own workspace data or file reader to import the data using the `io.reader` class. You can also write a custom reader to use instead of the built-in reader for supported file types. For examples, see:

- “Import Data Using a Custom File Reader”
- “Import Workspace Variables Using a Custom Data Reader”



To import data, select the **Import** button in the Simulation Data Inspector.

In the Import dialog, you can choose to import data from the workspace or from a file. The table below the options shows data available for import. If you do not see your workspace variable or file contents in the table, that means the Simulation Data Inspector does not have a built-in or registered reader that supports that data. You can select which data to import using the check boxes, and you can choose whether to import that data into an existing run or a new run. To select all or none of the data, use the check box next to **NAME**.



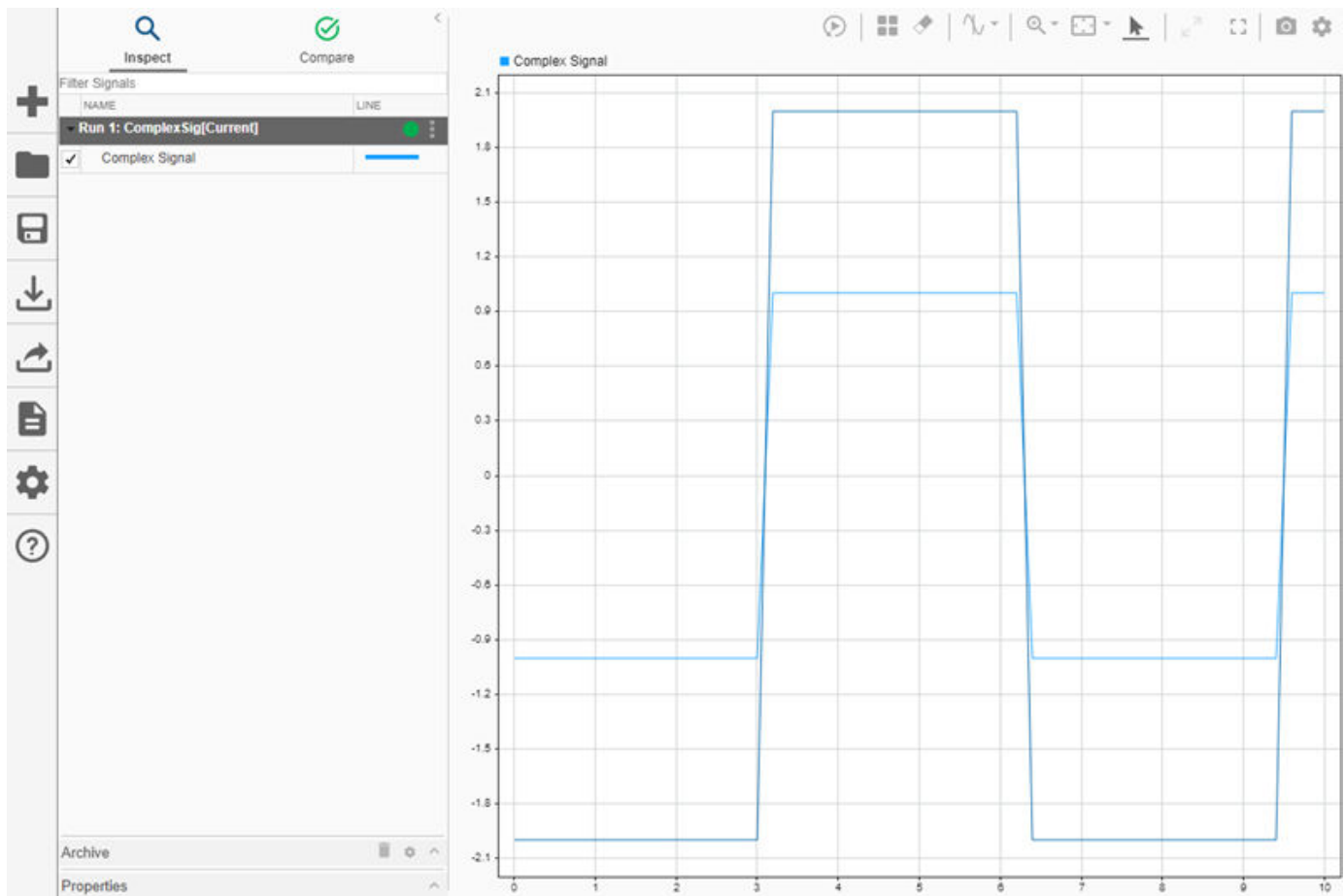
When you import data into a new run, the run always appears in the work area. You can manually move imported runs to the archive.

View Complex Data

To view complex data in the Simulation Data Inspector, import the data or log the signals to the Simulation Data Inspector. You can control how to visualize the complex signal using the **Properties** pane in the Simulation Data Inspector and in the **Instrumentation Properties** for the signal in the model. To access the **Instrumentation Properties** for a signal, right-click the logging badge for the signal and select **Properties**.

You can specify the **Complex Format** as Magnitude, Magnitude-Phase, Phase, or Real-Imaginary. If you select Magnitude-Phase or Real-Imaginary for the **Complex Format**, the Simulation Data Inspector plots both components of the signal when you select the check box for the signal. For signals in Real-Imaginary format, the **Line Color** specifies the color of the real component of the signal, and the imaginary component is a different shade of the **Line Color**. For example, the

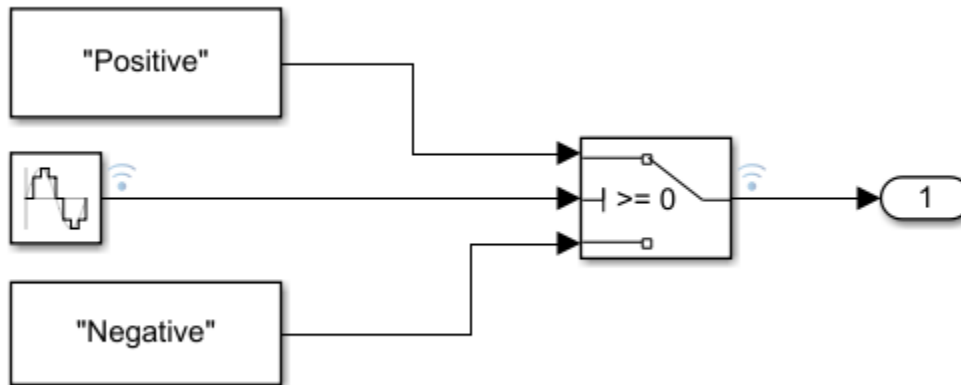
Complex Signal displays the real component of the signal in light blue, matching the **Line Color** parameter, and the imaginary component is shown in a darker shade of blue.



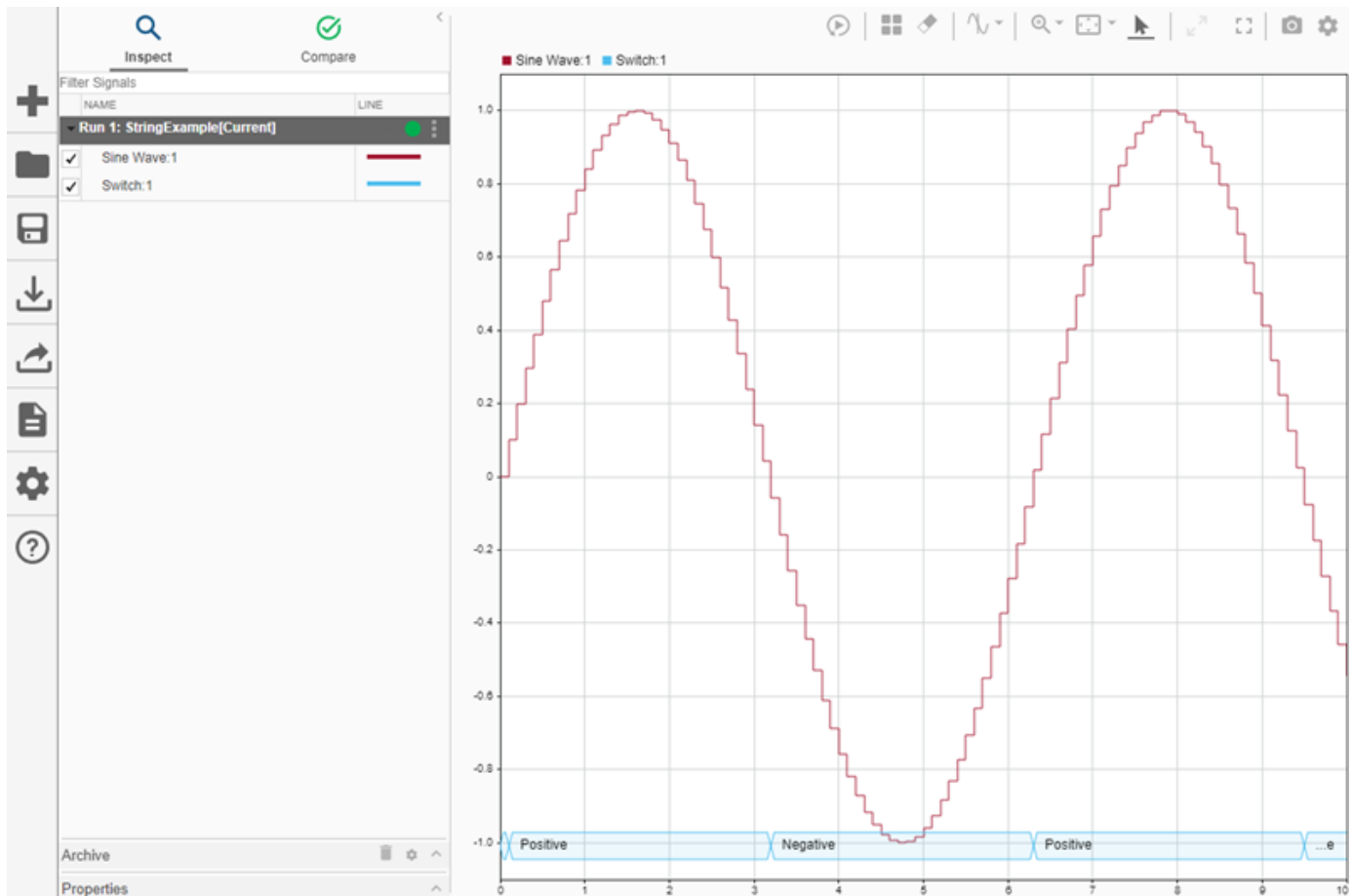
For signals in Magnitude-Phase format, the **Line Color** specifies the color of the magnitude component, and the phase is displayed in a different shade of the **Line Color**.

View String Data

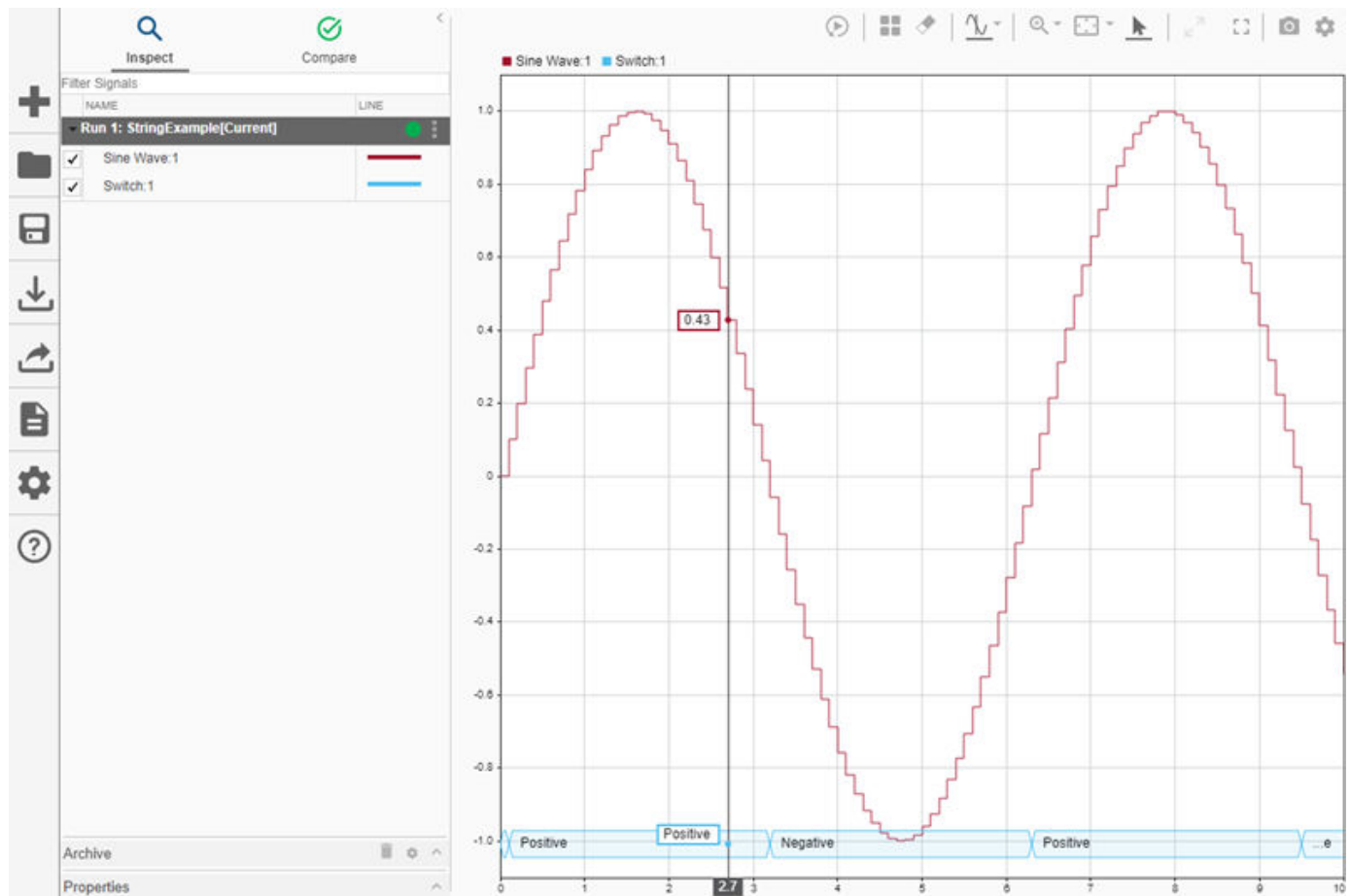
You can log and view string data with your signal data in the Simulation Data Inspector. For example, consider this simple model. The value of the sine wave block controls whether the switch sends a string reading Positive or Negative to the output.



