

Fixed-Point Toolbox

For Use with **MATLAB**®

- Computation
- Visualization
- Programming

User's Guide

Version 1



How to Contact The MathWorks:



www.mathworks.com Web
comp.soft-sys.matlab Newsgroup



support@mathworks.com Technical Support
suggest@mathworks.com Product enhancement suggestions
bugs@mathworks.com Bug reports
doc@mathworks.com Documentation error reports
service@mathworks.com Order status, license renewals, passcodes
info@mathworks.com Sales, pricing, and general information



508-647-7000 Phone



508-647-7001 Fax



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

For contact information about worldwide offices, see the MathWorks Web site.

Fixed-Point Toolbox User's Guide

© COPYRIGHT 2004–2005 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.

MATLAB, Simulink, Stateflow, Handle Graphics, Real-Time Workshop, and xPC TargetBox are registered trademarks of The MathWorks, Inc.

Other product or brand names are trademarks or registered trademarks of their respective holders.

Revision History:

June 2004	First printing	New for Version 1.0 (Release 14)
October 2004	Online only	Version 1.1 (Release 14SP1)
March 2005	Online only	Version 1.2 (Release 14SP2)

Getting Started

1

What Is the Fixed-Point Toolbox?	1-2
Features	1-2
Getting Help	1-3
Getting Help in This Document	1-3
Getting Help at the MATLAB Command Line	1-3
Display Settings	1-5
Demos	1-7

Fixed-Point Concepts

2

Fixed-Point Data Types	2-2
Scaling	2-4
Precision and Range	2-5
Range	2-5
Precision	2-6
Arithmetic Operations	2-8
Modulo Arithmetic	2-8
Two's Complement	2-9
Addition and Subtraction	2-10
Multiplication	2-11
Casts	2-16

fi Objects Compared to C Integer Data Types	2-20
Integer Data Types	2-20
Unary Conversions	2-22
Binary Conversions	2-23
Overflow Handling	2-25

Working with fi Objects

3

Constructing fi Objects	3-2
Examples of Constructing fi Objects	3-3
fi Object Properties	3-10
Data Properties	3-10
fimath Properties	3-10
numericType Properties	3-11
Setting Fixed-Point Properties at Object Creation	3-12
Using Direct Property Referencing with fi	3-12
fi Object Functions	3-14

Working with fimath Objects

4

Constructing fimath Objects	4-2
fimath Object Properties	4-4
Setting fimath Properties at Object Creation	4-4
Using Direct Property Referencing with fimath	4-5
Using fimath Objects to Perform Fixed-Point Arithmetic	4-6
Using fimath to Share Arithmetic Rules	4-8

Using fimath ProductMode and SumMode	4-10
FullPrecision	4-10
KeepLSB	4-11
KeepMSB	4-12
SpecifyPrecision	4-13
fimath Object Functions	4-15

Working with fipref Objects

5

Constructing fipref Objects	5-2
fipref Object Properties	5-3
Setting fipref Properties at Object Creation	5-3
Using Direct Property Referencing with fipref	5-3
Using fipref Objects to Set Display Preferences	5-5
Using fipref Objects to Set Logging Preferences	5-7
fipref Object Functions	5-10

Working with numerictype Objects

6

Constructing numerictype Objects	6-2
Examples of Constructing numerictype Objects	6-3
numerictype Object Properties	6-6
Setting numerictype Properties at Object Creation	6-6
Using Direct Property Referencing with numerictype Objects	6-7
Setting numerictype Properties in the Model Explorer ...	6-7

The numerictype Structure	6-10
Properties That Affect the Slope	6-11
Stored Integer Value and Real World Value	6-11
Using numerictype Objects to Share Data Type and Scaling Settings	6-12
numerictype Object Functions	6-15

Working with quantizer Objects

7

Constructing quantizer Objects	7-2
quantizer Object Properties	7-4
Settable quantizer Object Properties	7-4
Read-Only quantizer Object Properties	7-5
Quantizing Data with quantizer Objects	7-6
Transformations for Quantized Data	7-8
quantizer Object Functions	7-9

Interoperability with Other Products

8

Using fi Objects with Simulink	8-2
Reading Fixed-Point Data from the Workspace	8-2
Writing Fixed-Point Data to the Workspace	8-2
Logging Fixed-Point Signals	8-6
Accessing Fixed-Point Block Data During Simulation	8-6
Using fi Objects with Signal Processing Blockset	8-7

Reading Fixed-Point Signals from the Workspace	8-7
Writing Fixed-Point Signals to the Workspace	8-7
Using fi Objects with Filter Design Toolbox	8-12

Property Reference

9

fi Object Properties	9-2
bin	9-2
data	9-2
dec	9-2
double	9-2
fimath	9-2
hex	9-3
int	9-3
NumericType	9-3
oct	9-4
fimath Object Properties	9-5
CastBeforeSum	9-5
MaxProductWordLength	9-5
MaxSumWordLength	9-5
OverflowMode	9-5
ProductFractionLength	9-6
ProductMode	9-6
ProductWordLength	9-7
RoundMode	9-7
SumFractionLength	9-8
SumMode	9-8
SumWordLength	9-9
fipref Object Properties	9-10
FimathDisplay	9-10
LoggingMode	9-10
NumericTypeDisplay	9-10
NumberDisplay	9-11
numerictype Object Properties	9-12

Bias	9-12
DataType	9-12
DataTypeMode	9-12
FixedExponent	9-13
FractionLength	9-13
Scaling	9-14
Signed	9-14
Slope	9-14
SlopeAdjustmentFactor	9-14
WordLength	9-15
quantizer Object Properties	9-16
DataMode	9-16
Format	9-16
Max	9-17
Min	9-17
NOperations	9-18
NOverflows	9-18
NUnderflows	9-18
OverflowMode	9-18
RoundMode	9-19

Functions — Categorical List

10

Bitwise Functions	10-2
Constructor and Property Functions	10-2
Data Manipulation Functions	10-3
Data Type Functions	10-5
Data Quantizing Functions	10-6
Element-Wise Logical Operator Functions	10-6

Math Operation Functions	10-6
Matrix Manipulation Functions	10-8
Plotting Functions	10-9
Radix Conversion Functions	10-12
Relational Operator Functions	10-13
Statistics Functions	10-14
Subscripted Assignment and Reference Functions	10-15
fi Object Functions	10-16
fix Object Functions	10-18
fixpref Object Functions	10-19
numeric Object Functions	10-20
quantizer Object Functions	10-21

Functions — Alphabetical List

11

Glossary

Index

Getting Started

“What Is the Fixed-Point Toolbox?” (p. 1-2)	Describes the Fixed-Point Toolbox and its major features
“Getting Help” (p. 1-3)	Tells you how to get help on Fixed-Point Toolbox objects, properties, and functions
“Display Settings” (p. 1-5)	Describes the <code>fi</code> object display settings used in the code examples in this User’s Guide
“Demos” (p. 1-7)	Lists the Fixed-Point Toolbox demos

What Is the Fixed-Point Toolbox?

The Fixed-Point Toolbox provides fixed-point data types in MATLAB® and enables algorithm development by providing fixed-point arithmetic. The Fixed-Point Toolbox enables you to create the following types of objects:

- `fi` – Defines a fixed-point numeric object in the MATLAB workspace. Each `fi` object is composed of value data, a `fimath` object, and a `numericity` object.
- `fimath` – Governs how overloaded arithmetic operators work with `fi` objects
- `fioref` – Defines the display and logging preferences of `fi` objects
- `numericity` – Defines the data type and scaling attributes of `fi` objects
- `quantizer` – Quantizes data sets

Features

The Fixed-Point Toolbox provides you with

- The ability to define fixed-point data types, scaling, and rounding and overflow methods in the MATLAB workspace
- Bit-true real and complex simulation
- Basic fixed-point arithmetic with binary point-only signals
 - Arithmetic operators `+`, `-`, `*`, `.*`
 - Division using the `divide` function
- Arbitrary word length up to `intmax('uint16')`
- Overflow and underflow logging
- Relational, logical, and bitwise operators
- Statistics functions such as `max` and `min`
- Conversions between binary, hex, double, and built-in integers
- Interoperability with Simulink®, Signal Processing Blockset, and Filter Design Toolbox
- Compatibility with the Simulink To Workspace and From Workspace blocks

Getting Help

This section tells you how to get help for the Fixed-Point Toolbox in this document and at the MATLAB command line.

Getting Help in This Document

The objects of the Fixed-Point Toolbox are discussed in the following chapters:

- Chapter 3, “Working with fi Objects”
- Chapter 4, “Working with fimath Objects”
- Chapter 5, “Working with fipref Objects”
- Chapter 6, “Working with numericitype Objects”
- Chapter 7, “Working with quantizer Objects”

To get in-depth information about the properties of these objects, refer to the Property Reference in the online documentation.

To get in-depth information about the functions of these objects, refer to the Function Reference in the online documentation.

Getting Help at the MATLAB Command Line

To get command-line help for Fixed-Point Toolbox objects, type

```
help objectname
```

For example,

```
help fi
help fimath
help fipref
help numericitype
help quantizer
```

To invoke Help Browser documentation for Fixed-Point Toolbox functions from the MATLAB command line, type

```
doc fixedpoint/functionname
```

For example,

```
doc fixedpoint/int
doc fixedpoint/add
doc fixedpoint/savefipref
doc fixedpoint/quantize
```

Display Settings

In the Fixed-Point Toolbox, the display of `fi` objects is determined by the `fipref` object. Throughout this User's Guide, code examples of `fi` objects are usually shown as they appear when the `fipref` properties are set as follows:

- `NumberDisplay` – 'RealWorldValue'
- `NumericTypeDisplay` – 'full'
- `FimathDisplay` – 'none'

For example,

```
p = fipref('NumberDisplay', 'RealWorldValue',...
         'NumericTypeDisplay', 'full', 'FimathDisplay', 'none')

p =

        NumberDisplay: 'RealWorldValue'
    NumericTypeDisplay: 'full'
        FimathDisplay: 'none'
        LoggingMode: 'Off'

a = fi(pi)

a =

        3.1416

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

In other cases, it makes sense to also show the `fimath` object display:

- `NumberDisplay` – 'RealWorldValue'
- `NumericTypeDisplay` – 'full'

- `FimathDisplay` – 'full'

For example,

```
p = fipref('NumberDisplay', 'RealWorldValue',...  
         'NumericTypeDisplay', 'full', 'FimathDisplay', 'full')
```

```
p =
```

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

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

```
a =
```

```
3.1416
```

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

```
         RoundMode: round  
         OverflowMode: saturate  
         ProductMode: FullPrecision  
MaxProductWordLength: 128  
         SumMode: FullPrecision  
MaxSumWordLength: 128  
         CastBeforeSum: true
```

For more information, refer to Chapter 5, “Working with fipref Objects”

Demos

You can access demos in the **Demos** tab of the **Help Navigator** window. The Fixed-Point Toolbox includes the following demos:

- **fi Basics** – Demonstrates the basic use of the fixed-point object `fi`
- **fi Binary Point Scaling** – Explains binary point-only scaling
- **Fixed-Point Algorithm Development** – Shows the development and verification of a simple fixed-point algorithm
- **Fixed-Point C Development** – Shows how to use the parameters from a fixed-point MATLAB program in a fixed-point C program
- **Number Circle** – Illustrates the definitions of unsigned and signed two's complement integer and fixed-point numbers
- **Quantization Error** – Demonstrates the statistics of the error when signals are quantized using various rounding methods
- **Analysis of a Fixed-Point State-Space System with Limit Cycles** – Demonstrates a limit cycle detection routine applied to a state-space system

Fixed-Point Concepts

“Fixed-Point Data Types” (p. 2-2)	Defines fixed-point data types
“Scaling” (p. 2-4)	Discusses the types of scaling used in the Fixed-Point Toolbox; binary point-only and [Slope Bias]
“Precision and Range” (p. 2-5)	Discusses the concepts behind arithmetic operations in the Fixed-Point Toolbox.
“Arithmetic Operations” (p. 2-8)	Introduces the concepts behind arithmetic operations in the Fixed-Point Toolbox
“fi Objects Compared to C Integer Data Types” (p. 2-20)	Compares ANSI C integer data type ranges, conversions, and exception handling with those of fi objects

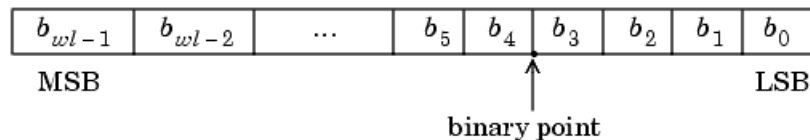
Fixed-Point Data Types

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. This chapter discusses many terms and concepts relating to fixed-point numbers, data types, and mathematics.

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, therefore, 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. Signed binary fixed-point numbers are typically represented in one of three ways:

- Sign/magnitude
- One's complement
- Two's complement

Two's complement is the most common representation of signed fixed-point numbers and is the only representation used by the Fixed-Point Toolbox. Refer to "Two's Complement" on page 2-9 for more information.

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

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

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 the Fixed-Point Toolbox, 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 fractional slope 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

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

The Fixed-Point Toolbox supports both binary point-only scaling and [Slope Bias] scaling.

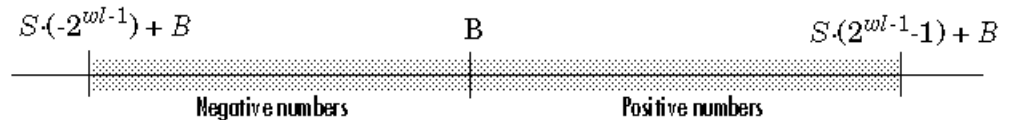
Note For examples of binary point-only scaling, see the Fixed-Point Toolbox demo "fi Binary Point Scaling."

Precision and Range

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

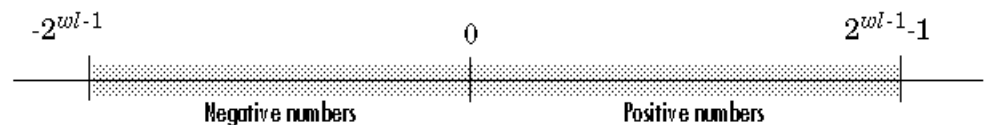
The range is the span of numbers that a fixed-point data type and scaling can represent. The range of representable numbers for a two's complement fixed-point number of word length wl , scaling S , and bias B is illustrated below:



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.

The Fixed-Point Toolbox 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. Refer to “Modulo Arithmetic” on page 2-8 for more information.

When you create a `fi` object in the Fixed-Point Toolbox, any overflows are saturated. The `OverflowMode` property of the default `fimath` object is `saturate`. You can log overflows and underflows as warnings by setting the `LoggingMode` property of the `fipref` object to `'overflowandunderflow'`. Refer to “LoggingMode” on page 9-10 for more information.

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

One of the limitations of representing numbers with finite precision is that not every number in the available range can be represented exactly. When the result of a fixed-point calculation is a number that cannot be represented exactly by the data type and scaling being used, precision is lost. A rounding method must be used to cast the result to a representable number. The Fixed-Point Toolbox currently supports the following rounding methods:

- `floor`, which is equivalent to truncation, rounds to the closest representable number in the direction of negative infinity.
- `ceil` rounds to the closest representable number in the direction of positive infinity.

- `fix` rounds to the closest representable integer in the direction of zero.
- `convergent` rounds to the closest representable integer. In the case of a tie, it rounds to the nearest even integer.
- `round` rounds to the closest representable integer. In the case of a tie, it rounds to the closest representable integer in the direction of positive infinity. This is the default rounding method for `fi` object creation and `fi` arithmetic.

Arithmetic Operations

The following sections describe the arithmetic operations used by the Fixed-Point Toolbox:

- “Modulo Arithmetic” on page 2-8
- “Two’s Complement” on page 2-9
- “Addition and Subtraction” on page 2-10
- “Multiplication” on page 2-11
- “Casts” on page 2-16

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:

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 1 using binary math.
- 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:

$$11010 \longrightarrow 00101$$

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):

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

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):

$$\begin{array}{r}
 010010.100 \text{ (18.5)} \\
 - 0110.110 \text{ (6.75)} \\
 \hline
 \end{array}
 \xrightarrow{\substack{\text{two's complement} \\ \text{and sign extension}}}
 \begin{array}{r}
 010010.100 \text{ (18.5)} \\
 +111001.010 \text{ (-6.75)} \\
 \hline
 1001011.110 \text{ (11.75)}
 \end{array}$$

Carry bit is discarded.

The default `fimath` object 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):

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

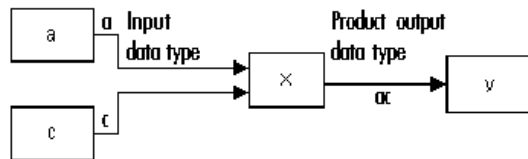
The extra 1 is the result of necessary sign extension.

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

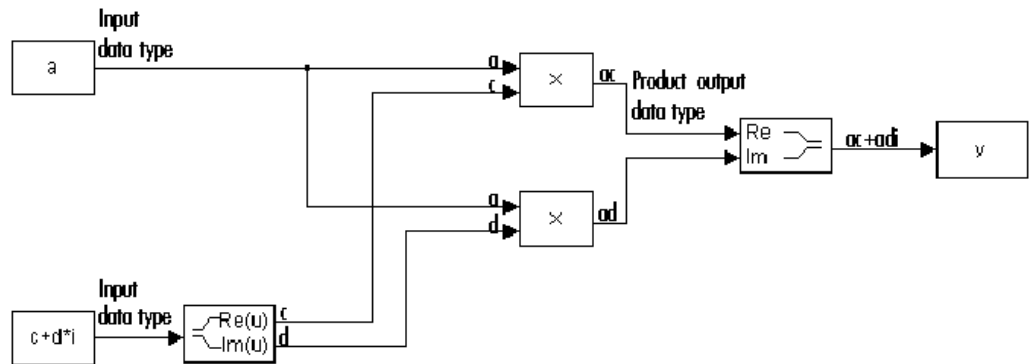
Multiplication Data Types

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

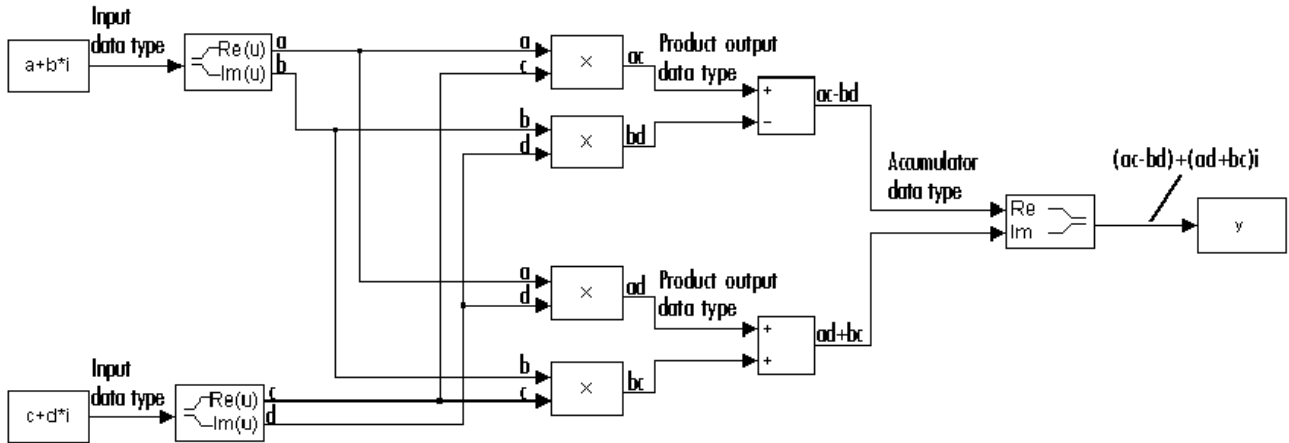
Real-Real Multiplication. The following diagram shows the data types used in the multiplication of two real numbers in the Fixed-Point Toolbox. The output of this multiplication is in the product data type, which is governed by the `fimath ProductMode` property:



Real-Complex Multiplication. The following diagram shows the data types used in the multiplication of a real and a complex fixed-point number in the Fixed-Point Toolbox. Real-complex and complex-real multiplication are equivalent. The output of this multiplication is in the product data type, which is governed by the `fimath ProductMode` property:



Complex-Complex Multiplication. The following diagram shows the multiplication of two complex fixed-point numbers in the Fixed-Point Toolbox. Note that the output of the multiplication is in the sum data type, which is governed by the `fimath SumMode` property. The product data type is determined by the `fimath ProductMode` property:



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

```
Signed: true
WordLength: 24
FractionLength: 20
```

```
y = fi(10, T2, F)
```

```
y =
```

```
10
```

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

```
z = x*y
```

```
z =
```

```
50
```

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

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

```
y =
```

```
10.0000 + 2.0000i
```

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

```
z = x*y
```

```
z =
```

```
50.0000 +10.0000i
```

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

Complex*Complex. Complex-complex multiplication involves an addition as well as multiplication, so 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  
        Signed: true  
        WordLength: 24  
    FractionLength: 20
```

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

```
y =
```

```
10.0000 + 2.0000i
```

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

```
z = x*y
```

```
z =
```

```
38.0000 +70.0000i
```

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

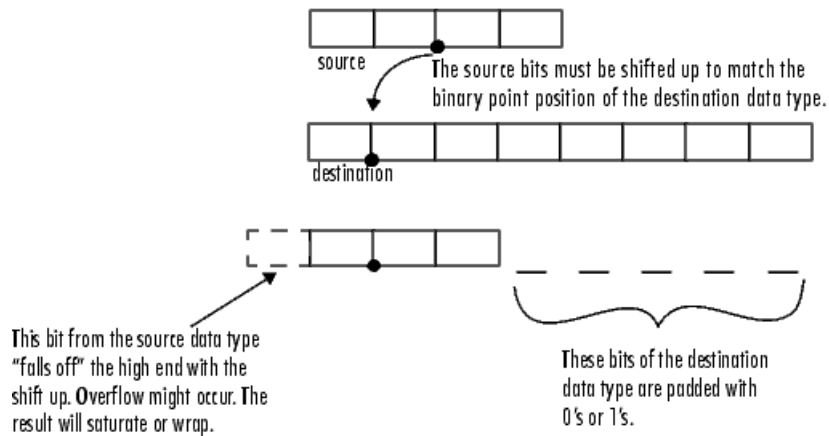
Casts

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.

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:

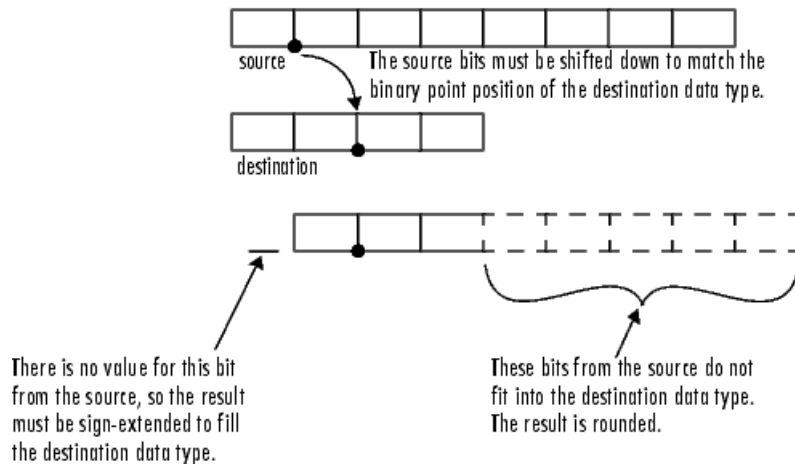
- If overflow does not occur, the empty bits are padded with 0's.
- If wrapping occurs, the empty bits are padded with 0's.
- If saturation occurs,
 - The empty bits of a positive number are padded with 1's.
 - The empty bits of a negative number are padded with 0's.

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 the result is sign-extended to fill the integer portion of the destination data type. 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 Compared to C Integer Data Types

The following sections compare the `fi` object with fixed-point data types and operations in C:

- “Integer Data Types” on page 2-20
- “Unary Conversions” on page 2-22
- “Binary Conversions” on page 2-23
- “Overflow Handling” on page 2-25

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 ranges of ANSI C integer data types.

ANSI C Integer Data Types

The following table shows the minimum ranges of ANSI C integer data types. The integer ranges can be larger than or equal to those 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`.

Note that the minimum ANSI C ranges are large enough to accommodate one’s complement or sign/magnitude representation, but not two’s complement representation. In the one’s complement and sign/magnitude representations, a signed integer with n bits has a range from $-2^{n-1} + 1$ to $2^{n-1} - 1$, inclusive. In both of these representations, an equal number of positive and negative numbers are represented, and zero is represented twice.

Integer Type	Minimum	Maximum
signed char	-127	127
unsigned char	0	255

Integer Type	Minimum	Maximum
short int	-32,767	32,767
unsigned short	0	65,535
int	-32,767	32,767
unsigned int	0	65,535
long int	-2,147,483,647	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 the Fixed-Point Toolbox. 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.

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}	$2^{n-1} - 1$	N/A
<code>fi(x,0,n,0)</code>	No	n (2 to 65,535)	0	0	$2^n - 1$	N/A
<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

Constructor	Signed	Word Length	Fraction Length	Minimum	Maximum	Closest ANSI C Equivalent
<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
!x	$\sim x = \text{not}(x)$	Result is logical.
~x	<code>bitcmp(x)</code>	Result is same numeric type as operand.
*x	No equivalent	N/A
x<<n	<code>bitshift(x,n)</code> positive n	Result is same numeric type as operand. Overflow mode is obeyed: wrap or saturate if 1-valued bits are shifted off the left, or into the sign bit if the operand is signed. 0-valued bits are shifted in on the right.
x>>n	<code>bitshift(x,-n)</code>	Result is same numeric type as operand. Round mode is obeyed if 1-valued bits are shifted off the right. 0-valued bits are shifted in on the left if the operand is either signed and positive or unsigned. 1-valued bits are shifted in on the left if the operand is signed and negative.
+x	+x	Result is same numeric type as operand.
-x	-x	Result is same numeric type as operand. Overflow mode is obeyed. For example, overflow might occur when you negate an unsigned fi or the most negative value of a signed fi.

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

Type of One Operand	Type of Other Operand	Properties of Other Operand After Conversion to a fi Object
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

Overflow Handling

The following sections compare how overflows are handled in ANSI C and the Fixed-Point Toolbox.

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 `fi`math object. You can specify certain overflow characteristics separately for sums (including differences) and products. Refer to the following table:

fi math Object Properties Related to Overflow Handling	Property Value	Description
OverflowMode	'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 "ProductMode" on page 9-6 in the online documentation.

fimath Object Properties Related to Overflow Handling	Property Value	Description
	'KeepLSB'	<p>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.</p> <p>The resulting word length is determined by the ProductWordLength property. If ProductWordLength is greater than is necessary for the full-precision product, then the result is stored in the least significant bits. If ProductWordLength is less than is necessary for the full-precision product, then overflow occurs.</p> <p>The rule for computing the resulting product fraction length is given in “ProductMode” on page 9-6 in the online documentation.</p>
	'KeepMSB'	<p>The most significant bits of the product are kept. Overflow is prevented, but precision may be lost.</p> <p>The resulting word length is determined by the ProductWordLength property. 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 “ProductMode” on page 9-6 in the online documentation.</p>
	'SpecifyPrecision'	<p>You can specify both the word length and the fraction length of the resulting product.</p>

fimath Object Properties Related to Overflow Handling	Property Value	Description
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 128 bits. The maximum 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 “SumMode” on page 9-8 in the online documentation.</p>
	'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 resulting word length is determined by the SumWordLength property. If SumWordLength is greater than is necessary for the full-precision sum, then the result is stored in the least significant bits. If SumWordLength is less than is necessary for the full-precision sum, then overflow occurs.</p>

fimath Object Properties Related to Overflow Handling	Property Value	Description
		The rule for computing the resulting sum fraction length is given in “SumMode” on page 9-8 in the online documentation.
	'KeepMSB'	<p>The most significant bits of the sum are kept. Overflow is prevented, but precision may be lost.</p> <p>The resulting word length is determined by the SumWordLength property. If SumWordLength is greater than is necessary for the full-precision sum, then the result is stored in the most significant bits. If SumWordLength 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 “SumMode” on page 9-8 in the online documentation.</p>
	'SpecifyPrecision'	You can specify both the word length and the fraction length of the resulting sum.
SumWordLength	Positive integer	The word length of sum results when SumMode is 'KeepLSB', 'KeepMSB', or 'SpecifyPrecision'.
MaxSumWordLength	Positive integer	The maximum sum word length allowed when SumMode is 'FullPrecision'. The default is 128 bits. The maximum is 65,535 bits. This property can help ensure that your simulation does not exceed your hardware requirements.
SumFractionLength	Integer	The fraction length of sum results when SumMode is 'SpecifyPrecision'.

Working with `fi` Objects

“Constructing `fi` Objects” (p. 3-2)

Teaches you how to create `fi` objects

“`fi` Object Properties” (p. 3-10)

Tells you how to find more information about the properties associated with `fi` objects, and shows you how to set these properties

“`fi` Object Functions” (p. 3-14)

Introduces the functions in the toolbox that operate directly on `fi` objects

Constructing `fi` Objects

You can create `fi` objects in the Fixed-Point Toolbox in one of two ways:

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

To get started, type

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

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

```
a =
```

```
0
```

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

A signed `fi` object is created with a value of 0, word length of 16 bits, and fraction length of 15 bits.

Note For information on the display format of `fi` objects, refer to “Display Settings” on page 1-5.

You can use the `fi` constructor function in the following ways:

- `fi(v)` returns a signed fixed-point object with value `v`, 16-bit word length, and best-precision fraction length.
- `fi(v,s)` returns a fixed-point object with value `v`, signedness `s`, 16-bit word length, and best-precision fraction length. `s` can be 0 (false) for unsigned or 1 (true) for signed.

- `fi(v,s,w)` returns a fixed-point object with value `v`, signedness `s`, word length `w`, and best-precision fraction length.
- `fi(v,s,w,f)` returns a fixed-point object with value `v`, signedness `s`, word length `w`, and fraction length `f`.
- `fi(v,s,w,slope,bias)` returns a fixed-point object with value `v`, signedness `s`, word length `w`, slope, and bias.
- `fi(v,s,w,slopeadjustmentfactor,fixedexponent,bias)` returns a fixed-point object with value `v`, signedness `s`, word length `w`, slope adjustment `slopeadjustmentfactor`, exponent `fixedexponent`, and bias `bias`.
- `fi(v,T)` returns a fixed-point object with value `v` and embedded.numericity `T`. Refer to Chapter 6, “Working with numericity Objects” for more information on numericity objects.
- `fi(v,T,F)` returns a fixed-point object with value `v`, embedded.numericity `T`, and embedded.fimath `F`. Refer to Chapter 4, “Working with fimath Objects” for more information on fimath objects.
- `fi(...'PropertyName',PropertyValue...)` and `fi('PropertyName',PropertyValue...)` allow you to set properties for a `fi` object using property name/property value pairs.