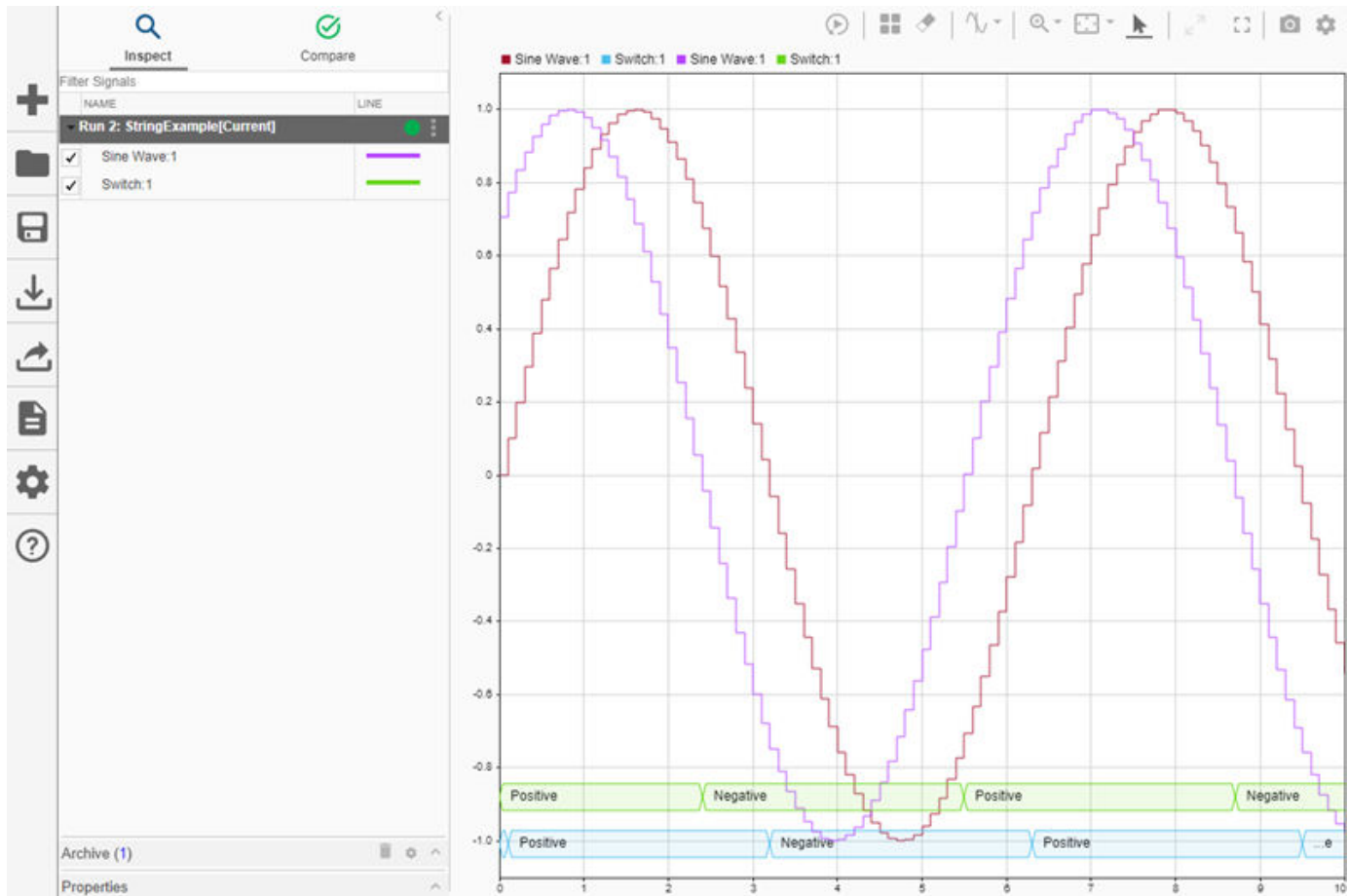
The plot shows the results of simulating the model. The string signal is shown at the bottom of the graphical viewing area. The value of the signal is displayed inside a band, and transitions in the string signal's value are marked with criss-crossed lines.



You can use cursors to inspect how the string signal values correspond with the sine signal's values.



When you plot multiple string signals on a plot, the signals stack in the order they were simulated or imported, with the most recent signal positioned at the top. For example, you might consider the effect of changing the phase of the sine wave controlling the switch.

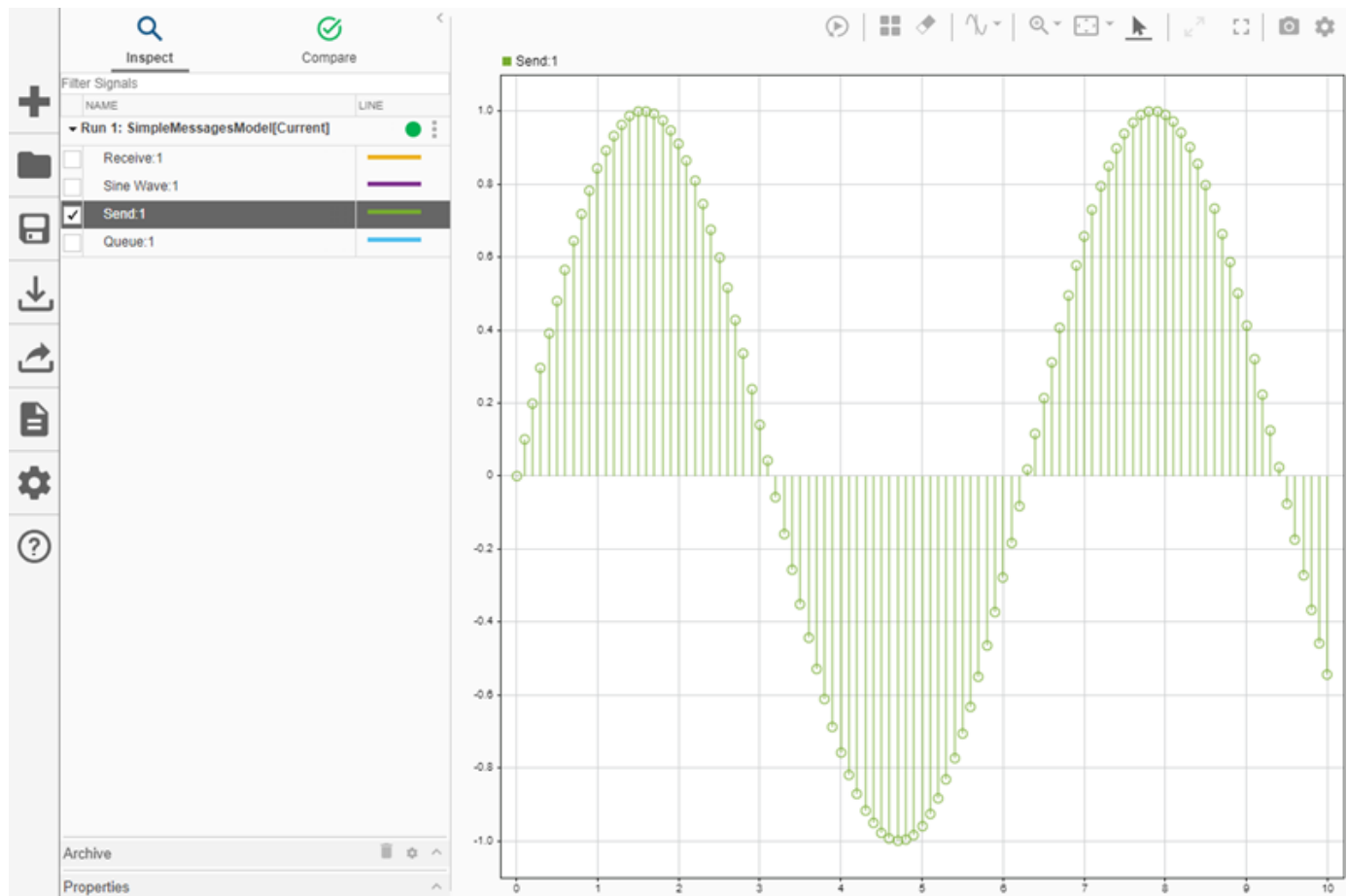


View Frame-Based Data

Processing data in frames rather than point by point provides a performance boost needed in some applications. To view frame-based data in the Simulation Data Inspector, you have to specify that the signal is frame-based in the **Instrumentation Properties** for the signal. To access the **Instrumentation Properties** dialog for a signal, right-click the signal's logging badge and select **Properties**. To specify a signal as frame-based, select **Columns as channels (frame based)** for **Input processing**.

View Event-Based Data

You can log or import event data to the Simulation Data Inspector. To view the logged event-based data, select the check box next to **Send: 1**. The Simulation Data Inspector displays the data as a stem plot, with each stem representing the number of events that occurred for a given sample time.



See Also

More About

- Inspect Simulation Data
- Compare Simulation Data
- Share Simulation Data Inspector Data and Views on page 55-36
- Decide How to Visualize Data
- Dataset Conversion for Logged Data

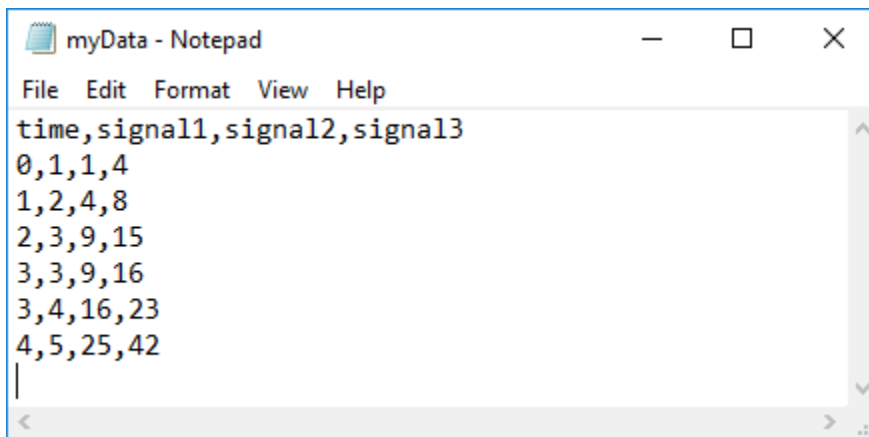
Import Data from a CSV File into the Simulation Data Inspector

To import data into the Simulation Data Inspector from a CSV file, format the data in the CSV file. Then, you can import the data using the Simulation Data Inspector UI or the `Simulink.sdi.createRun` function.

Tip When you want to import data from a CSV file where the data is formatted differently from the specification in this topic, you can write your own file reader for the Simulation Data Inspector using the `io.reader` class.

Basic File Format

In the simplest format, the first row in the CSV file is a header that lists the names of the signals in the file. The first column is time. The name for the time column must be `time`, and the time values must increase monotonically. The rows below the signal names list the signal values that correspond to each time step.



```
myData - Notepad
File Edit Format View Help
time,signal1,signal2,signal3
0,1,1,4
1,2,4,8
2,3,9,15
3,3,9,16
3,4,16,23
4,5,25,42
```

The import operation does not support time data that includes `Inf` or `NaN` values or signal data that includes `Inf` values. Empty or `NaN` signal values render as missing data. All built-in data types are supported.