Examples of Constructing fi Objects

For example, the following creates a `fi` object with a value of `pi`, a word length of 8 bits, and a fraction length of 3 bits.

```
a = fi(pi, 1, 8, 3)
```

```
a =
```

```
3.1250
```

```

      DataTypeMode: Fixed-point: binary point scaling
      Signed: true
      WordLength: 8
      FractionLength: 3

```

The value `v` can also be an array.

```
a = fi((magic(3)/10), 1, 16, 12)
```

```
a =
```

```
    0.8000    0.1001    0.6001  
    0.3000    0.5000    0.7000  
    0.3999    0.8999    0.2000
```

```
        DataTypeMode: Fixed-point: binary point scaling  
                Signed: true  
        WordLength: 16  
        FractionLength: 12
```

If you omit the argument `f`, it is set automatically to the best precision possible.

```
a = fi(pi, 1, 8)
```

```
a =
```

```
    3.1563
```

```
        DataTypeMode: Fixed-point: binary point scaling  
                Signed: true  
        WordLength: 8  
        FractionLength: 5
```

If you omit `w` and `f`, they are set automatically to 16 bits and the best precision possible, respectively.

```
a = fi(pi, 1)
```

```
a =
```

```
    3.1416
```

```
        DataTypeMode: Fixed-point: binary point scaling
```

```
Signed: true
WordLength: 16
FractionLength: 13
```

Constructing a fi Object with Property Name/Property Value Pairs

You can use property name/property value pairs to set fi properties when you create the object:

```
a = fi(pi, 'roundmode', 'floor', 'overflowmode', 'wrap')
a =
    3.1415
```

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

Constructing a fi Object Using a numerictype Object

You can use a numerictype object to define a fi object:

```
T = numerictype

T =

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

a = fi(pi, T)

a =
```

```
1.0000
```

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

```
      RoundMode: round  
    OverflowMode: saturate  
      ProductMode: FullPrecision  
MaxProductWordLength: 128  
      SumMode: FullPrecision  
MaxSumWordLength: 128  
    CastBeforeSum: true
```

You can also use a `fimath` object with a numeric type object to define a `fi` object:

```
F = fimath
```

```
F =
```

```
      RoundMode: round  
    OverflowMode: saturate  
      ProductMode: FullPrecision  
MaxProductWordLength: 128  
      SumMode: FullPrecision  
MaxSumWordLength: 128  
    CastBeforeSum: true
```

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

```
a =
```

```
1.0000
```

```

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

            RoundMode: round
            OverflowMode: saturate
            ProductMode: FullPrecision
    MaxProductWordLength: 128
            SumMode: FullPrecision
    MaxSumWordLength: 128
        CastBeforeSum: true

```

Determining Property Precedence

Note that the value of a property is taken from the last time it is set. For example, create a `numerictype` object with a value of `true` for the `'signed'` property:

```
T = numerictype('signed', true)
```

```
T =
```

```

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

```

Now create the following `fi` object in which the `numerictype` property is specified *after* the `signed` property, so that the resulting `fi` object is signed:

```
a = fi(pi, 'signed', false, 'numerictype', T)
```

```
a =
```

```
1.0000
```

```

        DataTypeMode: Fixed-point: binary point scaling
            Signed: true

```

```
        WordLength: 16
    FractionLength: 15

        RoundMode: round
    OverflowMode: saturate
    ProductMode: FullPrecision
MaxProductWordLength: 128
        SumMode: FullPrecision
    MaxSumWordLength: 128
    CastBeforeSum: true
```

Contrast this with the following `fi` object in which the `numericity` property is specified *before* the `signed` property, so the resulting `fi` object is unsigned:

```
b = fi(pi,'numericity',T,'signed',false)
```

```
b =
```

```
2.0000
```

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

```
        RoundMode: round
    OverflowMode: saturate
    ProductMode: FullPrecision
MaxProductWordLength: 128
        SumMode: FullPrecision
    MaxSumWordLength: 128
    CastBeforeSum: true
```

Copying a fi Object

To copy a `fi` object, use the `fi` constructor function:

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

```
a =
```


3.1416

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

b = fi(a)

b =

3.1416

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

fi Object Properties

The `fi` object has the following three general types of properties:

- “Data Properties” on page 3-10
- “`fimath` Properties” on page 3-10
- “`numerictype` Properties” on page 3-11

Data Properties

The data properties of a `fi` object are always writable:

- `bin` — Stored integer value of a `fi` object in binary
- `data` — Numerical real-world value of a `fi` object
- `dec` — Stored integer value of a `fi` object in decimal
- `double` — Real-world value of a `fi` object, stored as a MATLAB double
- `hex` — Stored integer value of a `fi` object in hexadecimal
- `int` — Stored integer value of a `fi` object, stored in a built-in MATLAB integer data type. You can also use `int8`, `int16`, `int32`, `uint8`, `uint16`, and `uint32` to get the stored integer value of a `fi` object in these formats
- `oct` — Stored integer value of a `fi` object in octal

fimath Properties

When you create a `fi` object, a `fimath` object is also automatically created as a property of the `fi` object:

- `fimath` — `fimath` object associated with a `fi` object

The following `fimath` properties are, by transitivity, also properties of a `fi` object. The properties of the `fimath` object listed below are always writable:

- `CastBeforeSum` — Whether both operands are cast to the sum data type before addition
- `MaxProductWordLength` — Maximum allowable word length for the product data type

- `MaxSumWordLength` — Maximum allowable word length for the sum data type
- `ProductFractionLength` — Fraction length, in bits, of the product data type
- `ProductMode` — Defines how the product data type is determined
- `ProductWordLength` — Word length, in bits, of the product data type
- `RoundMode` — Rounding mode
- `SumFractionLength` — Fraction length, in bits, of the sum data type
- `SumMode` — Defines how the sum data type is determined
- `SumWordLength` — The word length, in bits, of the sum data type

numericType Properties

When you create a `fi` object, a `numericType` object is also automatically created as a property of the `fi` object:

- `numericType` — Object containing all the numeric type attributes of a `fi` object

The following `numericType` properties are, by transitivity, also properties of a `fi` object. The properties of the `numericType` object listed below are not writable once the `fi` object has been created. However, you can create a copy of a `fi` object with new values specified for the `numericType` properties:

- `Bias` — Bias of a `fi` object
- `DataType` — Data type category associated with a `fi` object
- `DataTypeMode` — Data type and scaling mode of a `fi` object
- `FixedExponent` — Fixed-point exponent associated with a `fi` object
- `SlopeAdjustmentFactor` — Slope adjustment associated with a `fi` object
- `FractionLength` — Fraction length of the stored integer value of a `fi` object in bits
- `Scaling` — Fixed-point scaling mode of a `fi` object
- `Signed` — Whether a `fi` object is signed or unsigned
- `Slope` — Slope associated with a `fi` object

- `WordLength` — Word length of the stored integer value of a `fi` object in bits

These properties are described in detail in Chapter 9, “Property Reference” in the online documentation. There are two ways to specify properties for `fi` objects in the Fixed-Point Toolbox. Refer to the following sections:

- “Setting Fixed-Point Properties at Object Creation” on page 3-12
- “Using Direct Property Referencing with `fi`” on page 3-12

Setting 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 mode to wrap and the rounding mode to convergent,

```
a = fi(pi, 'OverflowMode', 'wrap', 'RoundMode', 'convergent')
```

```
a =
```

```
3.1416
```

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

```
      RoundMode: convergent  
      OverflowMode: wrap  
      ProductMode: FullPrecision  
MaxProductWordLength: 128  
      SumMode: FullPrecision  
MaxSumWordLength: 128  
      CastBeforeSum: true
```

Using Direct Property Referencing with `fi`

You can reference directly into a property for setting or retrieving `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 `DataTypeMode` of `a`,

```
a.DataTypeMode
```

```
ans =
```

```
Fixed-point: binary point scaling
```

To set the `OverflowMode` of `a`,

```
a.OverflowMode = 'wrap'
```

```
a =
```

```
3.1416
```

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

```
          RoundMode: convergent  
        OverflowMode: wrap  
          ProductMode: FullPrecision  
MaxProductWordLength: 128  
              SumMode: FullPrecision  
MaxSumWordLength: 128  
        CastBeforeSum: true
```

fi Object Functions

The functions in the following table operate directly on fi objects.

abs	all	and	any	area
bar	barh	bin	bitand	bitcmp
bitget	bitor	bitshift	bitxor	buffer
clabel	comet	comet3	compass	complex
coneplot	conj	contour	contour3	contourc
contourf	ctranspose	dec	diag	double
end	eps	eq	errorbar	etreeplot
ezcontour	ezcontourf	ezmesh	ezplot	ezplot3
ezpolar	ezsurf	ezsurfz	feather	fi
find	fplot	ge	get	gplot
gt	hankel	hex	hist	histc
horzcat	innerprodintbits	inspect	int	int8
int16	int32	intmax	intmin	ipermute
iscolumn	isequal	isfi	isnumeric	isobject
ispropequal	isrow	assigned	le	line
logical	lowerbound	lsb	lt	max
mesh	meshc	meshz	min	minus
mtimes	ne	not	numberofelements	oct
or	patch	pcolor	permute	plot
plot3	plotmatrix	plotyy	plus	polar
pow2	quiver	quiver3	range	realmax
realmin	rescale	rgbplot	ribbon	rose
scatter	scatter3	sdec	sign	single
slice	spy	stairs	stem	stem3
streamribbon	streamslice	streamtube	stripscaling	subsasgn

sum	surf	surfc	surf1	surfnorm
text	times	toeplitz	treemap	tril
trimesh	triplot	trisurf	triu	uint8
uint16	uint32	uminus	uplus	upperbound
vertcat	voronoi	voronoin	waterfall	xlim
ylim	zlim			

You can learn about the functions associated with `fi` objects in the Function Reference in the online documentation.

The following data-access functions can be also used to get the data in a `fi` object using dot notation.

- `bin`
- `data`
- `dec`
- `double`
- `hex`
- `int`
- `oct`

For example,

```
a = fi(pi);
n = int(a)

n =

    25736

a.int

ans =
```

```
25736
h = hex(a)

h =

6488

a.hex

ans =

6488
```


Working with fimath Objects

“Constructing fimath Objects” (p. 4-2)	Teaches you how to create fimath objects
“fimath Object Properties” (p. 4-4)	Tells you how to find more information about the properties associated with fimath objects, and shows you how to set these properties
“Using fimath Objects to Perform Fixed-Point Arithmetic” (p. 4-6)	Gives examples of using fimath objects to control the results of fixed-point arithmetic with fi objects
“Using fimath to Share Arithmetic Rules” (p. 4-8)	Gives an example of using a fimath object to share modular arithmetic information among multiple fi objects
“Using fimath ProductMode and SumMode” (p. 4-10)	Shows the differences among the different settings of the ProductMode and SumMode properties
“fimath Object Functions” (p. 4-15)	Introduces the functions in the toolbox that operate directly on fimath objects

Constructing `fimath` Objects

`fimath` objects define the arithmetic attributes of `fi` objects. You can create `fimath` objects in the Fixed-Point Toolbox in one of two ways:

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

To get started, type

```
F = fimath
```

to create a default `fimath` object.

```
F = fimath
```

```
F =
```

```
          RoundMode: round
          OverflowMode: saturate
          ProductMode: FullPrecision
MaxProductWordLength: 128
          SumMode: FullPrecision
MaxSumWordLength: 128
          CastBeforeSum: true
```

To copy a `fimath` object, use the `fimath` constructor function:

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

```
ans =
```

```
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 4-4.

fimath Object Properties

All the properties of `fimath` objects are writable:

- `CastBeforeSum` – Whether both operands are cast to the sum data type before addition
- `MaxProductWordLength` – Maximum allowable word length for the product data type
- `MaxSumWordLength` – Maximum allowable word length for the sum data type
- `OverflowMode` – Overflow-handling mode
- `ProductFractionLength` – Fraction length, in bits, of the product data type
- `ProductMode` – Defines how the product data type is determined
- `ProductWordLength` – Word length, in bits, of the product data type
- `RoundMode` – Rounding mode
- `SumFractionLength` – Fraction length, in bits, of the sum data type
- `SumMode` – Defines how the sum data type is determined
- `SumWordLength` – Word length, in bits, of the sum data type

These properties are described in detail in the Chapter 9, “Property Reference” in the online documentation. There are two ways to specify properties for `fimath` objects in the Fixed-Point Toolbox. Refer to the following sections:

- “Setting `fimath` Properties at Object Creation” on page 4-4
- “Using Direct Property Referencing with `fimath`” on page 4-5

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 mode to saturate and the rounding mode to convergent,

```
F = fimath('OverflowMode','saturate','RoundMode','convergent')  
  
F =
```

```
RoundMode: convergent
OverflowMode: saturate
ProductMode: FullPrecision
MaxProductWordLength: 128
SumMode: FullPrecision
MaxSumWordLength: 128
CastBeforeSum: true
```

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 this by using a period to index into a property by name.

For example, to get the RoundMode of F,

```
F.RoundMode

ans =

convergent
```

To set the OverflowMode of F,

```
F.OverflowMode = 'wrap'

F =
```

```
RoundMode: convergent
OverflowMode: wrap
ProductMode: FullPrecision
MaxProductWordLength: 128
SumMode: FullPrecision
MaxSumWordLength: 128
CastBeforeSum: true
```

Using fimath Objects to Perform Fixed-Point Arithmetic

The `fimath` object encapsulates the math properties of the Fixed-Point Toolbox, and is itself a property of the `fi` object. Every `fi` object has a `fimath` object as a property.

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

```
a =
```

```
3.1416
```

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

```
      RoundMode: round  
      OverflowMode: saturate  
      ProductMode: FullPrecision  
MaxProductWordLength: 128  
      SumMode: FullPrecision  
MaxSumWordLength: 128  
      CastBeforeSum: true
```

```
a.fimath
```

```
ans =
```

```
      RoundMode: round  
      OverflowMode: saturate  
      ProductMode: FullPrecision  
MaxProductWordLength: 128  
      SumMode: FullPrecision  
MaxSumWordLength: 128  
      CastBeforeSum: true
```

To perform arithmetic with `+`, `-`, `.*`, or `*`, two `fi` operands must have the same `fimath` properties.

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

```
ans =
```

```
1
```

```
a + b
```

```
ans =
```

```
11.1416
```

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

```
RoundMode: round  
OverflowMode: saturate  
ProductMode: FullPrecision  
MaxProductWordLength: 128  
SumMode: FullPrecision  
MaxSumWordLength: 128  
CastBeforeSum: true
```

Using fimath to Share Arithmetic Rules

You can use a `fimath` object to define common arithmetic rules that you would like to use for many `fi` objects. You can then create multiple `fi` objects, using the same `fimath` object for each. To do so, you also need to create a `numerictype` object to define a common data type and scaling. Refer to Chapter 6, “Working with `numerictype` Objects” for more information on `numerictype` objects. The following example shows the creation of a `numerictype` object and `fimath` object, which are then used 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
           Signed: true
        WordLength: 32
    FractionLength: 30

F = fimath('RoundMode', 'floor', 'OverflowMode', 'wrap')

F =

        RoundMode: floor
    OverflowMode: wrap
    ProductMode: FullPrecision
MaxProductWordLength: 128
        SumMode: FullPrecision
    MaxSumWordLength: 128
    CastBeforeSum: true

a = fi(pi, T, F)

a =

    -0.8584
```



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

          RoundMode: floor
          OverflowMode: wrap
          ProductMode: FullPrecision
        MaxProductWordLength: 128
          SumMode: FullPrecision
        MaxSumWordLength: 128
        CastBeforeSum: true

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

b =

    1.5708
```

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

          RoundMode: floor
          OverflowMode: wrap
          ProductMode: FullPrecision
        MaxProductWordLength: 128
          SumMode: FullPrecision
        MaxSumWordLength: 128
        CastBeforeSum: true
```

Using fimath ProductMode and SumMode

The following example shows the differences among the FullPrecision, KeepLSB, KeepMSB, and SpecifyPrecision settings of the ProductMode and SumMode properties. To follow along, first set the following display, overflow logging, and fixed-point math preferences:

```
p = fipref;
p.NumericTypeDisplay = 'short';
p.FimathDisplay = 'none';
p.LoggingMode = 'OverflowAndUnderflow';
F = fimath('OverflowMode','wrap','RoundMode','floor',...
    'CastBeforeSum',false);
warning off
format compact
```

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

```
a = fi(pi, true, 8)
a =
    3.1563
    s8,5
b = fi(exp(1), true, 8)
b =
    2.7188
    s8,5
```

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 =
    3.1563    %011.00101
```

```

        s8,5
b
b =
    2.7188    %010.10111
        s8,5
a*b
ans =
    8.5811    %001000.1001010011
        s16,10
a+b
ans =
    5.8750    %0101.11100
        s9,5

```

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 is aligned 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 =
    3.1563    %011.00101
        s8,5
b

```

```
b =
    2.7188    %010.10111
           s8,5
a*b
ans =
    0.5811    %00.1001010011
           s12,10
a+b
ans =
    5.8750    %0000101.11100
           s12,5
```

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 is aligned 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 =
    3.1563    %011.00101
         s8,5
b =
    2.7188    %010.10111
         s8,5
a*b
ans =
    8.5781    %001000.100101
         s12,6
a+b
ans =
    5.8750    %0101.11100000
         s12,8

```

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.

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 =
    3.1563    %011.00101
         s8,5
b
b =
    2.7188    %010.10111
         s8,5
a*b
ans =
    0.5781    %0.1001010
         s8,7
a+b
ans =
   -0.1250    %1.1110000
         s8,7
```

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 2 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.

fimath Object Functions

The following functions operate directly on `fimath` objects:

- `add`
- `disp`
- `fimath`
- `isequal`
- `isfimath`
- `mpy`
- `sub`

You can learn about the functions associated with `fimath` objects in the Function Reference in the Fixed-Point Toolbox online documentation.

Working with fipref Objects

“Constructing fipref Objects” (p. 5-2)	Teaches you how to create fipref objects
“fipref Object Properties” (p. 5-3)	Tells you how to find more information about the properties associated with fipref objects, and shows you how to set these properties
“Using fipref Objects to Set Display Preferences” (p. 5-5)	Gives examples of using fipref objects to set display preferences for fi objects
“Using fipref Objects to Set Logging Preferences” (p. 5-7)	Gives examples of using fipref objects to set logging preferences for fi objects
“fipref Object Functions” (p. 5-10)	Introduces the functions in the toolbox that operate directly on fipref objects

Constructing `fipref` Objects

The `fipref` object defines the display and logging attributes for all `fi` objects. You can use the `fipref` constructor function to create a new object.

To get started, type

```
P = fipref
```

to create a default `fipref` object.

```
P =
```

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

The syntax

```
P = fipref(...'PropertyName','PropertyValue'...)
```

allows you to set properties for a `fipref` object at object creation with property name/property value pairs.

Your `fipref` settings persist throughout your MATLAB session. Use `reset(fipref)` to return to the default settings during your session. Use `savefipref` to save your display preferences for subsequent MATLAB sessions.

fipref Object Properties

All the properties of fipref objects are writable:

- `FimathDisplay` – Display options for the `fimath` attributes of a `fi` object
- `NumericTypeDisplay` – Display options for the numeric type attributes of a `fi` object
- `NumberDisplay` – Display options for the value of a `fi` object
- `LoggingMode` – Logging options for operations performed on `fi` objects

These properties are described in detail in Chapter 9, “Property Reference” in the online documentation. There are two ways to specify properties for fipref objects in the Fixed-Point Toolbox. Refer to the following sections:

- “Setting fipref Properties at Object Creation” on page 5-3
- “Using Direct Property Referencing with fipref” on page 5-3

Setting fipref Properties at Object Creation

You can set properties of fipref objects at the time of object creation by including properties after the arguments of the fipref constructor function. For example, to set `NumberDisplay` to `bin` and `NumericTypeDisplay` to `short`,

```
P = fipref('NumberDisplay', 'bin', 'NumericTypeDisplay', 'short')
```

```
P =
```

```

    NumberDisplay: 'bin'
 NumericTypeDisplay: 'short'
    FimathDisplay: 'full'
    LoggingMode: 'Off'
```

Using Direct Property Referencing with fipref

You can reference directly into a property for setting or retrieving fipref 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 `NumberDisplay` of `P`,

```
P.NumberDisplay
```

```
ans =
```

```
bin
```

To set the NumericTypeDisplay of P,

```
P.NumericTypeDisplay = 'full'
```

```
P =
```

```
    NumberDisplay: 'bin'  
    NumericTypeDisplay: 'full'  
    FimathDisplay: 'full'  
    LoggingMode: 'Off'
```

Using fipref Objects to Set Display Preferences

You use the `fipref` object to dictate three aspects of the display of `fi` objects: how the value of a `fi` object is displayed, how the `fimath` properties are displayed, and how the `numericType` properties are displayed.

For example, the following shows the default `fipref` display for a `fi` object:

```
a = fi(pi)

a =

    3.1416

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

        RoundMode: round
    OverflowMode: saturate
        ProductMode: FullPrecision
MaxProductWordLength: 128
        SumMode: FullPrecision
    MaxSumWordLength: 128
    CastBeforeSum: true
```

Now, 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'

a

a =

0110010010001000
(two's complement bin)
s16,13

Using fipref Objects to Set Logging Preferences

Overflows and underflows are logged as warnings for all assignment, plus, minus, and multiplication operations when the `fipref` `LoggingMode` property is set to `OverflowAndUnderflow`. 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 = 'OverflowAndUnderflow';
```

- 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 thrown for overflows and underflows in assignment operations. For example, try:

```
a(1) = pi  
Warning: 1 overflow 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 occurred in the fi assignment operation.  
  
a =  
      0    1.9844    1.9844    1.9844    1.9844
```

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

```
a+a  
Warning: 12 overflows occurred in the fi + operation.  
  
ans =  
      0    1.0000    1.0000    1.0000    1.0000  
  
a-a  
Warning: 8 overflows occurred in the fi - operation.  
  
ans =  
      0    0    0    0    0
```

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

```
a.*a  
Warning: 4 product overflows occurred in the fi .* operation.  
  
ans =  
      0    1.0000    1.0000    1.0000    1.0000  
  
a*a'  
Warning: 4 product overflows occurred in the fi * operation.  
Warning: 3 sum overflows occurred in the fi * operation.
```



```
ans =
```

```
1.0000
```

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

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

```
dbstop if warning
```

to help you find the exact lines in an M-file that are causing overflows or underflows to occur.

Use

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

to only stop on lines that cause an underflow. Use

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

to only stop on lines that cause an overflow.

fipref Object Functions

The following functions operate directly on fipref objects:

- `disp`
- `fipref`
- `reset`
- `savefipref`

You can learn about the functions associated with fipref objects in the [Function Reference](#) in the online documentation.

Working with numerictype Objects

“Constructing numerictype Objects” (p. 6-2)	Teaches you how to create numerictype objects
“numerictype Object Properties” (p. 6-6)	Tells you how to find more information about the properties associated with numerictype objects, and shows you how to set these properties
“The numerictype Structure” (p. 6-10)	Presents the numerictype object as a MATLAB structure, and gives the valid fields and settings for those fields
“Using numerictype Objects to Share Data Type and Scaling Settings” (p. 6-12)	Gives an example of using a numerictype object to share modular data type and scaling information among multiple fi objects
“numerictype Object Functions” (p. 6-15)	Introduces the functions in the toolbox that operate directly on numerictype objects

Constructing numerictype Objects

numerictype objects define the data type and scaling attributes of fi objects. You can create numerictype objects in the Fixed-Point Toolbox 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 get started, type

```
T = numerictype
```

to create a default numerictype object.

```
T =
```

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

You can use the numerictype constructor function in the following ways:

- `T = numerictype` creates a default numerictype object.
- `T = numerictype(s)` creates a numerictype object with Fixed-point: binary point scaling, signedness `s`, 16-bit word length and 15-bit fraction length.
- `T = numerictype(s,w)` creates a numerictype object with Fixed-point: binary point scaling, signedness `s`, word length `w` and 15-bit fraction length.
- `T = numerictype(s,w,f)` creates a numerictype object with Fixed-point: binary point scaling, signedness `s`, word length `w` and fraction length `f`.
- `T = numerictype(s,w,slope,bias)` creates a numerictype object with Fixed-point: slope and bias scaling, signedness `s`, word length `w`, slope, and bias.

- `T = numerictype(s,w,slopeadjustmentfactor,fixedexponent,bias)` creates a numerictype object with Fixed-point: slope and bias scaling, signedness `s`, word length `w`, `slopeadjustmentfactor`, `fixedexponent`, and `bias`.
- `T = numerictype(property1,value1, ...)` allows you to set properties for a numerictype object using property name/property value pairs.
- `T = numerictype(T1, property1, value1, ...)` allows you to make a copy of an existing numerictype object, while modifying any or all of the property values.

Examples of Constructing numerictype Objects

For example, the following creates a signed numerictype object with a 32-bit word length and 30-bit fraction length.

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

```
T =
```

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

If you omit the argument `f`, it is automatically set to the best precision possible.

```
T = numerictype(1, 32)
```

```
T =
```

```
      DataTypeMode: Fixed-point: binary point scaling
              Signed: true
              WordLength: 32
      FractionLength: 15
```

If you omit `w` and `f`, they are set automatically to 16 bits and the best precision possible, respectively.

```
T = numerictype(1)
```

```
T =
```

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

Constructing a numerictype Object with Property Name/Property Value Pairs

You can use property name/property value pairs to set numerictype properties when you create the object.

```
T = numerictype('Signed', true, 'DataTypeMode', ...
    'Fixed-point: slope and bias', 'WordLength', 32, 'Slope', ...
    2^-2, 'Bias', 4)
```

```
T =
```

```
        DataTypeMode: Fixed-point: slope and bias scaling
              Signed: true
            WordLength: 32
              Slope: 0.25
              Bias: 4
```

Copying a numerictype Object

To copy a numerictype object, use the numerictype constructor function.

```
T = numerictype
```

```
T =
```

```
        DataTypeMode: Fixed-point: binary point scaling
```

```
Signed: true  
WordLength: 16  
FractionLength: 15
```

```
U = numerictype(T)
```

```
U =
```

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

numerictype Object Properties

All the properties of a numerictype object are writable. However, the numerictype properties of a fi object are not writable once the fi object has been created:

- Bias – Bias
- DataType – Data type category
- DataTypeMode – Data type and scaling mode
- FixedExponent – Fixed-point exponent
- SlopeAdjustmentFactor – Slope adjustment
- FractionLength – Fraction length of the stored integer value, in bits
- Scaling – Fixed-point scaling mode
- Signed – Signed or unsigned
- Slope – Slope
- WordLength – Word length of the stored integer value, in bits

These properties are described in detail in the Chapter 9, “Property Reference” in the online documentation. There are two ways to specify properties for numerictype objects in the Fixed-Point Toolbox. Refer to the following sections:

- “Setting numerictype Properties at Object Creation” on page 6-6
- “Using Direct Property Referencing with numerictype Objects” on page 6-7
- “Setting numerictype Properties in the Model Explorer” on page 6-7

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
            Signed: true
            WordLength: 32
        FractionLength: 30
    
```

Using 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
            Signed: true
            WordLength: 32
        FractionLength: 31
    
```

Setting numerictype Properties in the Model Explorer

You can view and change the properties for any numerictype object defined in the MATLAB workspace in the Model Explorer. Open the Model Explorer by selecting **View > Model Explorer** in any Simulink model.

The snapshot below shows the Model Explorer when you define the following numerictype objects in the MATLAB workspace:

```
T = numerictype
```

```
T =
```

```
    DataTypeMode: Fixed-point: binary point scaling
```

```
    Signed: true
```

```
    WordLength: 16
```

```
    FractionLength: 15
```

```
U = numerictype('DataTypeMode', 'Fixed-point: slope and bias')
```

```
U =
```

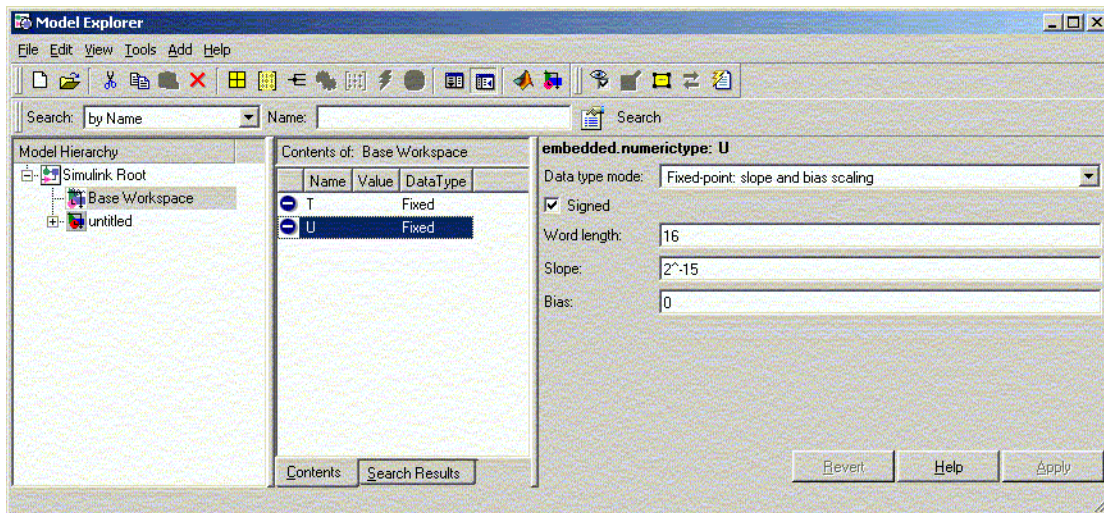
```
    DataTypeMode: Fixed-point: slope and bias scaling
```

```
    Signed: true
```

```
    WordLength: 16
```

```
    Slope: 2^-15
```

```
    Bias: 0
```



Select the **Base Workspace** node in the **Model Hierarchy** pane to view the current objects in the **Contents** pane. When you select a numerictype object in the **Contents** pane, you can view and change its properties in the **Dialog** pane.

The numerictype Structure

The numerictype object contains all the data type and scaling attributes of a `fi` object. The object acts the same as any MATLAB structure, except that it only lets you set valid values for defined fields. The following table shows the possible settings of each field of the structure that is valid for `fi` objects.

DataTypeMode	Data-Type	Scaling	Signed	Word-Length	Fraction-Length	SlopeBias	
<i>Fully specified fixed-point data types</i>							
Fixed-point: binary point scaling	fixed	BinaryPoint	1/0	w	f	1	0
Fixed-point: slope and bias scaling	fixed	SlopeBias	1/0	w	N/A	s	b
<i>Partially specified fixed-point data type</i>							
Fixed-point: unspecified scaling	fixed	Unspecified	1/0	w	N/A	N/A	N/A
<i>Built-in data types</i>							
int8	fixed	BinaryPoint	1	8	0	1	0
int16	fixed	BinaryPoint	1	16	0	1	0
int32	fixed	BinaryPoint	1	32	0	1	0
uint8	fixed	BinaryPoint	0	8	0	1	0
uint16	fixed	BinaryPoint	0	16	0	1	0
uint32	fixed	BinaryPoint	0	32	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 structure is related to the SlopeAdjustmentFactor and FixedExponent properties by

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

The FixedExponent and FractionLength properties are related by

$$\text{fixed exponent} = -\text{fraction length}$$

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

Stored Integer Value and Real World Value

The numerictype StoredIntegerValue and RealWorldValue properties are related according to

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

which is equivalent to

$$\begin{aligned} \text{real-world value} &= \text{stored integer value} \\ &\times (\text{slope adjustment factor} \times 2^{\text{fixed exponent}}) + \text{bias} \end{aligned}$$

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

Using numerictype Objects to Share Data Type and Scaling Settings

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. The following example shows the creation of a numerictype object, which is then used to create two fi objects with the same numerictype attributes:

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

T =

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

a = fi(pi,T)

a =

        3.1415926553309

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

        RoundMode: round
    OverflowMode: saturate
        ProductMode: FullPrecision
MaxProductWordLength: 128
        SumMode: FullPrecision
MaxSumWordLength: 128
    CastBeforeSum: true
```

```

b = fi(pi/2, T)

b =

    1.5707963258028

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

    RoundMode: round
    OverflowMode: saturate
    ProductMode: FullPrecision
    MaxProductWordLength: 128
    SumMode: FullPrecision
    MaxSumWordLength: 128
    CastBeforeSum: true

```

The following example shows the creation of a numerictype object with [Slope Bias] scaling, which is then used to create two fi objects with the same numerictype attributes:

```

T = numerictype('scaling','slopebias','slope', 2^2, 'bias', 0)

T =

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

c = fi(pi, T)

c =

    4

```

```
        DataTypeMode: Fixed-point: slope and bias scaling
          Signed: true
          WordLength: 16
            Slope: 2^2
            Bias: 0

          RoundMode: round
          OverflowMode: saturate
          ProductMode: FullPrecision
MaxProductWordLength: 128
          SumMode: FullPrecision
MaxSumWordLength: 128
          CastBeforeSum: true

d = fi(pi/2, T)

d =

    0
```

```
        DataTypeMode: Fixed-point: slope and bias scaling
          Signed: true
          WordLength: 16
            Slope: 2^2
            Bias: 0

          RoundMode: round
          OverflowMode: saturate
          ProductMode: FullPrecision
MaxProductWordLength: 128
          SumMode: FullPrecision
MaxSumWordLength: 128
          CastBeforeSum: true
```


numerictype Object Functions

The following functions operate directly on numerictype objects:

- `divide`
- `isequal`
- `isnumerictype`

You can learn about the functions associated with numerictype objects in the Function Reference in the online documentation.