Multiple Time Vectors

When your data includes signals with different time vectors, the file can include more than one time vector. Every time column must be named `time`. Time columns specify the sample times for signals to the right, up to the next time vector. For example, the first time column defines the time for `signal1` and `signal2`, and the second time column defines the time steps for `signal3`.

```

myData - Notepad
File Edit Format View Help
time,signal1,signal2,time,signal3
0,1,1,0,4
1,2,4,2,8
2,3,9,3,15
3,3,9,5,16
3,4,16
4,5,25

```

Signal columns must have the same number of data points as the associated time vector.

Signal Metadata

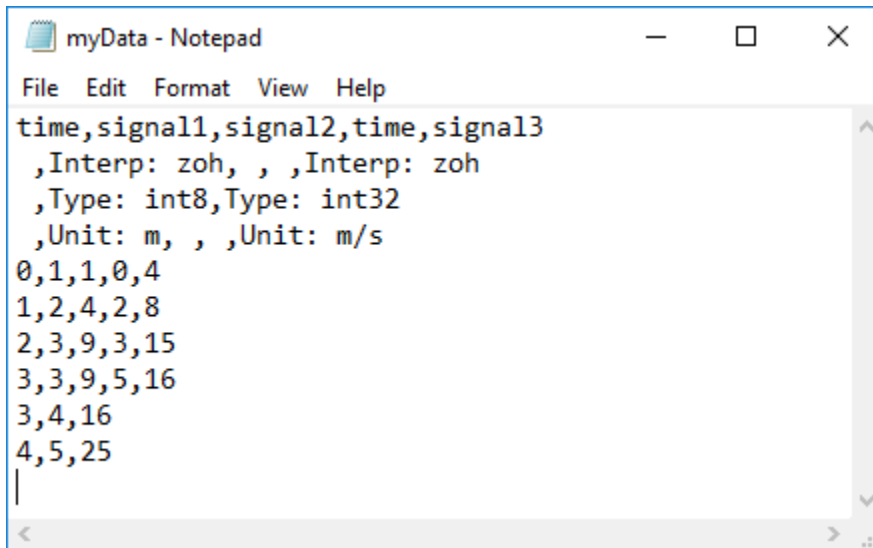
You can specify signal metadata in the CSV file to indicate the signal data type, units, interpolation method, block path, and port index. List metadata for each signal in rows between the signal name and the signal data. Label metadata according to this table.

| Signal Property | Label | Value |
|----------------------|------------|---|
| Data type | Type: | Built-in data type. |
| Units | Unit: | Supported unit. For example, Unit: m/s specifies units of meters per second. For a list of supported units, enter <code>showunitslist</code> in the MATLAB Command Window. |
| Interpolation method | Interp: | linear, zoh for zero order hold, or none. |
| Block Path | BlockPath: | Path to the block that generated the signal. |
| Port Index | PortIndex: | Integer. |

You can also import a signal with a data type defined by an enumeration class. Instead of using the Type: label, use the Enum: label and specify the value as the name of the enumeration class. The definition for the enumeration class must be saved on the MATLAB path.

When an imported file does not specify signal metadata, the Simulation Data Inspector assumes double data type and linear interpolation. You can specify the interpolation method as linear, zoh (zero-order hold), or none. If you do not specify units for the signals in your file, you can assign units to the signals in the Simulation Data Inspector after you import the file.

You can specify any combination of metadata for each signal. Leave a blank cell for signals with less specified metadata.



```

myData - Notepad
File Edit Format View Help
time,signal1,signal2,time,signal3
,Interp: zoh, , ,Interp: zoh
,Type: int8,Type: int32
,Unit: m, , ,Unit: m/s
0,1,1,0,4
1,2,4,2,8
2,3,9,3,15
3,3,9,5,16
3,4,16
4,5,25
|

```

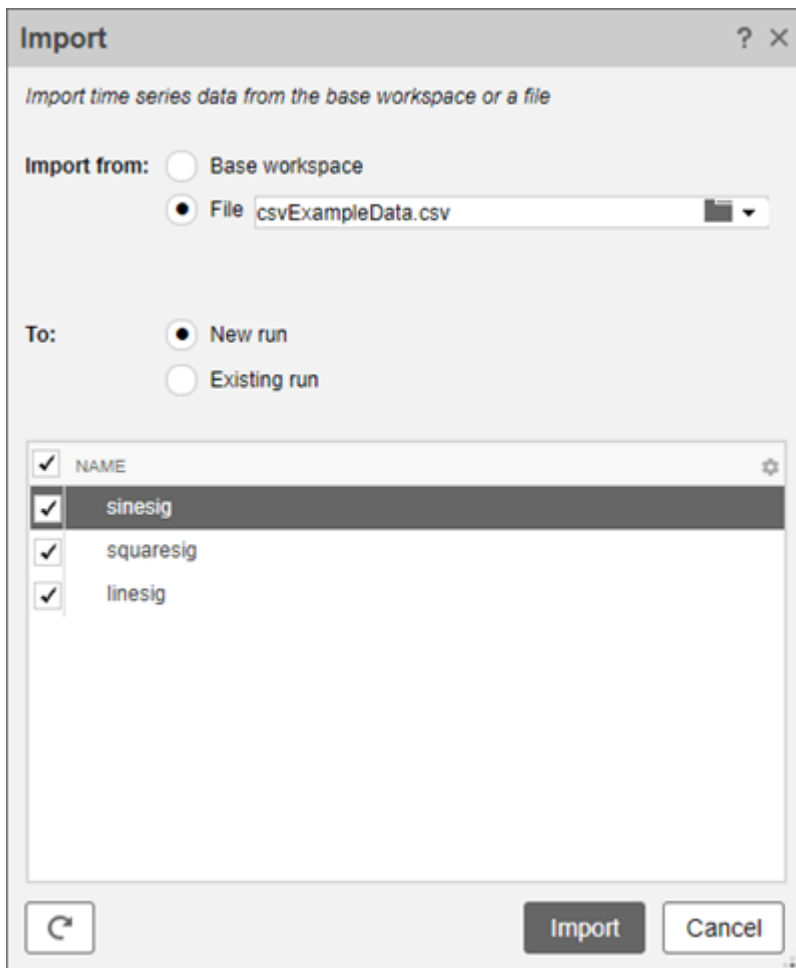
Import Data from a CSV File

You can import data from a CSV file using the Simulation Data Inspector UI or using the `Simulink.sdi.createRun` function.

To import data using the UI, open the Simulation Data Inspector using the `Simulink.sdi.view`

function or the **Data Inspector** button in the Simulink™ toolstrip. Then, click **Import** .

In the Import dialog, select the option to import data from a file and navigate in the file system to select the file. After you select the file, data available for import shows in the table. You can choose which signals to import and whether to import them to a new or existing run. This example imports all available signals to a new run. To select all or none of the signals, select or clear the check box next to NAME. After selecting the options, click the **Import** button.



When you import data into a new run using the UI, the new run name includes the run number followed by `Imported_Data`.

When you import data programmatically, you can specify the name of the imported run.

```
csvRunID = Simulink.sdi.createRun('CSV File Run', 'file', 'csvExampleData.csv');
```

See Also

Functions

`Simulink.sdi.createRun`

More About

- “View Data in the Simulation Data Inspector”
- “Microsoft Excel Import, Export, and Logging Format”
- “Import Data Using a Custom File Reader”

Microsoft Excel Import, Export, and Logging Format

Using the Simulation Data Inspector or Simulink Test, you can import data from a Microsoft Excel file or export data to a Microsoft Excel file. You can also log data to an Excel file using the Record block. The Simulation Data Inspector, Simulink Test, and the Record block all use the same file format, so you can use the same Microsoft Excel file with multiple applications.

Tip When the format of the data in your Excel file does not match the specification in this topic, you can write your own file reader to import the data using the `io.reader` class.

Basic File Format

In the simplest format, the first row in the Excel file is a header that lists the names of the signals in the file. The first column is time. The name for the time column must be `time`, and the time values must increase monotonically. The rows below the signal names list the signal values that correspond to each time step.

| | A | B | C | D |
|---|-------------------|----------------------|----------------------|----------------------|
| 1 | <code>time</code> | <code>signal1</code> | <code>signal2</code> | <code>signal3</code> |
| 2 | 0 | 1 | 1 | 4 |
| 3 | 1 | 2 | 4 | 8 |
| 4 | 2 | 3 | 9 | 15 |
| 5 | 3 | 3 | 9 | 16 |
| 6 | 3 | 4 | 16 | 23 |
| 7 | 4 | 5 | 25 | 42 |

The import operation does not support time data that includes `Inf` or `NaN` values or signal data that includes `Inf` values. Empty or `NaN` signal values imported from the Excel file render as missing data in the Simulation Data Inspector. All built-in data types are supported.

Multiple Time Vectors

When your data includes signals with different time vectors, the file can include more than one time vector. Every time column must be named `time`. Time columns specify the sample times for signals to the right, up to the next time vector. For example, the first time column defines the time for `signal1` and `signal2`, and the second time column defines the time steps for `signal3`.

| | A | B | C | D | E |
|---|-------------------|----------------------|----------------------|-------------------|----------------------|
| 1 | <code>time</code> | <code>signal1</code> | <code>signal2</code> | <code>time</code> | <code>signal3</code> |
| 2 | 0 | 1 | 1 | 0 | 4 |
| 3 | 1 | 2 | 4 | 2 | 8 |
| 4 | 2 | 3 | 9 | 3 | 15 |
| 5 | 3 | 3 | 9 | 5 | 16 |
| 6 | 3 | 4 | 16 | | |
| 7 | 4 | 5 | 25 | | |

Signal columns must have the same number of data points as the associated time vector.

Signal Metadata

The file can include metadata for signals such as data type, units, and interpolation method. The metadata is used to determine how to plot the data, how to apply unit and data conversions, and how to compute comparison results. For more information about how metadata is used in comparisons, see “How the Simulation Data Inspector Compares Data”.

Metadata for each signal is listed in rows between the signal names and the signal data. You can specify any combination of metadata for each signal. Leave a blank cell for signals with less specified metadata.

| | A | B | C | D | E |
|----|------|-------------|-------------|------|-------------|
| 1 | time | signal1 | signal2 | time | signal3 |
| 2 | | Interp: zoh | | | Interp: zoh |
| 3 | | Type: int8 | Type: int32 | | |
| 4 | | Unit: m | | | Unit: m/s |
| 5 | 0 | 1 | 1 | 0 | 4 |
| 6 | 1 | 2 | 4 | 2 | 8 |
| 7 | 2 | 3 | 9 | 3 | 15 |
| 8 | 3 | 3 | 9 | 5 | 16 |
| 9 | 3 | 4 | 16 | | |
| 10 | 4 | 5 | 25 | | |

Label each piece of metadata according to this table. The table also indicates which tools and operations support each piece of metadata. When an imported file does not specify signal metadata, double data type, linear interpolation, and union synchronization are used.

Property Descriptions

| Signal Property | Label | Values | Simulation Data Inspector Import | Record Block Logging and Simulation Data Inspector Export | Simulink Test Import and Export |
|------------------------|----------|--|----------------------------------|--|---------------------------------|
| Data type | Type: | Built-in data type. | Supported | Supported | Supported |
| Units | Unit: | Supported unit. For example, Unit: m/s specifies units of meters per second. For a list of supported units, enter showunitslist in the MATLAB Command Window. | Supported | Supported | Supported |
| Interpolation method | Interp: | linear, zoh for zero order hold, or none. | Supported | Supported | Supported |
| Synchronization method | Sync: | union or intersection. | Supported | Not Supported <i>Metadata not included in exported file.</i> | Supported |
| Relative tolerance | RelTol: | Percentage, represented as a decimal. For example, RelTol: 0.1 specifies a 10% relative tolerance. | Supported | Not Supported <i>Metadata not included in exported file.</i> | Supported |
| Absolute tolerance | AbsTol: | Numeric value. | Supported | Not Supported <i>Metadata not included in exported file.</i> | Supported |
| Time tolerance | TimeTol: | Numeric value, in seconds. | Supported | Not Supported <i>Metadata not included in exported file.</i> | Supported |

| Signal Property | Label | Values | Simulation Data Inspector Import | Record Block Logging and Simulation Data Inspector Export | Simulink Test Import and Export |
|-------------------|--------------|--|---|--|---------------------------------|
| Leading tolerance | LeadingTol : | Numeric value, in seconds. | Supported <i>Only visible in Simulink Test.</i> | Not Supported <i>Metadata not included in exported file.</i> | Supported |
| Lagging tolerance | LaggingTol : | Numeric Value, in seconds. | Supported <i>Only visible in Simulink Test.</i> | Not Supported <i>Metadata not included in exported file.</i> | Supported |
| Block Path | BlockPath : | Path to the block that generated the signal. | Supported | Supported | Supported |
| Port Index | PortIndex : | Integer. | Supported | Supported | Supported |
| Name | Name : | Signal name | Supported | Not Supported <i>Metadata not included in exported file.</i> | Supported |

User-Defined Data Types

In addition to built-in data types, you can use other labels in place of the `DataType: label` to specify fixed-point, enumerated, alias, and bus data types.