Working with quantizer Objects

“Constructing quantizer Objects” (p. 7-2)	Explains how to create quantizer objects.
“quantizer Object Properties” (p. 7-4)	Outlines the properties of the quantizer objects
“Quantizing Data with quantizer Objects” (p. 7-6)	Discusses using quantizer objects to quantize data –how and what quantizing data does
“Transformations for Quantized Data ” (p. 7-8)	Offers a brief explanation of transforming quantized data between representations
“quantizer Object Functions” (p. 7-9)	Introduces the functions in the toolbox that operate directly on quantizer objects

Constructing quantizer Objects

You can use quantizer objects to quantize data sets before you pass them to `fi` objects. You can create quantizer objects in the Fixed-Point Toolbox in one of two ways:

- You can use the quantizer constructor function to create a new object.
- You can use the quantizer constructor function to copy a quantizer object.

To create a quantizer object with default properties, type

```
q = quantizer

q =

    DataMode = fixed
    RoundMode = floor
    OverflowMode = saturate
    Format = [16 15]

    Max = reset
    Min = reset
    NOverflows = 0
    NUnderflows = 0
    NOperations = 0
```

To copy a quantizer object, use the quantizer constructor function:

```
r = quantizer(q)

r =

    DataMode = fixed
    RoundMode = floor
    OverflowMode = saturate
    Format = [16 15]

    Max = reset
    Min = reset
    NOverflows = 0
```

```
NUnderflows = 0  
NOperations = 0
```

A listing of all the properties of the quantizer object `q` you just created is displayed along with the associated property values. All property values are set to defaults when you construct a quantizer object this way. See “quantizer Object Properties” on page 7-4 for more details.

quantizer Object Properties

You can set the values of some quantizer object properties. However, some properties have read-only values. The following sections cover settable and read-only properties:

- “Settable quantizer Object Properties” on page 7-4
- “Read-Only quantizer Object Properties” on page 7-5

Settable quantizer Object Properties

You can set the following four quantizer object properties:

- `DataMode` – Type of arithmetic used in quantization
- `Format` – Data format of a quantizer object
- `OverflowMode` – Overflow-handling mode
- `RoundMode` – Rounding mode

See the Property Reference in the online documentation for more details about these properties, including their possible values.

For example, to create a fixed-point quantizer object with

- The `Format` property value set to `[16,14]`
- The `OverflowMode` property value set to `'saturate'`
- The `RoundMode` property value set to `'ceil'`

type

```
q =  
quantizer('datamode','fixed','format',[16,14],'overflowmode',...  
         'saturate','roundmode','ceil')
```

You do not have to include quantizer object property names when you set quantizer object property values.

For example, you can create quantizer object `q` from the previous example by typing

```
q = quantizer('fixed',[16,14],'saturate','ceil')
```

Note You do not have to include default property values when you construct a quantizer object. In this example, you could leave out 'fixed' and 'saturate'.

Read-Only quantizer Object Properties

quantizer objects have five read-only properties:

- `Max` – Maximum value data has before a quantizer object is applied, that is, before quantization using `quantize`
- `Min` – Minimum value data has before a quantizer object is applied, that is, before quantization using `quantize`
- `NOperations` – Number of quantization operations that occur during quantization when you use a quantizer object
- `NOverflows` – Number of overflows that occur during quantization using `quantize`
- `NUnderflows` – Number of underflows that occur during quantization using `quantize`

These properties log quantization information each time you use `quantize` to quantize data with a quantizer object. The associated property values change each time you use `quantize` with a given quantizer object. You can reset these values to the default value using `reset`.

For an example, see “Quantizing Data with quantizer Objects” on page 7-6.

Quantizing Data with quantizer Objects

You construct a quantizer object to specify the quantization parameters to use when you quantize data sets. You can use the `quantize` function to quantize data according to a quantizer object's specifications.

Once you quantize data with a quantizer object, its data-related, read-only property values might change.

The following example shows

- How you use `quantize` to quantize data
- How quantization affects read-only properties
- How you reset read-only properties to their default values using `reset`

1 Construct an example data set and a quantizer object.

```
randn('state',0);  
x = randn(100,4);  
q = quantizer([16,14]);
```

2 Retrieve the values of the `Max` and `Noverflows` properties.

```
q.max  
  
ans =  
reset  
  
q.noverflows  
  
ans =  
0
```

3 Quantize the data set according to the quantizer object's specifications.

```
y = quantize(q,x);
```

4 Check the quantizer object property values.

```
q.max
```



```
ans =  
2.3726
```

```
q.noverflows
```

```
ans =  
15
```

5 Reset the read-only properties and check them.

```
reset(q)  
q.max
```

```
ans =  
reset
```

```
q.noverflows
```

```
ans =  
0
```

Transformations for Quantized Data

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

Use

- `num2bin` to convert data to binary
- `num2hex` to convert data to hexadecimal
- `hex2num` to convert hexadecimal data to numeric
- `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 =  
011  
010  
001  
000  
111  
110  
101  
100
```

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

quantizer Object Functions

The functions in the table below operate directly on quantizer objects

bin2num	copyobj	denormalmax	denormalmin	disp
eps	exponentbias	exponentlength	exponentmax	exponentmin
fractionlength	get	hex2num	isequal	length
max	min	noperations	noverflows	num2bin
num2hex	num2int	nunderflows	quantize	quantizer
randquant	range	realmax	realmin	reset
round	set	tostring	wordlength	

You can learn about the functions associated with quantizer objects in the [Function Reference](#) in the online documentation.

Interoperability with Other Products

“Using `fi` Objects with Simulink”
(p. 8-2)

Describes how to pass fixed-point data back and forth between the MATLAB workspace and Simulink models using Simulink blocks

“Using `fi` Objects with Signal Processing Blockset” (p. 8-7)

Describes how to pass fixed-point data back and forth between the MATLAB workspace and Simulink models using Signal Processing Blockset blocks

“Using `fi` Objects with Filter Design Toolbox” (p. 8-12)

Provides a brief description of how to use `fi` objects with `dfilt` objects in the Filter Design Toolbox

Using fi Objects with Simulink

Fixed-Point Toolbox `fi` objects can be used to pass fixed-point data back and forth between the MATLAB workspace and Simulink models.

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 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 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
        Signed: true
        WordLength: 16
    FractionLength: 15
```

```
        RoundMode: round
        OverflowMode: saturate
        ProductMode: FullPrecision
    MaxProductWordLength: 128
        SumMode: FullPrecision
    MaxSumWordLength: 128
    CastBeforeSum: true
```

```
s.signals.values = a
```

```
s =
```

```
    signals: [1x1 struct]
```

```
s.signals.dimensions = 2
```

```
s =
```

```
    signals: [1x1 struct]
```

```
s.time = [0:10]'
```

```
s =
```

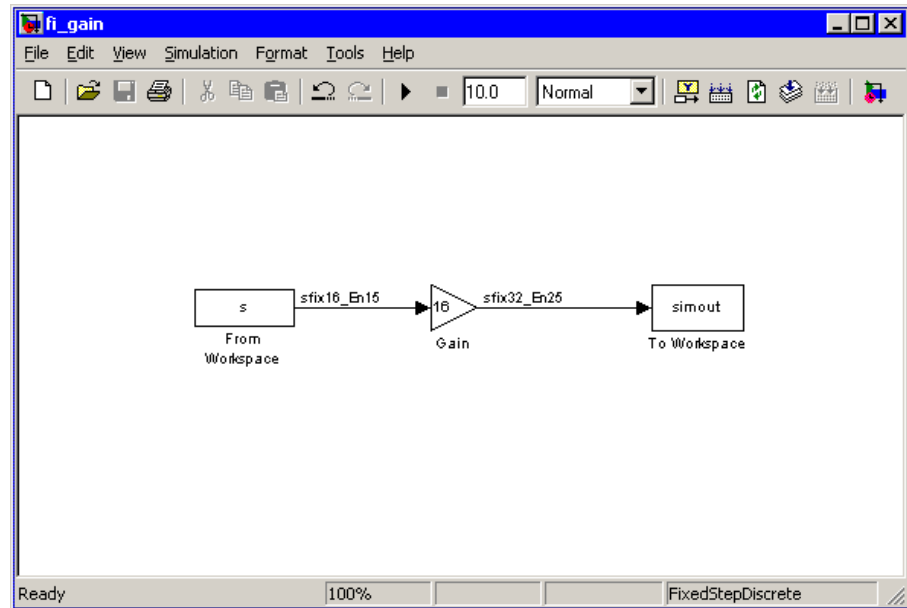
```
signals: [1x1 struct]
time: [11x1 double]
```

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 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 **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
    Signed: true
    WordLength: 32
    FractionLength: 25

    RoundMode: round
    OverflowMode: saturate
    ProductMode: FullPrecision
MaxProductWordLength: 128
    SumMode: FullPrecision
MaxSumWordLength: 128
    CastBeforeSum: true
```

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, select the **Log signal data** option in the signal's **Signal Properties** dialog box. For more information, refer to “Logging Signals” in the Using Simulink documentation.

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 Using Simulink documentation.

Using fi Objects with Signal Processing Blockset

Fixed-Point Toolbox `fi` objects can be used to pass fixed-point data back and forth between the MATLAB workspace and models using Signal Processing Blockset blocks.

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 the Signal Processing Blockset. 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 Signal To Workspace or Triggered To Workspace block from the Signal Processing Blockset. The fixed-point data is always written as a 2-D or 3-D array.

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 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
Signed: true
WordLength: 16
FractionLength: 15

RoundMode: round
OverflowMode: saturate
ProductMode: FullPrecision
MaxProductWordLength: 128
SumMode: FullPrecision
MaxSumWordLength: 128
CastBeforeSum: true

The Signal From Workspace block in the following model has the following settings:

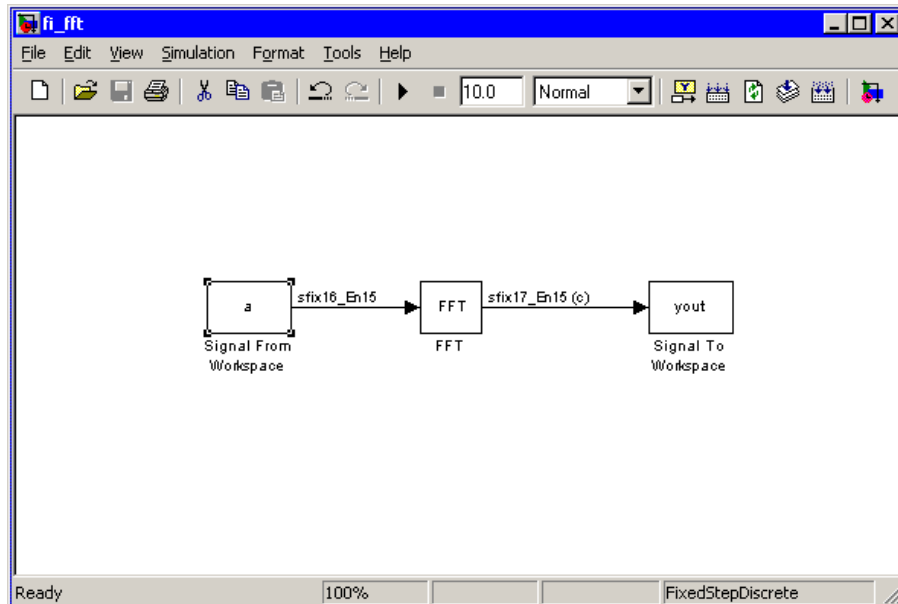
- **Signal** – a
- **Sample time** – 1
- **Samples per frame** – 2
- **Form output after final data value by** – Setting to zero

The following parameters in the **Solver** pane of the **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**

Remember, 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 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
Signed: true
WordLength: 17
FractionLength: 15

RoundMode: round
OverflowMode: saturate
ProductMode: FullPrecision
MaxProductWordLength: 128
SumMode: FullPrecision
MaxSumWordLength: 128

CastBeforeSum: true

Using `fi` Objects with Filter Design Toolbox

When the `Arithmetic` property is set to `'fixed'`, you can use an existing `fi` object as the input, states, or coefficients of a `dfilt` object in the Filter Design Toolbox. Also, fixed-point filters in the Filter Design Toolbox return `fi` objects as outputs. Refer to the Filter Design Toolbox documentation for more information.

Property Reference

“fi Object Properties” (p. 9-2)	Defines the <code>fi</code> object properties
“fimath Object Properties” (p. 9-5)	Defines the <code>fimath</code> object properties
“fipref Object Properties” (p. 9-10)	Defines the <code>fipref</code> object properties
“numerictype Object Properties” (p. 9-12)	Defines the <code>numerictype</code> object properties
“quantizer Object Properties” (p. 9-16)	Defines the <code>quantizer</code> object properties

fi Object Properties

The properties associated with `fi` objects are described in the following sections in alphabetical order.

Note The `fimath` properties and `numericity` properties are also properties of the `fi` object. Refer to “`fimath` Object Properties” on page 9-5 and “`numericity` Object Properties” on page 9-12 for more information.

bin

Stored integer value of a `fi` object in binary.

data

Numerical real-world value of a `fi` object

dec

Stored integer value of a `fi` object in decimal.

double

Real-world value of a `fi` object stored as a MATLAB `double`.

fimath

`fimath` object associated with a `fi` object. The default `fimath` object has the following settings:

```
RoundMode: round
OverflowMode: saturate
ProductMode: FullPrecision
MaxProductWordLength: 128
SumMode: FullPrecision
MaxSumWordLength: 128
CastBeforeSum: true
```

To learn more about `fi` math properties, refer to “`fi` math Object Properties” on page 9-5

hex

Stored integer value of a `fi` object in hexadecimal.

int

Stored integer value of a `fi` object, stored in a built-in MATLAB integer data type. You can also use `int8`, `int16`, `int32`, `uint8`, `uint16`, and `uint32` to get the stored integer value of a `fi` object in these formats.

NumericType

Structure containing all the data type and scaling attributes of a `fi` object. The `numericType` object acts the same as any MATLAB structure, except that it only lets you set valid values for defined fields. The following table shows the possible settings of each field of the structure that is valid for `fi` objects.

DataTypeMode	DataType	Scaling	Signed	Word- Length	Fraction- Length	Slope	Bias
<i>Fully specified fixed-point data types</i>							
Fixed-point: binary point scaling	fixed	BinaryPoint	1/0	w	f	1	0
Fixed-point: slope and bias scaling	fixed	SlopeBias	1/0	w	N/A	s	b
<i>Partially specified fixed-point data type</i>							
Fixed-point: unspecified scaling	fixed	Unspecified	1/0	w	N/A	N/A	N/A
<i>Built-in data types</i>							
int8	fixed	BinaryPoint	1	8	0	1	0
int16	fixed	BinaryPoint	1	16	0	1	0

DataTypeMode	DataType	Scaling	Signed	Word- Length	Fraction- Length	Slope	Bias
int32	fixed	BinaryPoint	1	32	0	1	0
uint8	fixed	BinaryPoint	0	8	0	1	0
uint16	fixed	BinaryPoint	0	16	0	1	0
uint32	fixed	BinaryPoint	0	32	0	1	0

You cannot change the numeric type properties of a `fi` object after `fi` object creation.

oct

Stored integer value of a `fi` object in octal.

fimath Object Properties

The properties associated with `fimath` objects are described in the following sections in alphabetical order.

CastBeforeSum

Whether both operands are cast to the sum data type before addition. Possible values of this property are 1 (cast before sum) and 0 (do not cast before sum).

The default value of this property is 1 (true).

MaxProductWordLength

Maximum allowable word length for the product data type.

The default value of this property is 128.

MaxSumWordLength

Maximum allowable word length for the sum data type.

The default value of this property is 128.

OverflowMode

Overflow-handling mode. The value of the `OverflowMode` property can be one of the following strings.

- `saturate` – 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.

The default value of this property is `saturate`.

ProductFractionLength

Fraction length, in bits, of the product data type. This value can be any positive or negative integer. The product data type defines the data type of the result of a multiplication of two `fi` objects.

The default value of this property is automatically set to the best precision possible based on the value of the product word length.

ProductMode

Defines how the product data type is determined. In the following descriptions, let A and B be real operands, with [word length, fraction length] pairs $[W_a F_a]$ and $[W_b F_b]$, respectively. W_p is the product data type word length and F_p is the product data type fraction length.

- **FullPrecision** – The full precision of the result is kept. An error is generated if the calculated word length is greater than `MaxProductWordLength`.

$$W_p = W_a + W_b$$

$$F_p = F_a + F_b$$

- **KeepLSB** – (keep least significant bits) You specify the product data type word length, while the fraction length is set to maintain the least significant bits of the product. In this mode, full precision is kept, but overflow is possible. This behavior models the C language integer operations.

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

$$F_p = F_a + F_b$$

- **KeepMSB** – (keep most significant bits) You specify the product data type word length, while the fraction length is set to maintain the most significant bits of the product. In this mode, overflow is prevented, but precision may be lost.

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

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

where

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

- `SpecifyPrecision` – You specify both the word length and fraction length of the product data type.

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

$$F_p = \text{specified in the } \text{ProductFractionLength} \text{ property}$$

The default value of this property is `FullPrecision`.

ProductWordLength

Word length, in bits, of the product data type. This value must be a positive integer. The product data type defines the data type of the result of a multiplication of two `fi` objects.

The default value of this property is 32.

RoundMode

The rounding mode. The value of the `RoundMode` property can be one of the following strings:

- `ceil` – Round toward positive infinity.
- `convergent` – Round toward nearest. Ties round to even numbers.
- `fix` – Round toward zero.
- `floor` – Round toward negative infinity.
- `round` – Round toward nearest. Ties round to the number toward positive infinity.

The default value of this property is `round`.

SumFractionLength

The fraction length, in bits, of the sum data type. This value can be any positive or negative integer. The sum data type defines the data type of the result of a sum of two `fi` objects.

The default value of this property is automatically set to the best precision possible based on the sum word length.

SumMode

Defines how the sum data type is determined. In the following descriptions, let A and B be real operands, with [word length, fraction length] pairs $[W_a, F_a]$ and $[W_b, F_b]$, respectively. W_s is the sum data type word length and F_s is the sum data type fraction length.

Note In the case where there are two operands, as in $A + B$, *NumberOfSummands* is 2, and $\text{ceil}(\log_2(\text{NumberOfSummands})) = 1$. In $\text{sum}(A)$, the *NumberOfSummands* is $\text{size}(A, 1)$.

- **FullPrecision** – The full precision of the result is kept. An error is generated if the calculated word length is greater than `MaxSumWordLength`.

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

where

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

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

- **KeepLSB** – (keep least significant bits) You specify the sum data type word length, while the fraction length is set to maintain the least significant bits of the sum. In this mode, full precision is kept, but overflow is possible. This behavior models the C language integer operations.

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

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

- **KeepMSB** – (keep most significant bits) You 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. In this mode, overflow is prevented, but precision may be lost.

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

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

where

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

- **SpecifyPrecision** – You specify both the word length and fraction length of the sum data type.

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

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

The default value of this property is `FullPrecision`.

SumWordLength

The word length, in bits, of the sum data type. This value must be a positive integer. The sum data type defines the data type of the result of a sum of two `fi` objects.

The default value of this property is 32.

fipref Object Properties

The properties associated with `fipref` objects are described in the following sections in alphabetical order.

FimathDisplay

Display options for the `fimath` attributes of a `fi` object

- `full` – Displays all of the `fimath` attributes of a fixed-point object
- `none` – None of the `fimath` attributes are displayed.

The default value of this property is `full`.

LoggingMode

Logging options for operations performed on `fi` objects

- `off` – No logging
- `overflowandunderflow` – Overflows and underflows are logged.

Overflows and underflows for assignment, plus, minus, and multiplication operations are logged as warnings when `LoggingMode` is set to `overflowandunderflow`.

The default value of this property is `off`.

NumericTypeDisplay

Display options for the `numerictype` attributes of a `fi` object

- `full` – Displays all the `numerictype` attributes of a fixed-point object
- `none` – None of the `numerictype` attributes are displayed.
- `short` – Displays an abbreviated notation of the fixed-point data type and scaling of a fixed-point object in the format `xWL,FL` where
 - `x` is `s` for signed and `u` for unsigned.
 - `WL` is the word length.

- FL is the fraction length.

The default value of this property is `full`.

NumberDisplay

Display options for the value of a `fi` object

- `bin` – Displays the stored integer value in binary format
- `dec` – Displays the stored integer value in unsigned decimal format
- `RealWorldValue` – Displays the stored integer value in the format specified by the MATLAB format function
- `hex` – Displays the stored integer value in hexadecimal format
- `int` – Displays the stored integer value in signed decimal format
- `none` – No value is displayed.

The default value of this property is `RealWorldValue`. In this mode, the value of a `fi` object is displayed in the format specified by the MATLAB format function: `+`, `bank`, `compact`, `hex`, `long`, `long e`, `long g`, `loose`, `rat`, `short`, `short e`, or `short g`. `fi` objects in `rat` format are displayed according to

$$1/(2^{\text{fixed-point exponent}}) \times \text{stored integer}$$

numerictype Object Properties

The properties associated with numerictype objects are described in the following sections in alphabetical order.

Bias

Bias associated with a fi object. The bias is part of the numerical representation used to interpret a fixed-point number. Along with the slope, the bias forms the scaling of the number. 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

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

DataType

Data type associated with a fi object. The only possible value of this property is Fixed – Fixed-point or integer data type.

DataTypeMode

Data type and scaling associated with a fi object. The possible values of this property are

- Fixed-point: binary point scaling – Fixed-point data type and scaling defined by the word length and fraction length
- Fixed-point: slope and bias scaling – Fixed-point data type and scaling defined by the slope and bias
- Fixed-point: unspecified scaling – A temporary setting that is only allowed at fi object creation, in order to allow for the automatic assignment of a binary point best-precision scaling
- int8 – Built-in signed 8-bit integer
- int16 – Built-in signed 16-bit integer

- `int32` – Built-in signed 32-bit integer
- `uint8` – Built-in unsigned 8-bit integer
- `uint16` – Built-in unsigned 16-bit integer
- `uint32` – Built-in unsigned 32-bit integer

The default value of this property is `Fixed-point: binary point scaling`.

FixedExponent

Fixed-point exponent associated with a `fi` object. The exponent is part of the numerical representation used to express a fixed-point number. 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

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

The exponent of a fixed-point number is equal to the negative of the fraction length:

$$\text{fixed exponent} = -\text{fraction length}$$

FractionLength

Value of the `FractionLength` property is the fraction length of the stored integer value of a `fi` object, in bits. The fraction length can be any integer value. If you do not specify the fraction length of a `fi` object, it is set to the best possible precision.

This property is automatically set by default to the best precision possible based on the value of the word length.

Scaling

Fixed-point scaling mode of a `fi` object. The possible values of this property are

- `BinaryPoint` – 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, in order to allow for the automatic assignment of a binary point best precision scaling
- `Integer` – The `fi` object is an integer; the binary point is understood to be at the far right of the word, making the fraction length zero.

The default value of this property is `BinaryPoint`.

Signed

Whether a `fi` object is signed.

The default value of this property is 1 (signed).

Slope

Slope associated with a `fi` object. The slope is part of the numerical representation used to express a fixed-point number. Along with the bias, the slope forms the scaling of a fixed-point number. 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

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

SlopeAdjustmentFactor

Slope adjustment associated with a `fi` object. The slope adjustment is equivalent to the fractional slope of a fixed-point number. The fractional slope

is part of the numerical representation used to express a fixed-point number. Fixed-point numbers can be represented as

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

where the slope can be expressed as

$$\textit{slope} = \textit{fractional slope} \times 2^{\textit{fixed exponent}}$$

WordLength

Value of the WordLength property is the word length of the stored integer value of a fixed-point object, in bits. The word length can be any positive integer value.

The default value of this property is 16.

quantizer Object Properties

The properties associated with quantizer objects are described in the following sections in alphabetical order.

DataMode

Type of arithmetic used in quantization. This property can have the following values:

- `fixed` – Signed fixed-point calculations
- `float` – User-specified floating-point calculations
- `double` – Double-precision floating-point calculations
- `single` – Single-precision floating-point calculations
- `ufixed` – Unsigned fixed-point calculations

The default value of this property is `fixed`.

When you set the `DataMode` property value to `double` or `single`, the `Format` property value becomes read only.

Format

Data format of a quantizer object. The interpretation of this property value depends on the value of the `DataMode` property.

For example, whether you specify the `DataMode` property with `fixed`- or floating-point arithmetic affects the interpretation of the data format property. For some `DataMode` property values, the data format property is read only.

The following table shows you how to interpret the values for the `Format` property value when you specify it, or how it is specified in read-only cases.

DataMode Property Value	Interpreting the Format Property Values
fixed or ufixed	<p>You specify the Format property value as a vector. The number of bits for the quantizer object word length is the first entry of this vector, and the number of bits for the quantizer object fraction length is the second entry.</p> <p>The word length can range from 2 to the limits of memory on your PC. The fraction length can range from 0 to one less than the word length.</p>
float	<p>You specify the Format property value as a vector. The number of bits you want for the quantizer object word length is the first entry of this vector, and the number of bits you want for the quantizer object exponent length is the second entry.</p> <p>The word length can range from 2 to the limits of memory on your PC. The exponent length can range from 0 to 11.</p>
double	<p>The Format property value is specified automatically (is read only) when you set the DataMode property to double. The value is [64 11], specifying the word length and exponent length, respectively.</p>
single	<p>The Format property value is specified automatically (is read only) when you set the DataMode property to single. The value is [32 8], specifying the word length and exponent length, respectively.</p>

Max

Maximum value data has before a quantizer object is applied to it, that is, before quantization using `quantize`. The value of `Max` accumulates if you use the same quantizer object to quantize several data sets. You can reset the value using `reset`.

The `Max` property is read only.

Min

Minimum value data has before a quantizer object is applied to it, that is, before quantization using `quantize`. The value of `Min` accumulates if you

use the same quantizer object to quantize several data sets. You can reset the value using `reset`.

The `Min` property is read only.

NOperations

Number of quantization operations that occur during quantization when you use a quantizer object. This value accumulates when you use the same quantizer object to process several data sets. You reset the value using `reset`.

The default value of this property is 0.

The `NOperations` property is read only.

NOverflows

Number of overflows that occur during quantization using `quantize`. This value accumulates if you use the same quantizer object to quantize several data sets. You can reset the value using `reset`.

The default value of this property is 0.

The `NOverflows` property is read only.

NUnderflows

Number of underflows that occur during quantization using `quantize`. This value accumulates when you use the same quantizer object to quantize several data sets. You can reset the value using `reset`.

The default value of this property is 0.

The `NUnderflows` property is read only.

OverflowMode

Overflow-handling mode. The value of the `OverflowMode` property can be one of the following strings:

- `saturate` – Overflows saturate.

When the values of data to be quantized lie outside the range of the largest and smallest representable numbers (as specified by the data format properties), these values are quantized to the value of either the largest or smallest representable value, depending on which is closest.

- `wrap` – Overflows wrap to the range of representable values.

When the values of data to be quantized lie outside the range of the largest and smallest representable numbers (as specified by the data format properties), these values are wrapped back into that range using modular arithmetic relative to the smallest representable number.

The default value of this property is `saturate`.

Note Floating-point numbers that extend beyond the dynamic range overflow to $\pm\text{inf}$.

The `OverflowMode` property value is set to `saturate` and becomes a read-only property when you set the value of the `DataMode` property to `float`, `double`, or `single`.

RoundMode

Rounding mode. The value of the `RoundMode` property can be one of the following strings:

- `ceil` – Round up to the next allowable quantized value.
- `convergent` – Round to the nearest allowable quantized value. Numbers that are exactly halfway between the two nearest allowable quantized values are rounded up only if the least significant bit (after rounding) would be set to 0.
- `fix` – Round negative numbers up and positive numbers down to the next allowable quantized value.
- `floor` – Round down to the next allowable quantized value.

- `round` – Round to the nearest allowable quantized value. Numbers that are halfway between the two nearest allowable quantized values are rounded up.

The default value of this property is `floor`.

Functions — Categorical List

- “Bitwise Functions” on page 10-2
- “Constructor and Property Functions” on page 10-2
- “Data Manipulation Functions” on page 10-3
- “Data Type Functions” on page 10-5