Property Descriptions

| Data Type | Label | Values | Simulation Data Inspector Import | Record Block Logging and Simulation Data Inspector Export | Simulink Test Import and Export |
|-------------|--------|--|--|---|--|
| Enumeration | Enum: | Name of the enumeration class. | Supported <i>Enumeration class definition must be saved on the MATLAB path.</i> | Supported <i>Enumeration class definition must be saved on the MATLAB path.</i> | Supported <i>Enumeration class definition must be saved on the MATLAB path.</i> |
| Alias | Alias: | Name of a Simulink.AliasType object in the MATLAB workspace. | Supported <i>For matrix and complex signals, specify the alias data type on the first channel.</i> | Not Supported | Supported <i>For matrix and complex signals, specify the alias data type on the first channel.</i> |
| Fixed-point | Fixdt: | <ul style="list-style-type: none"> fixdt constructor. Name of a Simulink.NumericType object in the MATLAB workspace. Name of a fixed-point data type as described in "Fixed-Point Numbers in Simulink" on page 35-13. | Supported | Not Supported | Supported |
| Bus | Bus: | Name of a Simulink.Bus object in the MATLAB workspace. | Supported | Not Supported | Supported |

When you specify the type using the name of a Simulink.Bus object and the object is not in the MATLAB workspace, the data still imports from the file. However, individual signals in the bus use data types described in the file rather than data types defined in the Simulink.Bus object.

Complex, Multidimensional, and Bus Signals

You can import and export complex, multidimensional, and bus signals using an Excel file. The signal name for a column of data indicates whether that data is part of a complex, multidimensional, or bus signal. Excel file import and export do not support array of bus signals.

Note When you export data from a nonvirtual bus with variable-size signals to an Excel file, the variable-size signal data is expanded to individual channels, and the hierarchical nature of the data is lost. Data imported from this file is returned as a flat list.

Multidimensional signal names include index information in parentheses. For example, the signal name for a column might be `signal1(2,3)`. When you import data from a file that includes multidimensional signal data, elements in the data not included in the file take zero sample values with the same data type and complexity as the other elements.

Complex signal data is always in real-imaginary format. Signal names for columns containing complex signal data include `(real)` and `(imag)` to indicate which data each column contains. When you import data from a file that includes imaginary signal data without specifying values for the real component of that signal, the signal values for the real component default to zero.

Multidimensional signals can contain complex data. The signal name includes the indication for the index within the multidimensional signal and the real or imaginary tag. For example, `signal1(1,3)(real)`.

Dots in signal names specify the hierarchy for bus signals. For example:

- `bus.y.a`
- `bus.y.b`
- `bus.x`

| | A | B | C | D | E |
|----|------|-------------|-------------|------|-------------|
| 1 | time | bus.y.a | bus.y.b | time | bus.x |
| 2 | | Interp: zoh | | | Interp: zoh |
| 3 | | Type: int8 | Type: int32 | | |
| 4 | | Unit: m | | | Unit: m/s |
| 5 | 0 | 1 | 1 | 0 | 4 |
| 6 | 1 | 2 | 4 | 2 | 8 |
| 7 | 2 | 3 | 9 | 3 | 15 |
| 8 | 3 | 3 | 9 | 5 | 16 |
| 9 | 3 | 4 | 16 | | |
| 10 | 4 | 5 | 25 | | |

Tip When the name of your signal includes characters that could make it appear as though it were part of a matrix, complex signal, or bus, use the Name metadata option to specify the name you want the imported signal to use in the Simulation Data Inspector and Simulink Test.

Function-Call Signals

Signal data specified in columns before the first time column is imported as one or more function-call signals. The data in the column specifies the times at which the function-call signal was enabled. The imported signals have a value of 1 for the times specified in the column. The time values for function-call signals must be double, scalar, and real, and must increase monotonically.

When you export data from the Simulation Data Inspector, function-call signals are formatted the same as other signals, with a time column and a column for signal values.

Simulation Parameters

You can import data for parameter values used in simulation. In the Simulation Data Inspector, the parameter values are shown as signals. Simulink Test uses imported parameter values to specify values for those parameters in the tests it runs based on imported data.

Parameter data is specified using two or three columns. The first column specifies the parameter names, with the cell in the header row for that column labeled **Parameter:**. The second column specifies the value used for each parameter, with the cell in the header row labeled **Value:**. Parameter data may also include a third column that contains the block path associated with each parameter, with the cell in the header row labeled **BlockPath:**. Specify names, values, and block paths for parameters starting in the first row that contains signal data, below rows used to specify signal metadata. For example, this file specifies values for two parameters, X and Y.

| | A | B | C | D | E | F | G |
|----|------|-------------|-------------|------|-------------|-------------------|-----|
| 1 | time | signal1 | signal2 | time | signal3 | Parameter: Value: | |
| 2 | | Interp: zoh | | | Interp: zoh | | |
| 3 | | Type: int8 | Type: int32 | | | | |
| 4 | | Unit: m | | | Unit: m/s | | |
| 5 | 0 | 1 | 1 | 0 | 4 X | | 2 |
| 6 | 1 | 2 | 4 | 2 | 8 Y | | 1.2 |
| 7 | 2 | 3 | 9 | 3 | 15 | | |
| 8 | 3 | 3 | 9 | 5 | 16 | | |
| 9 | 3 | 4 | 16 | | | | |
| 10 | 4 | 5 | 25 | | | | |

Multiple Runs

You can include data for multiple runs in a single file. Within a sheet, you can divide data into runs by labeling data with a simulation number and a source type, such as **Input** or **Output**. Specify the simulation number and source type as additional signal metadata, using the label **Simulation:** for the simulation number and the label **Source:** for the source type. The Simulation Data Inspector uses the simulation number and source type only to determine which signals belong in each run. Simulink Test uses the information to define inputs, parameters, and acceptance criteria for tests to run based on imported data.

You do not need to specify the simulation number and output type for every signal. Signals to the right of a signal with a simulation number and source use the same simulation number and source

until the next signal with a different source or simulation number. For example, this file defines data for two simulations and imports into four runs in the Simulation Data Inspector:

- **Run 1** contains signal1 and signal2.
- **Run 2** contains signal3, X, and Y.
- **Run 3** contains signal4.
- **Run 4** contains signal5.

| | A | B | C | D | E | F | G | H | I | J |
|----|------|---------------|-------------|------|----------------|------------|---------|------|---------------|----------------|
| 1 | time | signal1 | signal2 | time | signal3 | Parameter: | Values: | time | signal4 | signal5 |
| 2 | | Interp: zoh | | | Interp: zoh | | | | | |
| 3 | | Type: int8 | Type: int32 | | | | | | | |
| 4 | | Unit: m | | | Unit: m/s | | | | | |
| 5 | | Simulation: 1 | | | | | | | Simulation: 2 | |
| 6 | | Source: Input | | | Source: Output | | | | Source: Input | Source: Output |
| 7 | 0 | 1 | 1 | 0 | 4 X | | 2 | 1 | 2 | 1 |
| 8 | 1 | 2 | 4 | 2 | 8 Y | | 1.2 | 2 | 6 | 3 |
| 9 | 2 | 3 | 9 | 3 | 15 | | | 3 | 4 | 5 |
| 10 | 3 | 3 | 9 | 5 | 16 | | | 4 | 8 | 7 |
| 11 | 3 | 4 | 16 | | | | | 5 | 10 | 2 |
| 12 | 4 | 5 | 25 | | | | | | | |

You can also use sheets within the Microsoft Excel file to divide the data into runs and tests. When you do not specify simulation number and source information, the data on each sheet is imported into a separate run in the Simulation Data Inspector. When you export multiple runs from the Simulation Data Inspector, the data for each run is saved on a separate sheet. When you import a Microsoft Excel file that contains data on multiple sheets into Simulink Test, you are prompted to specify how to import the data.

See Also

`Simulink.sdi.createRun` | `Simulink.sdi.exportRun`

More About

- “View Data in the Simulation Data Inspector”
- “Import Data from a CSV File into the Simulation Data Inspector”
- “Import Data Using a Custom File Reader”

Configure the Simulation Data Inspector

The Simulation Data Inspector supports a wide range of use cases for analyzing and visualizing data. You can modify preferences in the Simulation Data Inspector to match your visualization and analysis requirements. The preferences that you specify persist between MATLAB sessions.

By specifying preferences in the Simulation Data Inspector, you can configure options such as:

- How signals and metadata are displayed.
- Which data automatically imports from parallel simulations.
- Where prior run data is retained and how much prior data to store.
- How much memory is used during save operations.
- The system of units used to display signals.



To open the Simulation Data Inspector preferences, click Preferences.

Note You can restore all preferences in the Simulation Data Inspector to default values by clicking **Restore Defaults** in the Preferences menu or by using the `Simulink.sdi.clearPreferences` function.

Logged Data Size and Location

By default, simulation data logs to disk with data loaded into memory on demand, and the maximum size of logged data is constrained only by available disk space. You can use the **Disk Management** settings in the Simulation Data Inspector to directly control the size and location of logged data.

The **Record mode** setting specifies whether logged data is retained after simulation. When you change the **Record mode** setting to **View during simulation only**, no logged data is available in the Simulation Data Inspector or the workspace after the simulation completes. Only use this mode when you do not want to save logged data. The **Record mode** setting reverts to **View and record data** each time you start MATLAB. Changing the **Record mode** setting can affect other applications, such as visualization tools. For details, see “View Data Only During Simulation”.

To directly limit the size of logged data, you can specify a minimum amount of free disk space or a maximum size for the logged data. By default, logged data must leave at least 100 MB of free disk space with no maximum size limit. Specify the required disk space and maximum size in GB, and specify 0 to apply no disk space requirement or no maximum size limit.

When you specify a minimum disk space requirement or a maximum size for logged data, you can also specify whether to prioritize retaining data from the current simulation or data from prior simulations when approaching the limit. By default, the Simulation Data Inspector prioritizes retaining data for the current run by deleting data for prior runs. To prioritize retaining prior data, change the **When low on disk space** setting to **Keep prior runs and stop recording**. You see a warning message when prior runs are deleted and when recording is disabled. If recording is disabled due to the size of logged data, you need to change the **Record Mode** back to **View and**

record data to continue logging data, after you have freed up disk space. For more information, see “Specify a Minimum Disk Space Requirement or Maximum Size for Logged Data”.

The **Storage Mode** setting specifies whether to log data to disk or to memory. By default, data logs to disk. When you configure a parallel worker to log data to memory, data transfer back to the host is not supported. Logging data to memory is not supported for rapid accelerator simulations or models deployed using Simulink Compiler.

You can also specify the location of the temporary file that stores logged data. By default, data logs to the temporary files directory on your computer. You may change the file location when you need to log large amounts of data and a secondary drive provides more storage capacity. Logging data to a network location can degrade performance.

Programmatic Use

You can programmatically configure and check each preference value.

| Preference | Functions |
|-------------------------------|--|
| Record mode | Simulink.sdi.setRecordData Simulink.sdi.getRecordData |
| Required Free Space | Simulink.sdi.setRequiredFreeSpace Simulink.sdi.getRequiredFreeSpace |
| Max Disk Usage | Simulink.sdi.setMaxDiskUsage Simulink.sdi.getMaxDiskUsage |
| When low on disk space | Simulink.sdi.setDeleteRunsOnLowSpace Simulink.sdi.getDeleteRunsOnLowSpace |
| Storage Mode | Simulink.sdi.setStorageMode Simulink.sdi.getStorageMode |
| Storage Location | Simulink.sdi.setStorageLocation Simulink.sdi.getStorageLocation |

Archive Behavior and Run Limit

When you run multiple simulations in a single MATLAB session, the Simulation Data Inspector retains results from each simulation so you can analyze the results together. Use the Simulation Data Inspector archive to manage runs in the user interface and control the number of runs the Simulation Data Inspector retains.


You can configure a limit for the number of runs to retain in the archive and whether the Simulation Data Inspector automatically moves prior runs into the archive.

Manage Runs Using the Archive

By default, the Simulation Data Inspector automatically archives simulation runs. When you simulate a model, the prior simulation run moves to the archive, and the Simulation Data Inspector updates the view to show data for aligned signals in the current run.

The archive does not impose functional limitations on the runs and signals it contains. You can plot signals from the archive, and you can use runs and signals in the archive in comparisons. You can drag runs of interest from the archive to the work area and vice versa whether **Automatically Archive** is selected or disabled.

To prevent the Simulation Data Inspector from automatically moving prior simulations runs to the archive, clear the **Automatically archive** setting. With automatic archiving disabled, the Simulation Data Inspector does not move prior runs into the **Archive** pane or automatically update plots to display data from the current simulation.

Tip To manually delete the contents of the archive, click Delete archived runs .

Control Number of Runs Retained in Simulation Data Inspector

You can specify a limit for the number of runs to retain in the archive. When the number of runs in the archive reaches the limit, the Simulation Data Inspector deletes runs in the archive on a first-in, first-out basis.

The run limit applies only to runs in the archive. For the Simulation Data Inspector to automatically limit the data it retains by deleting old runs, select **Automatically archive** and specify a size limit.

By default, the Simulation Data Inspector retains the last 20 runs moved to the archive. To remove the limit, select **No limit**. To specify the maximum number of runs to store in the archive, select **Last n runs** and enter the limit. A warning occurs if you specify a limit that would delete runs already in the archive.

Programmatic Use

You can programmatically configure and check the archive behavior and run limit.

| Preference | Functions |
|------------------------------|--|
| Automatically archive | <code>Simulink.sdi.setAutoArchiveMode</code> <code>Simulink.sdi.getAutoArchiveMode</code> |
| Size | <code>Simulink.sdi.setArchiveRunLimit</code> <code>Simulink.sdi.getArchiveRunLimit</code> |

Incoming Run Names and Location

You can configure how the Simulation Data Inspector handles incoming runs from import or simulation. You can choose whether new runs are added at the top of the work area or the bottom and specify a naming rule to use for runs created from simulation.

By default, the Simulation Data Inspector adds new runs below prior runs in the work area. The **Archive** settings also affect the location of runs. By default, prior runs are moved to the archive when a new simulation run is created.

The run naming rule is used to name runs created from simulation. You can create the run naming rule using a mix of literal text that is used in the run name as-is and one or more tokens that represent metadata about the run. By default, the Simulation Data Inspector names runs using the run index and model name: Run <run_index>: <model_name>.

Tip To rename an existing run, double-click the name in the work area and enter the new name, or modify the run name in the **Properties** pane.

Programmatic Use

You can programmatically configure and check incoming run names and locations.

| Preference | Functions |
|---------------------|---|
| Add New Runs | Simulink.sdi.appendRunToTop Simulink.sdi.getAppendRunToTop |
| Naming Rule | Simulink.sdi.setRunNamingRule Simulink.sdi.getRunNamingRule Simulink.sdi.resetRunNamingRule |

Signal Metadata to Display

You can control which signal metadata is displayed in the work area of the **Inspect** pane and in the results section on the **Compare** pane in the Simulation Data Inspector. You specify the metadata to display separately for each pane using the **Table Columns** preferences in the **Inspect** and **Compare** sections of the Preferences dialog, respectively.

Inspect Pane

By default, the signal name and the line style and color used to plot the signal are displayed on the **Inspect** pane. To display different or additional metadata in the work area on the **Inspect** pane, select the check box next to each piece of metadata you want to display in the **Table Columns** preference in the **Inspect** section. You can always view complete metadata for the selected signal in the **Inspect** pane using the **Properties** pane.

Note Metadata displayed in the work area on **Inspect** pane is included when you generate a report of plotted signals. You can also specify metadata to include in the report regardless of what is displayed in the work area when you create the report programmatically using the `Simulink.sdi.report` function.

Compare Pane

By default, the **Compare** pane shows the signal name, the absolute and relative tolerances used in the signal comparison, and the maximum difference from the comparison result. To display different or additional metadata in the results on the **Compare** pane, select the check box next to each piece of metadata you want to display in the **Table Columns** preference in the **Compare** section. You can always view complete metadata for the signals compared for a selected signal result using the **Properties** pane, where metadata that differs between the compared signals is highlighted. Signal metadata displayed on the **Compare** pane does not affect the contents of comparison reports.

Signal Selection on the Inspect Pane

You can configure how you select signals to plot on the selected subplot in the Simulation Data Inspector. By default, you use check boxes next to each signal to plot. You can also choose to plot signals based on selection in the work area. Use **Check Mode** when creating views and visualizations that represent findings and analysis of a data set. Use **Browse Mode** to quickly view and analyze data sets with a large number of signals.

For more information about creating visualizations using **Check Mode**, see “Create Plots Using the Simulation Data Inspector”.

For more information about using **Browse Mode**, see “Visualize Many Logged Signals”.

Note To use **Browse Mode**, your layout must include only **Time Plot** visualizations.

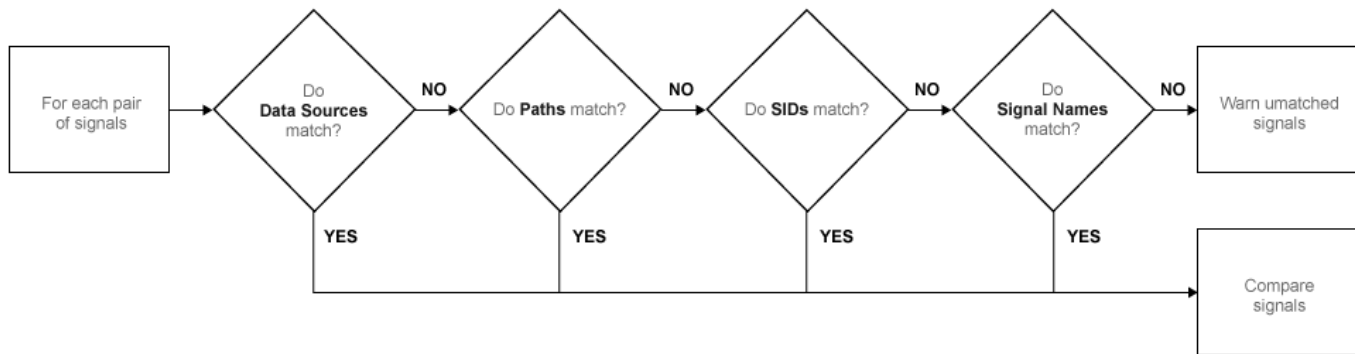
How Signals Are Aligned for Comparison

When you compare runs using the Simulation Data Inspector, the comparison algorithm pairs signals for signal comparison through a process called alignment. You can align signals between the compared runs using one or more of the signal properties shown in the table.

| Property | Description |
|-------------|---|
| Data Source | Path of the variable in the MATLAB workspace for data imported from the workspace |
| Path | Block path for the source of the data in its model |
| SID | Automatically assigned Simulink identifier |
| Signal Name | Name of the signal |

You can specify the priority for each piece of metadata used for alignment. The **Align By** field specifies the highest priority property used to align signals. The priority drops with each subsequent **Then By** field. You must specify a primary alignment property in the **Align By** field, but you can leave any number of **Then By** fields blank.

By default, the Simulation Data Inspector aligns signals between runs according to this flow chart.



For more information about configuring comparisons in the Simulation Data Inspector, see “How the Simulation Data Inspector Compares Data”.

Colors Used to Display Comparison Results

You can configure the colors used to display comparison results using the Simulation Data Inspector preferences. You can specify whether to use the signal color from the **Inspect** pane or a fixed color for the baseline and compared signals. You can also choose colors for the tolerance and the difference signal.

By default, the Simulation Data Inspector displays comparison results using fixed colors for the baseline and compared signals. Using a fixed color allows you to avoid the baseline signal color and compared signal color being either the same or too similar to distinguish.

Signal Grouping

You can specify how to group signals within a run in the Simulation Data Inspector. The preferences apply to both the **Inspect** and **Compare** panes and comparison reports. You can group signals by:

- **Domain** — Signal type. For example, signals created by signal logging have a domain of **Signal**, while signals created from logging model outputs have a domain of **Outputs**.
- **Physical System Hierarchy** — Signal Simscape™ physical system hierarchy. The option to group by physical system hierarchy is available when you have a Simscape license.
- **Data Hierarchy** — Signal location within structured data. For example, data hierarchy grouping reflects the hierarchy of a bus.
- **Model Hierarchy** — Signal location within model hierarchy. Grouping by model hierarchy can be helpful when you log data from a model that includes model or subsystem references.

Grouping signals adds rows for the hierarchical nodes, which you can expand to show the signals within that node. By default, the Simulation Data Inspector groups signals by domain, then by physical system hierarchy (if you have a Simscape license), and then by data hierarchy.

To remove grouping and display a flat list of signals in each run, select **None** for all grouping options.

Programmatic Use

To specify how to group signals programmatically, use the `Simulink.sdi.setTableGrouping` function.

Data to Stream from Parallel Simulations

When you run parallel simulations using the `parsim` function, you can stream logged simulation data to the Simulation Data Inspector. A dot next to the run name in the **Inspect** pane indicates the status of the simulation that corresponds to the run, so you can monitor simulation progress while visualizing the streamed data. You can control whether data streams from a parallel simulation based on the type of worker the data comes from.

By default, the Simulation Data Inspector is configured for manual import of data from parallel workers. You can use the Simulation Data Inspector programmatic interface to inspect the data on the worker and decide whether to send it to the client Simulation Data Inspector for further analysis. To manually move data from a parallel worker to the Simulation Data Inspector, use the `Simulink.sdi.sendWorkerRunToClient` function.

You may want to automatically stream data from parallel simulations that run on local workers or on local and remote workers. Streaming data from both local and remote workers may affect simulation performance, depending on how many simulations you run and how much data you log. When you choose to stream data from local workers or all parallel workers, all logged simulation data automatically shows in the Simulation Data Inspector.

Programmatic Use

You can configure Simulation Data Inspector support for parallel worker data programmatically using the `Simulink.sdi.enablePCTSupport` function.

Options for Saving and Loading Session Files

You can specify a maximum amount of memory to use while loading or saving a session file. By default, the Simulation Data Inspector uses a maximum of 100 MB of memory when you load or save a session file. You can specify a memory use limit as low as 50 MB.

To reduce the size of the saved session file, you can specify a compression option.

- **None** — Do not compress saved data.
- **Normal** — Compress the saved file as much as possible.
- **Fastest** — Compress the saved file less than **Normal** compression for faster save time.

Signal Display Units

Signals in the Simulation Data Inspector have two units properties: stored units and display units. The stored units represent the units of the data saved to disk. The display units specify how the Simulation Data Inspector displays the data. You can configure the Simulation Data Inspector to use a system of units to define the display units for all signals. You can choose either the **SI** or **US Customary** system of units, or you can display data using its stored units.

When you use a system of units to define display units for signals in the Simulation Data Inspector, the display units update for any signal with display units that are not valid for that unit system. For example, if you select **SI** units, the display units for a signal may update from ft to m.

Note The system of units you choose to use in the Simulation Data Inspector does not affect the stored units for any signal. You can convert the stored units for a signal using the `convertUnits` function. Conversion may result in loss of precision.

In addition to selecting a system of units, you can specify override units so that all signals of a given measurement type are displayed using consistent units. For example, if you want to visualize all signals that represent weight using units of kg, specify kg as an override unit.

Tip For a list of units supported by Simulink, enter `showunitslist` in the MATLAB Command Window.

You can also modify the display units for a specific signal using the **Properties** pane. For more information, see “Modify Signal Properties in the Simulation Data Inspector”.

Programmatic Use

Configure the unit system and override units using the `Simulink.sdi.setUnitSystem` function. You can check the current units preferences using the `Simulink.sdi.getUnitSystem` function.

See Also

Functions

`Simulink.sdi.clearPreferences` | `Simulink.sdi.setRunNamingRule` |
`Simulink.sdi.setTableGrouping` | `Simulink.sdi.enablePCTSupport` |
`Simulink.sdi.setArchiveRunLimit` | `Simulink.sdi.setAutoArchiveMode`

More About

- “Iterate Model Design Using the Simulation Data Inspector”
- “How the Simulation Data Inspector Compares Data”
- “Compare Simulation Data”
- “Create Plots Using the Simulation Data Inspector”
- “Modify Signal Properties in the Simulation Data Inspector”

How the Simulation Data Inspector Compares Data

You can tailor the Simulation Data Inspector comparison process to fit your requirements in multiple ways. When comparing runs, the Simulation Data Inspector:

- 1 Aligns signal pairs in the **Baseline** and **Compare To** runs based on the **Alignment** settings.




The Simulation Data Inspector does not compare signals that it cannot align.

- 2 Synchronizes aligned signal pairs according to the specified **Sync Method**.

Values for time points added in synchronization are interpolated according to the specified **Interpolation Method**.

- 3 Computes the difference of the signal pairs.
- 4 Compares the difference result against specified tolerances.

When the comparison run completes, the results of the comparison are displayed in the navigation pane.

| Status | Comparison Result |
|---|---|
|  | Difference falls within the specified tolerance. |
|  | Difference violates specified tolerance. |
|  | The signal does not align with a signal from the Compare To run. |

When you compare signals with differing time intervals, the Simulation Data Inspector compares the signals on their overlapping interval.