- “Data Quantizing Functions” on page 10-6
- “Element-Wise Logical Operator Functions” on page 10-6
- “Math Operation Functions” on page 10-6
- “Matrix Manipulation Functions” on page 10-8
- “Plotting Functions” on page 10-9
- “Radix Conversion Functions” on page 10-12
- “Relational Operator Functions” on page 10-13
- “Statistics Functions” on page 10-14
- “Subscripted Assignment and Reference Functions” on page 10-15
- “fi Object Functions” on page 10-16
- “fimath Object Functions” on page 10-18
- “fipref Object Functions” on page 10-19
- “numerictype Object Functions” on page 10-20
- “quantizer Object Functions” on page 10-21

Bitwise Functions

<code>bitand</code>	Return the bitwise AND of two <code>fi</code> objects
<code>bitcmp</code>	Return the bitwise complement of a <code>fi</code> object
<code>bitget</code>	Return the bit at a certain position
<code>bitor</code>	Return the bitwise OR of two <code>fi</code> objects
<code>bitset</code>	Set the bit at a certain position
<code>bitshift</code>	Shift bits specified number of places
<code>bitxor</code>	Return the bitwise exclusive OR of two <code>fi</code> objects

Constructor and Property Functions

<code>copyobj</code>	Make an independent copy of a quantizer object
<code>fi</code>	Construct a <code>fi</code> object
<code>fimath</code>	Construct a <code>fimath</code> object
<code>fipref</code>	Construct a <code>fipref</code> object
<code>get</code>	Return the property values of a quantizer object
<code>inspect</code>	Display Property Inspector
<code>numerictype</code>	Construct a <code>numerictype</code> object
<code>quantizer</code>	Construct a quantizer object
<code>reset</code>	Reset one or more objects to their initial conditions
<code>savefipref</code>	Save <code>fi</code> preferences for the next MATLAB session

<code>set</code>	Set or display property values for quantizer objects
<code>stripscaling</code>	Return the stored integer of a <code>fi</code> object
<code>tostring</code>	Convert a quantizer object to a string

Data Manipulation Functions

<code>denormalmax</code>	Return the largest denormalized quantized number for a quantizer object
<code>denormalmin</code>	Return the smallest denormalized quantized number for a quantizer object
<code>eps</code>	Return the quantized relative accuracy for <code>fi</code> objects or quantizer objects
<code>exponentbias</code>	Return the exponent bias for a quantizer object
<code>exponentlength</code>	Return the exponent length of a quantizer object
<code>exponentmax</code>	Return the maximum exponent for a quantizer object
<code>exponentmin</code>	Return the minimum exponent for a quantizer object
<code>fractionlength</code>	Return the fraction length of a quantizer object

<code>isequal</code>	Determine whether the real-world values of two <code>fi</code> objects are equal, or determine whether the properties of two <code>fimath</code> , <code>numericType</code> , or <code>quantizer</code> objects are equal
<code>isfi</code>	Determine whether a variable is a <code>fi</code> object
<code>isfimath</code>	Determine whether a variable is a <code>fimath</code> object
<code>isnumericType</code>	Determine whether a variable is a <code>numericType</code> object
<code>ispropequal</code>	Determine whether the properties of two <code>fi</code> objects are equal
<code>issigned</code>	Determine whether a <code>fi</code> object is signed
<code>lowerbound</code>	Return lower bound of range of <code>fi</code> object
<code>lsb</code>	Return the scaling of the least significant bit of a <code>fi</code> object
<code>range</code>	Return the numerical range of a <code>fi</code> object or <code>quantizer</code> object
<code>realmax</code>	Return the largest positive fixed-point value or quantized number
<code>realmin</code>	Return the smallest positive normalized fixed-point value or quantized number
<code>rescale</code>	Change the scaling of a <code>fi</code> object
<code>upperbound</code>	Return upper bound of range of <code>fi</code> object
<code>wordlength</code>	Return the word length of a <code>quantizer</code> object

Data Type Functions

<code>double</code>	Return the double-precision floating-point real-world value of a <code>fi</code> object
<code>int</code>	Return the smallest built-in integer in which the stored integer value of a <code>fi</code> object will fit
<code>int16</code>	Return the stored integer value of a <code>fi</code> object as a built-in <code>int16</code>
<code>int32</code>	Return the stored integer value of a <code>fi</code> object as a built-in <code>int32</code>
<code>int8</code>	Return the stored integer value of a <code>fi</code> object as a built-in <code>int8</code>
<code>intmax</code>	Return the largest positive stored integer value representable by the <code>numerictype</code> of a <code>fi</code> object
<code>intmin</code>	Return smallest stored integer value representable by <code>numerictype</code> of <code>fi</code> object
<code>logical</code>	Convert numeric values to logical
<code>single</code>	Return the single-precision floating-point real-world value of a <code>fi</code> object
<code>uint16</code>	Return the stored integer value of a <code>fi</code> object as a built-in <code>uint16</code>
<code>uint32</code>	Return the stored integer value of a <code>fi</code> object as a built-in <code>uint32</code>
<code>uint8</code>	Return the stored integer value of a <code>fi</code> object as a built-in <code>uint8</code>

Data Quantizing Functions

convergent	Apply convergent rounding
quantize	Apply a quantizer object to data
randquant	Generate a uniformly distributed, quantized random number using a quantizer object
round	Round input data using a quantizer object without checking for overflow

Element-Wise Logical Operator Functions

all	Determine if all array elements are nonzero
and	Find logical AND of array or scalar inputs
any	Determine if any array elements are nonzero
not	Find logical NOT of array or scalar input
or	Find logical OR of array or scalar inputs

Math Operation Functions

abs	Return the absolute value of a <code>fi</code> object
add	Add two objects using a <code>fimath</code> object

<code>complex</code>	Construct a complex <code>fi</code> object from real and imaginary parts
<code>conj</code>	Return the complex conjugate of a <code>fi</code> object
<code>divide</code>	Divide two objects using a <code>numeric</code> object
<code>imag</code>	Return the imaginary part
<code>innerprodintbits</code>	Return the number of integer bits needed for a fixed-point inner product
<code>minus</code>	Return the matrix difference between <code>fi</code> objects
<code>mpy</code>	Multiply two objects using a <code>fimath</code> object
<code>mtimes</code>	Return the matrix product of <code>fi</code> objects
<code>plus</code>	Return the matrix sum of <code>fi</code> objects
<code>pow2</code>	Multiply by a power of 2
<code>real</code>	Return real part of complex number
<code>sign</code>	Perform signum function on array
<code>sub</code>	Subtract two objects using a <code>fimath</code> object
<code>sum</code>	Return sum of array elements
<code>times</code>	Return the result of element-by-element multiplication of <code>fi</code> objects
<code>uminus</code>	Negate the elements of a <code>fi</code> object array
<code>uplus</code>	Unary plus

Matrix Manipulation Functions

<code>buffer</code>	Buffer signal vector into matrix of data frames
<code>ctranspose</code>	Return the complex conjugate transpose of a <code>fi</code> object
<code>diag</code>	Return diagonal matrices or the diagonals of a matrix
<code>disp</code>	Display an object
<code>end</code>	Indicate last index of array
<code>find</code>	Find indices and values of nonzero elements
<code>hankel</code>	Return a Hankel matrix
<code>horzcat</code>	Horizontally concatenate two or more <code>fi</code> objects
<code>ipermute</code>	Inverse permute the dimensions of a multidimensional array
<code>iscolumn</code>	Determine whether a <code>fi</code> object is a column vector
<code>isempty</code>	Determine if array is empty
<code>isnumeric</code>	Determine if input is numeric array
<code>isobject</code>	Determine if input is MATLAB OOPS object
<code>isreal</code>	Determine if all array elements are real numbers
<code>isrow</code>	Determine whether a <code>fi</code> object is a row vector
<code>isscalar</code>	Determine if input is scalar
<code>isvector</code>	Determine if input is vector
<code>length</code>	Return the length of a vector
<code>ndims</code>	Return number of array dimensions

<code>permute</code>	Rearrange the dimensions of a multidimensional array
<code> repmat</code>	Replicate and tile an array
<code> reshape</code>	Reshape array
<code> size</code>	Return array dimensions
<code> squeeze</code>	Remove singleton dimensions
<code> toeplitz</code>	Create Toeplitz matrix
<code> transpose</code>	Return the transpose
<code> tril</code>	Return the lower triangular part of a matrix
<code> vertcat</code>	Vertically concatenate two or more <code>f i</code> objects

Plotting Functions

<code> area</code>	Create a filled area 2-D plot
<code> bar</code>	Create a vertical bar graph
<code> barh</code>	Create a horizontal bar graph
<code> clabel</code>	Create contour plot elevation labels
<code> comet</code>	Create a 2-D comet plot
<code> comet3</code>	Create a 3-D comet plot
<code> compass</code>	Plot arrows emanating from the origin
<code> coneplot</code>	Plot velocity vectors as cones in a 3-D vector field
<code> contour</code>	Create a contour graph of a matrix
<code> contour3</code>	Create a 3-D contour plot

<code>contourc</code>	Create a two-level contour plot computation
<code>contourf</code>	Create a filled 2-D contour plot
<code>errorbar</code>	Plot error bars along a curve
<code>etreeplot</code>	Plot elimination tree
<code>ezcontour</code>	Easy-to-use contour plotter
<code>ezcontourf</code>	Easy-to-use filled contour plotter
<code>ezmesh</code>	Easy-to-use 3-D mesh plotter
<code>ezplot</code>	Easy-to-use function plotter
<code>ezplot3</code>	Easy-to-use 3-D parametric curve plotter
<code>ezpolar</code>	Easy-to-use polar coordinate plotter
<code>ezsurf</code>	Easy-to-use 3-D colored surface plotter
<code>ezsurfz</code>	Easy-to-use combination surface/contour plotter
<code>feather</code>	Plot velocity vectors
<code>fplot</code>	Plot a function between specified limits
<code>gplot</code>	Plot set of nodes using an adjacency matrix
<code>hist</code>	Create histogram plot
<code>histc</code>	Return histogram count
<code>line</code>	Create line object
<code>loglog</code>	Create log-log scale plot
<code>mesh</code>	Create mesh plot
<code>meshc</code>	Create mesh plot with contour plot
<code>meshz</code>	Create mesh plot with curtain plot
<code>patch</code>	Create patch graphics object

<code>pcolor</code>	Create pseudocolor plot
<code>plot</code>	Create linear 2-D plot
<code>plot3</code>	Create 3-D line plot
<code>plotmatrix</code>	Draw scatter plots
<code>plotyy</code>	Create graph with y-axes on both right and left sides
<code>polar</code>	Plot polar coordinates
<code>quiver</code>	Create quiver or velocity plot
<code>quiver3</code>	Create 3-D quiver or velocity plot
<code>rgbplot</code>	Plot colormap
<code>ribbon</code>	Create ribbon plot
<code>rose</code>	Create angle histogram
<code>scatter</code>	Create a scatter or bubble plot
<code>scatter3</code>	Create a 3-D scatter or bubble plot
<code>semilogx</code>	Create semilogarithmic plot with logarithmic x-axis
<code>semilogy</code>	Create semilogarithmic plot with logarithmic y-axis
<code>slice</code>	Create volumetric slice plot
<code>spy</code>	Visualize sparsity pattern
<code>stairs</code>	Create staircase graph
<code>stem</code>	Plot discrete sequence data
<code>stem3</code>	Plot 3-D discrete sequence data
<code>streamribbon</code>	Create a 3-D stream ribbon plot
<code>streamslice</code>	Draw streamlines in slice planes
<code>streamtube</code>	Create a 3-D stream tube plot
<code>surf</code>	Create 3-D shaded surface plot
<code>surfc</code>	Create 3-D shaded surface plot with contour plot

<code>surf1</code>	Create a surface plot with colormap-based lighting
<code>surfnorm</code>	Compute and display 3–D surface normals
<code>text</code>	Create text object in current axes
<code>treeplot</code>	Plot picture of tree
<code>trimesh</code>	Create triangular mesh plot
<code>triplot</code>	Create 2–D triangular plot
<code>trisurf</code>	Create triangular surface plot
<code>triu</code>	Return the upper triangular part of a matrix
<code>voronoi</code>	Create Voronoi diagram
<code>voronoin</code>	Create n-dimensional Voronoi diagram
<code>waterfall</code>	Create waterfall plot
<code>xlim</code>	Set or query x-axis limits
<code>ylim</code>	Set or query y-axis limits
<code>zlim</code>	Set or query z-axis limits

Radix Conversion Functions

<code>bin</code>	Return the binary representation of the stored integer of a <code>fi</code> object as a string
<code>bin2num</code>	Convert a two’s complement binary string to a number using a quantizer object
<code>dec</code>	Return the unsigned decimal representation of the stored integer of a <code>fi</code> object as a string

hex	Return the hexadecimal representation of the stored integer of a <code>fi</code> object as a string
hex2num	Convert a hexadecimal string to a number using a quantizer object
num2bin	Convert a number to a binary string using a quantizer object
num2hex	Convert a number to its hexadecimal equivalent using a quantizer object
num2int	Convert a number to a signed integer
oct	Return the octal representation of the stored integer of a <code>fi</code> object as a string
sdec	Return signed decimal representation of stored integer of <code>fi</code> object as string

Relational Operator Functions

eq	Determine whether the real-world values of two <code>fi</code> objects are equal
ge	Determine whether the real-world value of one <code>fi</code> object is greater than or equal to another
gt	Determine whether the real-world value of one <code>fi</code> object is greater than another
le	Determine whether the real-world value of a <code>fi</code> object is less than or equal to another

<code>lt</code>	Determine whether the real-world value of a <code>fi</code> object is less than another
<code>ne</code>	Determine whether the real-world values of two <code>fi</code> objects are not equal

Statistics Functions

<code>max</code>	Return the largest element in an array of <code>fi</code> objects or the maximum value of a quantizer object before quantization
<code>min</code>	Return the smallest element in an array of <code>fi</code> objects or the minimum value of a quantizer object before quantization
<code>noperations</code>	Return the number of quantization operations performed by a quantizer object
<code>noverflows</code>	Return the number of overflows from quantization operations performed by a quantizer object
<code>numberofelements</code>	Return number of data elements in <code>fi</code> array
<code>nunderflows</code>	Return the number of underflows from quantization operations performed by a quantizer object

Subscripted Assignment and Reference Functions

subsasgn

Subscripted assignment

subsref

Subscripted reference

fi Object Functions

The functions in the table below operate directly on `fi` objects.

<code>abs</code>	<code>all</code>	<code>and</code>	<code>any</code>	<code>area</code>
<code>bar</code>	<code>barh</code>	<code>bin</code>	<code>bitand</code>	<code>bitcmp</code>
<code>bitget</code>	<code>bitor</code>	<code>bitshift</code>	<code>bitxor</code>	<code>buffer</code>
<code>clabel</code>	<code>comet</code>	<code>comet3</code>	<code>compass</code>	<code>complex</code>
<code>coneplot</code>	<code>conj</code>	<code>contour</code>	<code>contour3</code>	<code>contourc</code>
<code>contourf</code>	<code>ctranspose</code>	<code>dec</code>	<code>diag</code>	<code>double</code>
<code>end</code>	<code>eps</code>	<code>eq</code>	<code>errorbar</code>	<code>etreeplot</code>
<code>ezcontour</code>	<code>ezcontourf</code>	<code>ezmesh</code>	<code>ezplot</code>	<code>ezplot3</code>
<code>ezpolar</code>	<code>ezsurf</code>	<code>ezsurfz</code>	<code>feather</code>	<code>fi</code>
<code>find</code>	<code>fplot</code>	<code>ge</code>	<code>get</code>	<code>gplot</code>
<code>gt</code>	<code>hankel</code>	<code>hex</code>	<code>hist</code>	<code>histc</code>
<code>horzcat</code>	<code>innerprodintbits</code>	<code>inspect</code>	<code>int</code>	<code>int8</code>
<code>int16</code>	<code>int32</code>	<code>intmax</code>	<code>intmin</code>	<code>ipermute</code>
<code>iscolumn</code>	<code>isequal</code>	<code>isfi</code>	<code>isnumeric</code>	<code>isobject</code>
<code>ispropequal</code>	<code>isrow</code>	<code>assigned</code>	<code>le</code>	<code>line</code>
<code>logical</code>	<code>lowerbound</code>	<code>lsb</code>	<code>lt</code>	<code>max</code>
<code>mesh</code>	<code>meshc</code>	<code>meshz</code>	<code>min</code>	<code>minus</code>
<code>mtimes</code>	<code>ne</code>	<code>not</code>	<code>numberofelements</code>	<code>oct</code>
<code>or</code>	<code>patch</code>	<code>pcolor</code>	<code>permute</code>	<code>plot</code>
<code>plot3</code>	<code>plotmatrix</code>	<code>plotyy</code>	<code>plus</code>	<code>polar</code>
<code>pow2</code>	<code>quiver</code>	<code>quiver3</code>	<code>range</code>	<code>realmax</code>
<code>realmin</code>	<code>rescale</code>	<code>rgbplot</code>	<code>ribbon</code>	<code>rose</code>
<code>scatter</code>	<code>scatter3</code>	<code>sdec</code>	<code>sign</code>	<code>single</code>
<code>slice</code>	<code>spy</code>	<code>stairs</code>	<code>stem</code>	<code>stem3</code>
<code>streamribbon</code>	<code>streamslice</code>	<code>streamtube</code>	<code>stripscaling</code>	<code>subsasgn</code>

sum	surf	surfc	surfl	surfnorm
text	times	toeplitz	treeplot	tril
trimesh	triplot	trisurf	triu	uint8
uint16	uint32	uminus	uplus	upperbound
vertcat	voronoi	voronoin	waterfall	xlim
ylim	zlim			

fimath Object Functions

The following functions operate directly on `fimath` objects.

- `add`
- `disp`
- `fimath`
- `isequal`
- `isfimath`
- `mpy`
- `sub`

fipref Object Functions

The following functions operate directly on fipref objects.

- disp
- fipref
- reset
- savefipref

numerictype Object Functions

The following functions operate directly on `numerictype` objects.

- `divide`
- `isequal`
- `isnumerictype`

quantizer Object Functions

The functions in the table below operate directly on quantizer objects.

bin2num	copyobj	denormalmax	denormalmin	disp
eps	exponentbias	exponentlength	exponentmax	exponentmin
fractionlength	get	hex2num	isequal	length
max	min	noperations	noverflows	num2bin
num2hex	num2int	nunderflows	quantize	quantizer
randquant	range	realmax	realmin	reset
round	set	tostring	wordlength	

Functions — Alphabetical List

abs

Purpose Return the absolute value of a `fi` object

Syntax `abs(a)`

Description `abs(a)` returns the absolute value of `fi` object `a`.

When the object `a` is real and has a signed data type, the absolute value of the most negative value is problematic since it is not representable. In this case, the absolute value saturates to the most positive value representable by the data type if the `OverflowMode` property is set to `saturate`. If `OverflowMode` is `wrap`, the absolute value of the most negative value has no effect.

`abs` does not support complex inputs.

Examples The following example shows the difference between the absolute value results for the most negative value representable by a signed data type when `OverflowMode` is `saturate` or `wrap`.