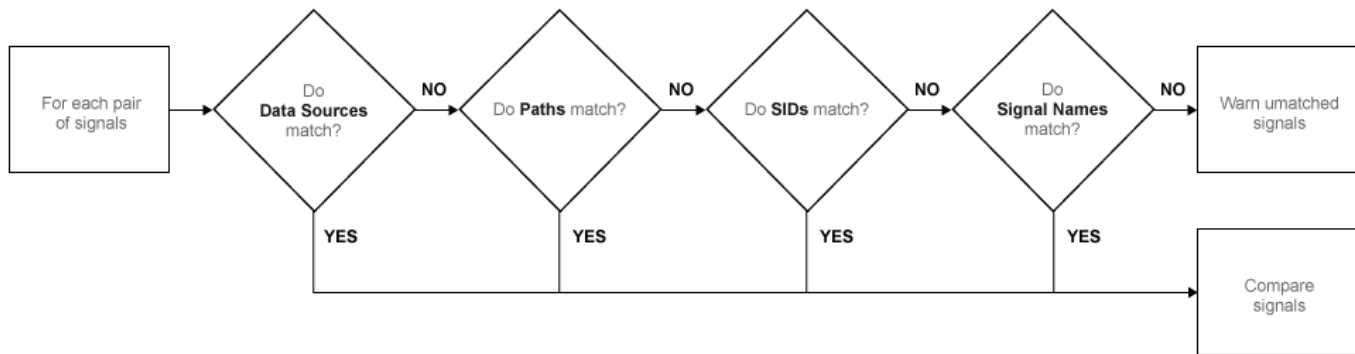
Signal Alignment

In the alignment step, the Simulation Data Inspector decides which signal from the **Compare To** run pairs with a given signal in the **Baseline** run. When you compare signals with the Simulation Data Inspector, you complete the alignment step by selecting the **Baseline** and **Compare To** signals.

The Simulation Data Inspector aligns signals using a combination of their Data Source, Path, SID, and Signal Name properties.

| Property | Description |
|-------------|---|
| Data Source | Path of the variable in the MATLAB workspace for data imported from the workspace |
| Path | Block path for the source of the data in its model |
| SID | Automatically assigned Simulink identifier |
| Signal Name | Name of the signal in the model |

With the default alignment settings, the Simulation Data Inspector aligns signals between runs according to this flow chart.

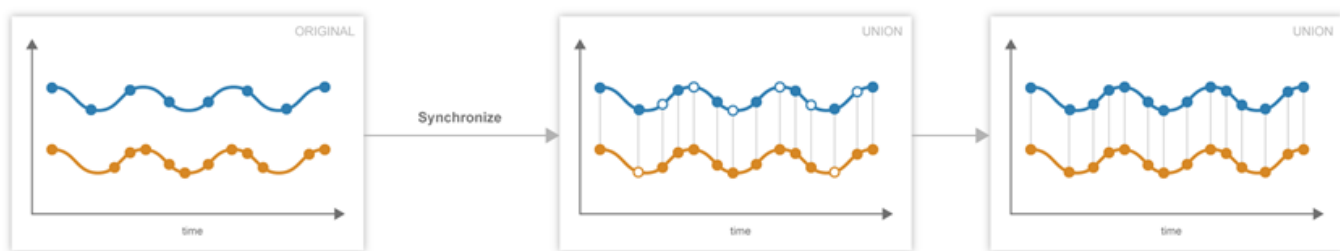


You can specify the priority for each of the signal properties used for alignment in the Simulation Data Inspector **Preferences**. The **Align By** field specifies the highest priority property used to align signals. The priority drops with each subsequent **Then By** field. You must specify a primary alignment property in the **Align By** field, but you can leave any number of the **Then By** fields blank.

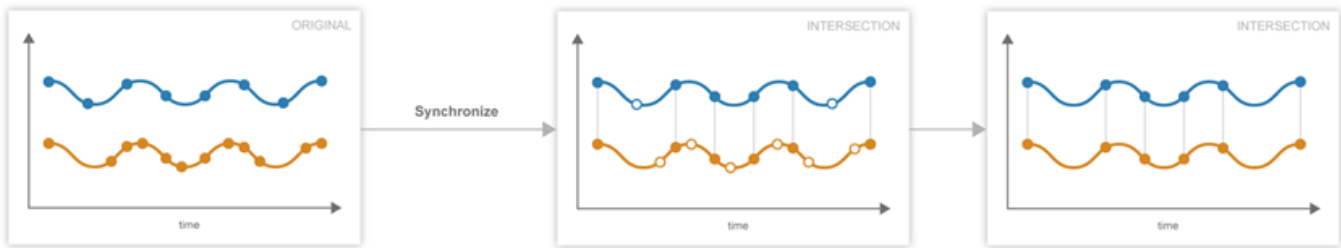
Synchronization

Often, signals that you want to compare don't contain the exact same set of time points. The synchronization step in Simulation Data Inspector comparisons resolves discrepancies in signals' time vectors. You can choose union or intersection as the synchronization method.

When you specify union synchronization, the Simulation Data Inspector builds a time vector that includes every sample time between the two signals. For each sample time not originally present in either signal, the Simulation Data Inspector interpolates the value. The second graph in the illustration shows the union synchronization process, where the Simulation Data Inspector identifies samples to add in each signal, represented by the unfilled circles. The final plot shows the signals after the Simulation Data Inspector has interpolated values for the added time points. The Simulation Data Inspector computes the difference using the signals in the final graph, so that the computed difference signal contains all the data points between the signals.



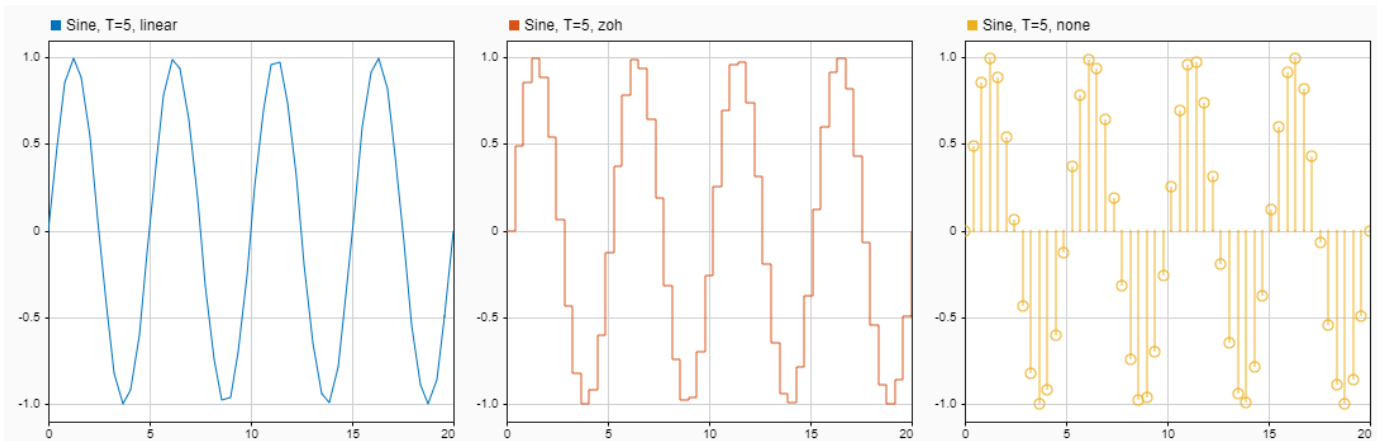
When you specify intersection synchronization, the Simulation Data Inspector uses only the sample times present in both signals in the comparison. In the second graph, the Simulation Data Inspector identifies samples that do not have a corresponding sample for comparison, shown as unfilled circles. The final graph shows the signals used for the comparison, without the samples identified in the second graph.



The choice between the synchronization options involves a trade off between speed and accuracy. The interpolation required by `union` synchronization takes time, but provides a more precise result. When you use `intersection` synchronization, the comparison finishes quickly because the Simulation Data Inspector computes the difference for fewer data points and does not interpolate. However, some data is discarded and precision is lost with `intersection` synchronization.

Interpolation

The interpolation property of a signal determines how the Simulation Data Inspector displays the signal and how additional data values are computed in synchronization. You can choose to interpolate your data with a zero-order hold (`zoh`) or a linear approximation. You can also specify no interpolation.



When you specify `zoh` or `none` for the **Interpolation Method**, the Simulation Data Inspector replicates the data of the previous sample for interpolated sample times. When you specify `linear` interpolation, the Simulation Data Inspector uses samples on either side of the interpolated point to linearly approximate the interpolated value. Typically, discrete signals use `zoh` interpolation and continuous signals use `linear` interpolation. You can specify the **Interpolation Method** for your signals in the signal properties.

Tolerance Specification

The Simulation Data Inspector allows you to specify the scope and value of the tolerance for your signal. You can define a tolerance band using any combination of absolute, relative, and time tolerance values, and you can specify whether the specified tolerance applies to an individual signal or to all the signals in a run.

Tolerance Scope

In the Simulation Data Inspector, you can specify the tolerance for your data globally or for an individual signal. Global tolerance values apply to all signals in a run that do not have **Override Global Tol** set to yes. You can specify global tolerance values for your data at the top of the graphical viewing area in the **Compare** view. To specify signal specific tolerance values, edit the signal properties and ensure the **Override Global Tol** property is set to yes.

Tolerance Computation

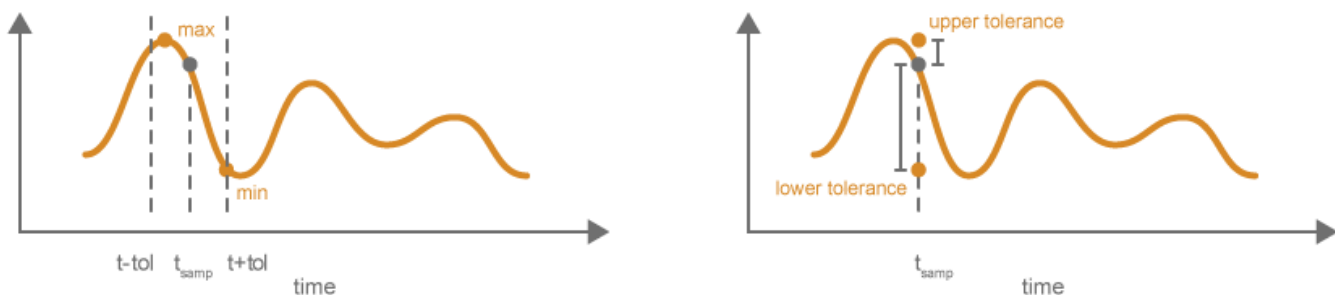
In the Simulation Data Inspector, you can specify a tolerance band for your run or signal using a combination of absolute, relative, and time tolerance values. When you specify the tolerance for your run or signal using multiple types of tolerances, each tolerance can yield a different answer for the tolerance at each point. The Simulation Data Inspector computes the overall tolerance band by selecting the most lenient tolerance result for each data point.

When you define your tolerance using only the absolute and relative tolerance properties, the Simulation Data Inspector computes the tolerance for each point as a simple maximum.

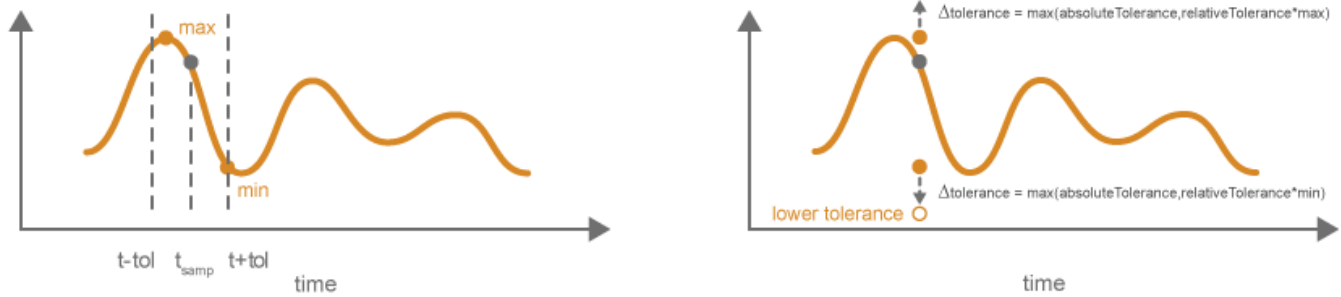
```
tolerance = max(absoluteTolerance, relativeTolerance*abs(baselineData));
```

The upper boundary of the tolerance band is formed by adding **tolerance** to the **Baseline** signal. Similarly, the Simulation Data Inspector computes the lower boundary of the tolerance band by subtracting **tolerance** from the **Baseline** signal.

When you specify a time tolerance, the Simulation Data Inspector evaluates the time tolerance first, over a time interval defined as $[(t_{\text{samp}} - \text{tol}), (t_{\text{samp}} + \text{tol})]$ for each sample. The Simulation Data Inspector builds the lower tolerance band by selecting the minimum point on the interval for each sample. Similarly, the maximum point on the interval defines the upper tolerance for each sample.



If you specify a tolerance band using an absolute or relative tolerance in addition to a time tolerance, the Simulation Data Inspector applies the time tolerance first, and then applies the absolute and relative tolerances to the maximum and minimum points selected with the time tolerance.



`upperTolerance = max + max(absoluteTolerance, relativeTolerance*max)`

`lowerTolerance = min - max(absoluteTolerance, relativeTolerance*min)`

Limitations

The Simulation Data Inspector does not support comparing:

- Signals of data types `int64` or `uint64`.
- Variable-size signals.

See Also

Related Examples

- “Compare Simulation Data”

Save and Share Simulation Data Inspector Data and Views

After you inspect, analyze, or compare your data in the Simulation Data Inspector, you can share your results with others. The Simulation Data Inspector provides several options for sharing and saving your data and results, depending on your needs. With the Simulation Data Inspector, you can:

- Save your data and layout modifications in a Simulation Data Inspector session.
- Share your layout modifications in a Simulation Data Inspector view.
- Share images and figures of plots you create in the Simulation Data Inspector.
- Create a Simulation Data Inspector report.
- Export data to the workspace.
- Export data to a file.

Save and Load Simulation Data Inspector Sessions

If you want to save or share data along with a configured view in the Simulation Data Inspector, save your data and settings in a Simulation Data Inspector session. You can save sessions as MAT- or MLDATX-files. The default format is MLDATX. When you save a Simulation Data Inspector session, the session file contains:

- All runs, data, and properties from the **Inspect** pane, including which run is the current run and which runs are in the archive.
- Plot display selection for signals in the **Inspect** pane.
- Subplot layout and line style and color selections.

Note Comparison results and global tolerances are not saved in Simulation Data Inspector sessions.

To save a Simulation Data Inspector session:

- 1 Hover over the save icon on the left side bar. Then, click **Save As**.



- 2 Name the file.
- 3 Browse to the location where you want to save the session, and click **Save**.


For large datasets, a status overlay in the bottom right of the graphical viewing area displays information about the progress of the save operation and allows you to cancel the save operation.

The **Save** tab of the Simulation Data Inspector preferences menu on the left side bar allows you to configure options related to save operations for MLDATX-files. You can set a limit as low as 50MB on the amount of memory used for the save operation. You can also select one of three **Compression** options:

- **None**, the default, applies no compression during the save operation.

- **Normal** creates the smallest file size.
- **Fastest** creates a smaller file size than you would get by selecting **None**, but provides a faster save time than **Normal**.



To load a Simulation Data Inspector session, click the open icon  on the left side bar. Then, browse to select the MLDATX-file you want to open, and click **Open**.

Alternatively, you can double-click the MLDATX-file. MATLAB and the Simulation Data Inspector open if they are not already open.

When the Simulation Data Inspector already contains runs and you open a session, all of the runs in the session move to the archive. The view updates to show plotted signals from the session file. You can drag runs between the work area and archive as desired.


When the Simulation Data Inspector does not contain runs and you open a session, the Simulation Data Inspector puts runs in the work area and archive as specified in the file.

Share Simulation Data Inspector Views


When you have different sets of data that you want to visualize the same way, you can save a view. A view saves the layout and appearance characteristics of the Simulation Data Inspector without saving the data. Specifically, a view saves:

- Plot visualization type, layout, axis ranges, linking characteristics, and normalized axes
- Location of signals in the plots, including plotted signals in the archive
- Signal grouping and columns on display in the **Inspect** pane
- Signal color and line styling

To save a view:

- 1 Click Visualizations and layouts .
- 2 In **Saved Views**, click **Save current view**.
- 3 In the dialog box, specify a name for the view and browse to the location where you want to save the MLDATX-file.
- 4 Click **Save**.


To load a view:

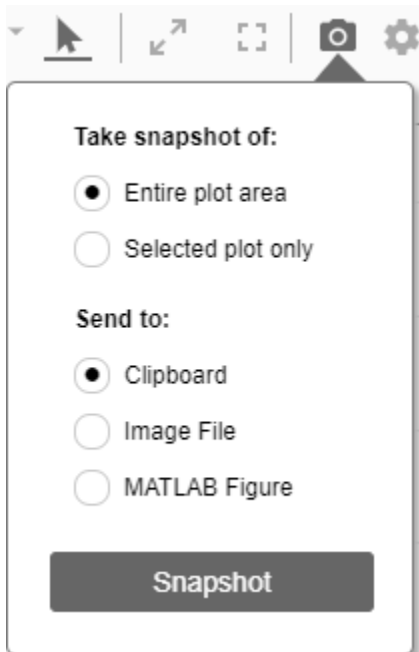
- 1 Click Visualizations and layouts .
- 2 In **Saved Views**, click **Open saved view**.
- 3 Browse to the view you would like to load, and click **Open**.

Share Simulation Data Inspector Plots

Use the snapshot feature to share the plots you generate in the Simulation Data Inspector. You can export your plots to the clipboard to paste into a document, as an image file, or to a MATLAB figure.

You can choose to capture the entire plot area, including all subplots in the plot area, or to capture only the selected subplot.

Click the camera icon  on the toolbar to access the snapshot menu. Use the radio buttons to select the area you want to share and how you want to share the plot. After you make your selections, click **Snapshot** to export the plot.



If you create an image, select where you would like to save the image in the file browser.

You can create snapshots of your plots in the Simulation Data Inspector programmatically using `Simulink.sdi.snapshot`.

Create Simulation Data Inspector Report

To generate documentation of your results quickly, create a Simulation Data Inspector report. You can create a report of your data in either the **Inspect** or the **Compare** pane. The report is an HTML file that includes information about all the signals and plots in the active pane. The report includes all signal information displayed in the signal table in the navigation pane. For more information about configuring the table, see “Inspect Metadata”.

To generate a Simulation Data Inspector Report:

1

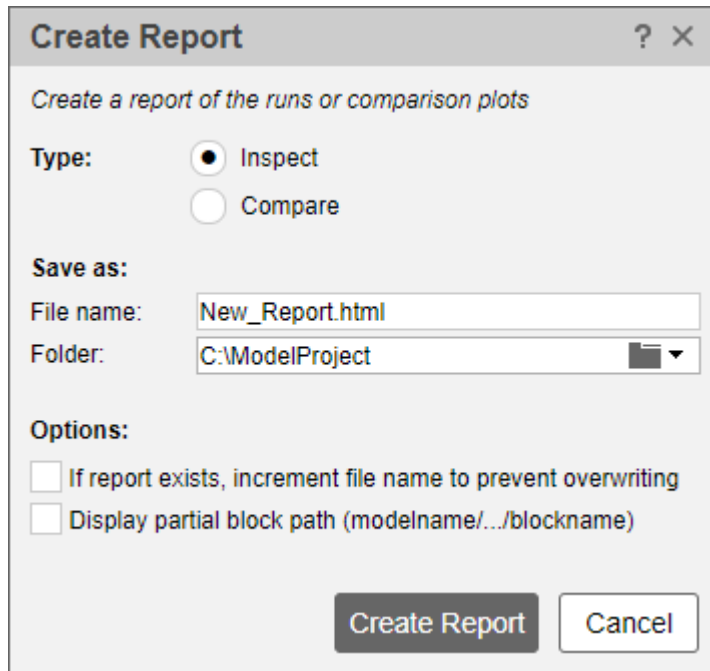


Click the create report icon on the left bar.

2 Specify the type of report you want to create.

- Select **Inspect** to include the plots and signals from the **Inspect** pane.

- Select **Compare** to include the data and plots from the **Compare** pane. When you generate a **Compare Runs** report, you can choose to **Report only mismatched signals** or to **Report all signals**. If you select **Report only mismatched signals**, the report shows only signal comparisons that are not within the specified tolerances.



- 3 Specify a **File name** for the report, and navigate to the **Folder** where you want to save the report.
- 4 Click **Create Report**.

The generated report automatically opens in your default browser.

Export Data to the Workspace or a File

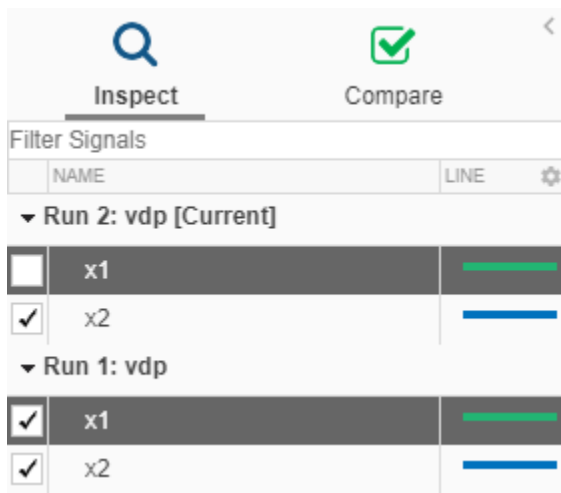
You can use the Simulation Data Inspector to export data to the base workspace, a MAT file, or a Microsoft Excel file. You can export a selection of runs and signals, runs in the work area, or all runs in the **Inspect** pane, including the **Archive**.

When you export a selection of runs and signals, make the selection of data to export before clicking



the export button .

Only the selected runs and signals are exported. In this example, only the x1 signals from Run 1 and Run 2 are exported. The check box selections for plotting data do not affect whether a signal is exported.



When you export a single signal to the workspace or a MAT file, the signal is exported to a `timeseries` object. Data exported to the workspace or a MAT file for a run or multiple signals is stored as a `Simulink.SimulationData.Dataset` object.

To export data to a file, select the **File** option in the **Export** dialog. You can specify a file name and browse to the location where you want to save the exported file. When you export data to a MAT file, a single exported signal is stored as a `timeseries` object, and runs or multiple signals are stored as a `Simulink.SimulationData.Dataset` object. When you export data to a Microsoft Excel file, the data is stored using the format described in “Microsoft Excel Import, Export, and Logging Format”.

To export to a Microsoft Excel file, select the XLSX extension from the drop-down. When you export data to a Microsoft Excel file, you can specify additional options for the format of the data in the exported file. If the file name you provided already exists, you can choose to overwrite the entire file or to only overwrite sheets containing data that corresponds to the exported data. You can also choose which metadata to include and whether signals with identical time data share a time column in the exported file.

Export Video Signal to an MP4 File

You can export a 2D or 3D signal that contains RGB or monochrome video data to an MP4 file using the Simulation Data Inspector. For example, when you log a video signal in a simulation, you can export the data to an MP4 file and view the video using a video player. To export a video signal to an MP4 file:

- 1 Select the signal you want to export.

2



Click **Export** in the toolbar on the left or right-click the signal and select **Export**.

- 3 In the Export dialog box, choose to export **Selected runs and signals** to a file.
- 4 Specify a file name and the path to the location where you want to save the file.
- 5 Select **MP4 video file** from the list and click **Export**.

For the option to export to an MP4 file to be available:

- You must export only one signal at a time.
- The selected signal must be 2D or 3D and contain RGB or monochrome video data.
- The selected signal must be represented in the Simulation Data Inspector as a single signal with multidimensional sample values.

You may need to convert the signal representation before exporting the signal data. For more information, see “Analyze Multidimensional Signal Data”.

- The data type for the signal values must be `double`, `single`, or `uint8`.

Exporting a video signal to an MP4 file is not supported for Linux operating systems.

See Also

Functions

`Simulink.sdi.saveView`

Related Examples

- “View Data in the Simulation Data Inspector”
- “Inspect Simulation Data”
- “Compare Simulation Data”

Inspect and Compare Data Programmatically

You can harness the capabilities of the Simulation Data Inspector from the MATLAB command line using the Simulation Data Inspector API.

The Simulation Data Inspector organizes data in runs and signals, assigning a unique numeric identification to each run and signal. Some Simulation Data Inspector API functions use the run and signal IDs to reference data, rather than accepting the run or signal itself as an input. To access the run IDs in the workspace, you can use `Simulink.sdi.getAllRunIDs` or `Simulink.sdi.getRunIDByIndex`. You can access signal IDs through a `Simulink.sdi.Run` object using the `getSignalIDByIndex` method.

The `Simulink.sdi.Run` and `Simulink.sdi.Signal` classes provide access to your data and allow you to view and modify run and signal metadata. You can modify the Simulation Data Inspector preferences using functions like `Simulink.sdi.setSubPlotLayout`, `Simulink.sdi.setRunNamingRule`, and `Simulink.sdi.setMarkersOn`. To restore the Simulation Data Inspector's default settings, use `Simulink.sdi.clearPreferences`.

Create a Run and View the Data

This example shows how to create a run, add data to it, and then view the data in the Simulation Data Inspector.

Create Data for the Run

Create `timeseries` objects to contain data for a sine signal and a cosine signal. Give each `timeseries` object a descriptive name.

```
time = linspace(0,20,100);

sine_vals = sin(2*pi/5*time);
sine_ts = timeseries(sine_vals,time);
sine_ts.Name = 'Sine, T=5';

cos_vals = cos(2*pi/8*time);
cos_ts = timeseries(cos_vals,time);
cos_ts.Name = 'Cosine, T=8';
```

Create a Run and Add Data

Use the `Simulink.sdi.view` function to open the Simulation Data Inspector.

```
Simulink.sdi.view
```

To import data into the Simulation Data Inspector from the workspace, create a `Simulink.sdi.Run` object using the `Simulink.sdi.Run.create` function. Add information about the run to its metadata using the `Name` and `Description` properties of the `Run` object.

```
sinusoidsRun = Simulink.sdi.Run.create;
sinusoidsRun.Name = 'Sinusoids';
sinusoidsRun.Description = 'Sine and cosine signals with different frequencies';
```

Use the `add` function to add the data you created in the workspace to the empty run.


```
add(sinusoidsRun, 'vars', sine_ts, cos_ts);
```

Plot the Data in the Simulation Data Inspector

Use the `getSignalByIndex` function to access `Simulink.sdi.Signal` objects that contain the signal data. You can use the `Simulink.sdi.Signal` object properties to specify the line style and color for the signal and plot it in the Simulation Data Inspector. Specify the `LineColor` and `LineDashed` properties for each signal.

```
sine_sig = getSignalByIndex(sinusoidsRun,1);
sine_sig.LineColor = [0 0 1];
sine_sig.LineDashed = '-.';
```

```
cos_sig = sinusoidsRun.getSignalByIndex(2);
cos_sig.LineColor = [0 1 0];
cos_sig.LineDashed = '--';
```

Use the `Simulink.sdi.setSubPlotLayout` function to configure a 2-by-1 subplot layout in the Simulation Data Inspector plotting area. Then use the `plotOnSubPlot` function to plot the sine signal on the top subplot and the cosine signal on the lower subplot.

```
Simulink.sdi.setSubPlotLayout(2,1);
```

```
plotOnSubPlot(sine_sig,1,1,true);
plotOnSubPlot(cos_sig,2,1,true);
```

Close the Simulation Data Inspector and Save Your Data

When you have finished inspecting the plotted signal data, you can close the Simulation Data Inspector and save the session to an MLDATX file.

```
Simulink.sdi.close('sinusoids.mldatx')
```

Compare Two Signals in the Same Run

You can use the Simulation Data Inspector programmatic interface to compare signals within a single run. This example compares the input and output signals of an aircraft longitudinal controller.

First, load the session that contains the data.

```
Simulink.sdi.load('AircraftExample.mldatx');
```

Use the `Simulink.sdi.Run.getLatest` function to access the latest run in the data.

```
aircraftRun = Simulink.sdi.Run.getLatest;
```

Then, you can use the `Simulink.sdi.getSignalsByName` function to access the `Stick` signal, which represents the input to the controller, and the `alpha, rad` signal that represents the output.

```
stick = getSignalsByName(aircraftRun, 'Stick');
alpha = getSignalsByName(aircraftRun, 'alpha, rad');
```

Before you compare the signals, you can specify a tolerance value to use for the comparison. Comparisons use tolerance values specified for the baseline signal in the comparison, so set an absolute tolerance value of 0.1 on the `Stick` signal.

```
stick.AbsTol = 0.1;
```

Now, compare the signals using the `Simulink.sdi.compareSignals` function. The `Stick` signal is the baseline, and the `alpha_rad` signal is the signal to compare against the baseline.

```
comparisonResults = Simulink.sdi.compareSignals(stick.ID,alpha.ID);
match = comparisonResults.Status
```

```
match =
    ComparisonSignalStatus enumeration
        OutOfTolerance
```

The comparison result is out of tolerance. You can use the `Simulink.sdi.view` function to open the Simulation Data Inspector to view and analyze the comparison results.

Compare Runs with Global Tolerance

You can specify global tolerance values to use when comparing two simulation runs. Global tolerance values are applied to all signals within the run. This example shows how to specify global tolerance values for a run comparison and how to analyze and save the comparison results.

First, load the session file that contains the data to compare. The session file contains data for four simulations of an aircraft longitudinal controller. This example compares data from two runs that use different input filter time constants.

```
Simulink.sdi.load('AircraftExample.mldatx');
```

To access the run data to compare, use the `Simulink.sdi.getAllRunIDs` function to get the run IDs that correspond to the last two simulation runs.

```
runIDs = Simulink.sdi.getAllRunIDs;
runID1 = runIDs(end - 1);
runID2 = runIDs(end);
```

Use the `Simulink.sdi.compareRuns` function to compare the runs. Specify a global relative tolerance value of `0.2` and a global time tolerance value of `0.5`.

```
runResult = Simulink.sdi.compareRuns(runID1,runID2,'reltol',0.2,'timetol',0.5);
```

Check the `Summary` property of the returned `Simulink.sdi.DiffRunResult` object to see whether signals were within the tolerance values or out of tolerance.

```
runResult.Summary
ans = struct with fields:
    OutOfTolerance: 0
    WithinTolerance: 3
    Unaligned: 0
    UnitsMismatch: 0
    Empty: 0
    Canceled: 0
    EmptySynced: 0
    DataTypeMismatch: 0
```

```

    TimeMismatch: 0
  StartStopMismatch: 0
    Unsupported: 0

```

All three signal comparison results fell within the specified global tolerance.

You can save the comparison results to an MLDATX file using the `saveResult` function.

```
saveResult(runResult, 'InputFilterComparison');
```

Analyze Simulation Data Using Signal Tolerances

You can programmatically specify signal tolerance values to use in comparisons performed using the Simulation Data Inspector. In this example, you compare data collected by simulating a model of an aircraft longitudinal flight control system. Each simulation uses a different value for the input filter time constant and logs the input and output signals. You analyze the effect of the time constant change by comparing results using the Simulation Data Inspector and signal tolerances.

First, load the session file that contains the simulation data.

```
Simulink.sdi.load('AircraftExample.mldatx');
```

The session file contains four runs. In this example, you compare data from the first two runs in the file. Access the `Simulink.sdi.Run` objects for the first two runs loaded from the file.

```
runIDs = Simulink.sdi.getAllRunIDs;
runIDs1 = runIDs(end-3);
runIDs2 = runIDs(end-2);
```

Now, compare the two runs without specifying any tolerances.

```
noTolDiffResult = Simulink.sdi.compareRuns(runIDs1, runIDs2);
```

Use the `getResultByIndex` function to access the comparison results for the `q` and `alpha` signals.

```
qResult = getResultByIndex(noTolDiffResult, 1);
alphaResult = getResultByIndex(noTolDiffResult, 2);
```

Check the `Status` of each signal result to see whether the comparison result fell within our out of tolerance.

```
qResult.Status
```

```
ans =
  ComparisonSignalStatus enumeration
    OutOfTolerance

```

```
alphaResult.Status
```

```
ans =
  ComparisonSignalStatus enumeration

```

OutOfTolerance

The comparison used a value of 0 for all tolerances, so the `OutOfTolerance` result means the signals are not identical.

You can further analyze the effect of the time constant by specifying tolerance values for the signals. Specify the tolerances by setting the properties for the `Simulink.sdi.Signal` objects that correspond to the signals being compared. Comparisons use tolerances specified for the baseline signals. This example specifies a time tolerance and an absolute tolerance.

To specify a tolerance, first access the `Signal` objects from the baseline run.

```
runTs1 = Simulink.sdi.getRun(runIDTs1);
qSig = getSignalsByName(runTs1, 'q, rad/sec');
alphaSig = getSignalsByName(runTs1, 'alpha, rad');
```

Specify an absolute tolerance of 0.1 and a time tolerance of 0.6 for the `q` signal using the `AbsTol` and `TimeTol` properties.

```
qSig.AbsTol = 0.1;
qSig.TimeTol = 0.6;
```

Specify an absolute tolerance of 0.2 and a time tolerance of 0.8 for the `alpha` signal.

```
alphaSig.AbsTol = 0.2;
alphaSig.TimeTol = 0.8;
```

Compare the results again. Access the results from the comparison and check the `Status` property for each signal.

```
tolDiffResult = Simulink.sdi.compareRuns(runIDTs1, runIDTs2);
qResult2 = getResultByIndex(tolDiffResult, 1);
alphaResult2 = getResultByIndex(tolDiffResult, 2);
```

```
qResult2.Status
```

```
ans =
    ComparisonSignalStatus enumeration
        WithinTolerance
```

```
alphaResult2.Status
```

```
ans =
    ComparisonSignalStatus enumeration
        WithinTolerance
```

See Also
Simulation Data Inspector

Related Examples

- [“Compare Simulation Data”](#)
- [“How the Simulation Data Inspector Compares Data”](#)
- [“Create Plots Using the Simulation Data Inspector”](#)

Limit the Size of Logged Data

In this section...

“Limit the Number of Runs Retained in the Simulation Data Inspector Archive” on page 55-48

“Specify a Minimum Disk Space Requirement or Maximum Size for Logged Data” on page 55-48

“View Data Only During Simulation” on page 55-49

“Reduce the Number of Data Points Logged from Simulation” on page 55-49

Logging simulation data can produce large amounts of data that fill up disk space. Such situations include logging many signals, logging data for long simulations, and running many simulations without deleting run data from the Simulation Data Inspector. You can choose among several options to limit the size of logged simulation data. You can:

- Limit the number of runs retained in the Simulation Data Inspector archive.
- Reduce the number of data points logged in each simulation.
- Specify a minimum disk space requirement or maximum size for logged data.
- Configure logging for only viewing data during simulation.

Depending on your requirements, you can use more than one strategy to limit the size of logged data.

Limit the Number of Runs Retained in the Simulation Data Inspector Archive

When you run multiple simulations in a single MATLAB session, logged simulation data accumulates in the Simulation Data Inspector even if you overwrite the logging data in the MATLAB workspace. To reduce the amount of data the Simulation Data Inspector retains, you can configure a limit for the number of runs stored in the archive. When the number of runs in the archive reaches the size limit, the Simulation Data Inspector starts to delete runs from the archive on a first-in, first-out basis.

Configure the archive **Size** setting in the Simulation Data Inspector preferences. The size limit only applies to runs in the archive. For the Simulation Data Inspector to automatically limit data retention, select **Automatically archive** and specify the maximum number of runs to retain in the archive. By default, **Automatically archive** is enabled with an archive size limit of twenty runs. If you experience issues with logged data consuming too much disk space, consider adjusting the size limit for the archive in the Simulation Data Inspector preferences.

Specify a Minimum Disk Space Requirement or Maximum Size for Logged Data

You can use preferences in the Simulation Data Inspector to directly limit the size of logged data by specifying a minimum amount of disk space to leave free or by specifying a maximum size for logged data on disk. Each setting accounts for all kinds of logged data. By default, logged data must leave at least 100 MB of free disk space with no maximum size limit. Specify the required disk space and maximum size in GB, and specify 0 to apply no disk space requirement or no maximum size limit.

When you specify a minimum disk space requirement or a maximum size for logged data, you can also specify whether to prioritize retaining data from the current simulation or data from prior simulations when approaching the limit. By default, the Simulation Data Inspector prioritizes

retaining data for the current run. As the free disk space or logged data size approaches the limit, prior runs are deleted first to free up space for data being logged in the current run. If deleting runs does not free up enough space, recording is disabled. To prioritize retaining prior data, change the **When low on disk space** setting to **Keep prior runs and stop recording**. You see a warning message when prior runs are deleted and when recording is disabled. If recording is disabled due to the size of logged data, you need to change the **Record Mode** back to **View and record data** to continue logging data, after you have freed up disk space.

View Data Only During Simulation

In some situations, you may want to only view the data for logged signals and not save the values. For example, when using the Simulation Data Inspector to visualize data streaming from hardware, you may only want to view the data live and not record it. You can change the **Record mode** in the Simulation Data Inspector preferences to **View during simulation only** so that logged data is not saved and you can still view the data during simulation. The **Record mode** is reset to **View and record data** at the start of each MATLAB session.

When you change the **Record mode** to **View during simulation only**:

- Logged data is not available in the Simulation Data Inspector or workspace after simulation.
- You can view data using dashboard blocks, scopes, and the Simulation Data Inspector, but plots clear when you pan or zoom.
- You cannot access logged data during simulation using the Simulation Data Inspector programmatic interface.

Reduce the Number of Data Points Logged from Simulation

Model configuration parameters and signal properties allow you to limit the number of data points logged in a simulation. Be sure to carefully consider data requirements when limiting logged data points. Limiting data can skip critical time points, and can lead to aliasing, if your effective sample rate is too low.

You can reduce the number of data points using:

- Decimation — Log every n th signal value.
- Limit data points to last — Only log the last n signal values.
- Logging intervals — Specify specific time intervals in which to log data.

For details, see “Specify Signal Values to Log”.

See Also

Tools

Simulation Data Inspector

Related Examples

- “Specify Signal Values to Log”
- “Configure the Simulation Data Inspector”