```
P = fipref('NumericTypeDisplay','full','FimathDisplay','full');  
a = fi(-128)
```

```
a =
```

```
-128
```

```
        DataTypeMode: Fixed-point: binary point scaling  
             Signed: true  
        WordLength: 16  
    FractionLength: 8
```

```
        RoundMode: round  
    OverflowMode: saturate  
        ProductMode: FullPrecision  
MaxProductWordLength: 128  
             SumMode: FullPrecision  
MaxSumWordLength: 128
```

```
CastBeforeSum: true
abs(a)
ans =
    127.9961
```

```
DataTypeMode: Fixed-point: binary point scaling
Signed: true
WordLength: 16
FractionLength: 8

RoundMode: round
OverflowMode: saturate
ProductMode: FullPrecision
MaxProductWordLength: 128
SumMode: FullPrecision
MaxSumWordLength: 128
CastBeforeSum: true
a.OverflowMode = 'wrap'

a =
    -128
```

```
DataTypeMode: Fixed-point: binary point scaling
Signed: true
WordLength: 16
FractionLength: 8

RoundMode: round
OverflowMode: wrap
ProductMode: FullPrecision
MaxProductWordLength: 128
SumMode: FullPrecision
```

abs

```
MaxSumWordLength: 128
CastBeforeSum: true
abs(a)
ans =
-128
```

```
DataTypeMode: Fixed-point: binary point scaling
Signed: true
WordLength: 16
FractionLength: 8

RoundMode: round
OverflowMode: wrap
ProductMode: FullPrecision
MaxProductWordLength: 128
SumMode: FullPrecision
MaxSumWordLength: 128
CastBeforeSum: true
```

Purpose Add two objects using a `fimath` object

Syntax `c = F.add(a,b)`

Description `c = F.add(a,b)` adds objects `a` and `b` using `fimath` object `F`. This is helpful in cases when you want to override the `fimath` objects of `a` and `b`, or if the `fimath` objects of `a` and `b` are different.

`a` and `b` must have the same dimensions unless one is a scalar. If either `a` or `b` is scalar, then `c` has the dimensions of the nonscalar object.

If either `a` or `b` is a `fi` object, and the other is a MATLAB built-in numeric type, then the built-in object is cast to the word length of the `fi` object, preserving best-precision fraction length.

Examples In this example, `c` is the 32-bit sum of `a` and `b` with fraction length 16:

```
a = fi(pi);
b = fi(exp(1));
F = fimath('SumMode','SpecifyPrecision','SumWordLength',
          32,'SumFractionLength',16);
c = F.add(a,b)
```

```
c =
```

```
5.8599
```

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

```
RoundMode: round
OverflowMode: saturate
ProductMode: FullPrecision
MaxProductWordLength: 128
SumMode: SpecifyPrecision
```

add

```
SumWordLength: 32
SumFractionLength: 16
CastBeforeSum: true
```

Algorithm

`c = F.add(a,b)` is equivalent to

```
a.fimath = F;
b.fimath = F;
c = a + b;
```

except that the `fimath` properties of `a` and `b` are not modified when you use the functional form.

See Also

`divide`, `fi`, `fimath`, `mpy`, `numericType`, `sub`, `sum`

Purpose Determine if all array elements are nonzero

Description Refer to the MATLAB `all` reference page for more information.

and

Purpose Find logical AND of array or scalar inputs

Description Refer to the MATLAB and reference page for more information.

Purpose Determine if any array elements are nonzero

Description Refer to the MATLAB `any` reference page for more information.

area

Purpose Create a filled area 2-D plot

Description Refer to the MATLAB area reference page for more information.

Purpose Create a vertical bar graph

Description Refer to the MATLAB `bar` reference page for more information.

barh

Purpose Create a horizontal bar graph

Description Refer to the MATLAB `barh` reference page for more information.

Purpose Return the binary representation of the stored integer of a `fi` object as a string

Syntax `bin(a)`

Description Fixed-point numbers can be represented as

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

or, equivalently,

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

The stored integer is the raw binary number, in which the binary point is assumed to be at the far right of the word.

`bin(a)` returns the stored integer of `fi` object `a` in unsigned binary format as a string.

Examples **Example 1**

The following code

```
a = fi([-1 1],1,8,7);  
bin(a)
```

returns

```
10000000 01111111
```

See Also `dec`, `hex`, `int`, `oct`

bin2num

Purpose Convert a two's complement binary string to a number using a quantizer object

Syntax `y = bin2num(q,b)`

Description `y = bin2num(q,b)` uses the properties of quantizer object `q` to convert binary string `b` to numeric array `y`. When `b` is a cell array containing binary strings, `y` is a cell array of the same dimension containing numeric arrays. The fixed-point binary representation is two's complement. The floating-point binary representation is in IEEE Standard 754 style.

`bin2num` and `num2bin` are inverses of one another. Note that `num2bin` always returns the strings in a column.

Examples Create a quantizer object and an array of numeric strings. Convert the numeric strings to binary strings, then use `bin2num` to convert them back to numeric strings.

```
q=quantizer([4 3]);  
[a,b]=range(q);  
x=(b:-eps(q):a)';  
b = num2bin(q,x)
```

```
b =
```

```
0111  
0110  
0101  
0100  
0011  
0010  
0001  
0000  
1111  
1110  
1101
```



```
1100
1011
1010
1001
1000
```

bin2num performs the inverse operation of num2bin.

```
y=bin2num(q,b)
```

```
y =
```

```
0.8750
0.7500
0.6250
0.5000
0.3750
0.2500
0.1250
0
-0.1250
-0.2500
-0.3750
-0.5000
-0.6250
-0.7500
-0.8750
-1.0000
```

See Also

hex2num, num2bin, num2hex, num2int

bitand

Purpose Return the bitwise AND of two `fi` objects

Syntax `c = bitand(a, b)`

Description `c = bitand(a, b)` returns the bitwise AND of `fi` objects `a` and `b`. The `numericType` of `a` and `b` must be identical. If the `numericType` is signed, then the bit representation of the stored integer is in two's complement representation.

See Also `bitcmp`, `bitget`, `bitor`, `bitset`, `bitxor`

Purpose Return the bitwise complement of a `fi` object

Syntax `c = bitcmp(a)`

Description `c = bitcmp(a)` returns the bitwise complement of `fi` object `a` as an `n`-bit nonnegative integer. If `a` has a signed `numericType`, then the bit representation of the stored integer is in two's complement representation.

See Also `bitand`, `bitget`, `bitor`, `bitset`, `bitxor`

bitget

Purpose Return the bit at a certain position

Syntax `c = bitget(a, bit)`

Description `c = bitget(a, bit)` returns the value of the bit at position `bit` in `a`. `a` must be a nonnegative integer, and `bit` must be a number between 1 and the number of bits in the floating-point integer representation of `a`. If `a` has a signed `numericType`, then the bit representation of the stored integer is in two's complement representation.

See Also `bitand`, `bitcmp`, `bitor`, `bitset`, `bitxor`

Purpose Return the bitwise OR of two `fi` objects

Syntax `c = bitor(a, b)`

Description `c = bitor(a, b)` returns the bitwise OR of `fi` objects `a` and `b`. The `numerictype` of `a` and `b` must be identical. If the `numerictype` is signed, then the bit representation of the stored integer is in two's complement representation.

See Also `bitand`, `bitcmp`, `bitget`, `bitset`, `bitxor`

bitset

Purpose Set the bit at a certain position

Syntax
`c = bitset(a, bit)`
`c = bitset(a, bit, v)`

Description `c = bitset(a, bit)` sets bit position `bit` in `a` to 1 (on).
`c = bitset(a, bit, v)` sets bit position `bit` in `a` to `v`. `v` must be either 0 (off) or 1 (on).
`a` must be a nonnegative integer, and `bit` must be a number between 1 and the number of bits in the floating-point integer representation of `a`. If `a` has a signed `numerictype`, then the bit representation of the stored integer is in two's complement representation.

See Also `bitand`, `bitcmp`, `bitget`, `bitor`, `bitxor`

Purpose Shift bits specified number of places

Syntax `c = bitshift(a, k)`

Description `c = bitshift(a, k)` returns the value of `a` shifted by `k` bits.
`a` can be any fixed-point numeric type. The `OverflowMode` and `RoundingMode` properties are obeyed.

See Also `bitand`, `bitcmp`, `bitget`, `bitor`, `bitset`, `bitxor`

bitxor

Purpose Return the bitwise exclusive OR of two `fi` objects

Syntax `c = bitxor(a, b)`

Description `c = bitxor(a, b)` returns the bitwise exclusive OR of `fi` objects `a` and `b`. The `numericType` of `a` and `b` must be identical. If the `numericType` is signed, then the bit representation of the stored integer is in two's complement representation.

See Also `bitand`, `bitcmp`, `bitget`, `bitor`, `bitset`

Purpose

Buffer signal vector into matrix of data frames

Description

Refer to the Signal Processing Toolbox `buffer` reference page for more information.

clabel

Purpose Create contour plot elevation labels

Description Refer to the MATLAB `clabel` reference page for more information.

Purpose Create a 2-D comet plot

Description Refer to the MATLAB comet reference page for more information.

comet3

Purpose Create a 3-D comet plot

Description Refer to the MATLAB `comet3` reference page for more information.

Purpose Plot arrows emanating from the origin

Description Refer to the MATLAB compass reference page for more information.

complex

Purpose Construct a complex `fi` object from real and imaginary parts

Syntax
`c = complex(a,b)`
`c = complex(a)`

Description The `complex` function constructs a complex `fi` object from real and imaginary parts.

`c = complex(a,b)` returns the complex result $a + bi$, where a and b are identically sized real N-D arrays, matrices, or scalars of the same data type. When b is all zero, c is complex with an all-zero imaginary part. This is in contrast to the addition of $a + 0i$, which returns a strictly real result.

`c = complex(a)` for a real `fi` object a returns the complex result $a + bi$ with real part a and an all-zero imaginary part. Even though its imaginary part is all zero, c is complex.

The `numericType` and `fiMath` objects of the leftmost input that is a `fi` object are applied to the output c .

See Also `imag`, `real`

Purpose Plot velocity vectors as cones in a 3-D vector field

Description Refer to the MATLAB coneplot reference page for more information.

conj

Purpose Return the complex conjugate of a fi object

Syntax `conj(a)`

Description `conj(a)` is the complex conjugate of fi object `a`.
When `a` is complex,

$$\text{conj}(a) = \text{real}(a) - i \times \text{imag}(a)$$

The `numericType` and `fiMath` objects of the input `a` are applied to the output.

See Also `complex`, `imag`, `real`

Purpose Create a contour graph of a matrix

Description Refer to the MATLAB contour reference page for more information.

contour3

Purpose Create a 3-D contour plot

Description Refer to the MATLAB contour3 reference page for more information.

Purpose Create a two-level contour plot computation

Description Refer to the MATLAB `contourc` reference page for more information.

contourf

Purpose Create a filled 2-D contour plot

Description Refer to the MATLAB `contourf` reference page for more information.

Purpose Apply convergent rounding

Syntax `convergent(x)`

Description `convergent(x)` rounds the elements of `x` to the nearest integer, except in a tie, then rounds to the nearest even integer.

Examples MATLAB `round` and `convergent` differ in the way they treat values whose fractional part is 0.5. In `round`, every tie is rounded up in absolute value. `convergent` rounds ties to the nearest even integer.

```
x=[ -3.5:3.5]';  
[x convergent(x) round(x)]  
ans =  
  
-3.5000 -4.0000 -4.0000  
-2.5000 -2.0000 -3.0000  
-1.5000 -2.0000 -2.0000  
-0.5000 0 -1.0000  
0.5000 0 1.0000  
1.5000 2.0000 2.0000  
2.5000 2.0000 3.0000  
3.5000 4.0000 4.0000
```

copyobj

Purpose Make an independent copy of a quantizer object

Syntax `q1 = copyobj(q)`
 `[q1,q2,...] = copyobj(obja,objb,...)`

Description `q1 = copyobj(q)` makes a copy of quantizer object `q` and returns it in `q1`.

`[q1,q2,...] = copyobj(obja,objb,...)` copies `obja` into `q1`, `objb` into `q2`, and so on.

Using `copyobj` to copy a quantizer object is not the same as using the command syntax `q1 = q` to copy a quantizer object. quantizer objects have memory (their read-only properties). When you use `copyobj`, the resulting copy is independent of the original item—it does not share the original object’s memory, such as the values of the properties `min`, `max`, `noverflows`, or `noperations`. Using `q1 = q` creates a new object that is an alias for the original and shares the original object’s memory, and thus its property values.

Examples `q = quantizer('CoefficientFormat',[8 7]);`
 `q1 = copyobj(q);`

See Also `quantizer`, `get`, `set`

Purpose Return the complex conjugate transpose of a `fi` object

Syntax `ctranspose(a)`

Description `ctranspose(a)` returns the complex conjugate transpose of `fi` object `a`. It is also called for the syntax `'`.

See Also `transpose`

dec

Purpose Return the unsigned decimal representation of the stored integer of a `fi` object as a string

Syntax `dec(a)`

Description Fixed-point numbers can be represented as

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

or, equivalently,

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

The stored integer is the raw binary number, in which the binary point is assumed to be at the far right of the word.

`dec(a)` returns the stored integer of `fi` object `a` in unsigned decimal format as a string.

Examples

Example 1

The code

```
a = fi([-1 1],1,8,7);  
dec(a)
```

returns

```
128 127
```

See Also

`bin`, `hex`, `int`, `oct`, `sdec`

Purpose Return the largest denormalized quantized number for a quantizer object

Syntax `x = denormalmax(q)`

Description `x = denormalmax(q)` is the largest positive denormalized quantized number where `q` is a quantizer object. Anything larger than `x` is a normalized number. Denormalized numbers apply only to floating-point format. When `q` represents fixed-point numbers, this function returns `eps(q)`.

Examples

```
q = quantizer('float',[6 3]);
x = denormalmax(q)

x =

    0.1875
```

Algorithm When `q` is a floating-point quantizer object,

$$\text{denormalmax}(q) = \text{realmin}(q) - \text{denormalmin}(q)$$

When `q` is a fixed-point quantizer object,

$$\text{denormalmax}(q) = \text{eps}(q)$$

See Also `denormalmin`, `eps`, `quantizer`

denormalmin

Purpose Return the smallest denormalized quantized number for a quantizer object

Syntax `x = denormalmin(q)`

Description `x = denormalmin(q)` is the smallest positive denormalized quantized number where `q` is a quantizer object. Anything smaller than `x` underflows to zero with respect to the quantizer object `q`. Denormalized numbers apply only to floating-point format. When `q` represents a fixed-point number, `denormalmin` returns `eps(q)`.

Examples

```
q = quantizer('float',[6 3]);
denormalmin(q)

ans =

    0.0625
```

Algorithm When `q` is a floating-point quantizer object,

$$x = 2^{E_{\min} - f}$$

where E_{\min} is equal to `exponentmin(q)`.

When `q` is a fixed-point quantizer object,

$$x = \text{eps}(q) = 2^{-f}$$

where f is equal to `fractionlength(q)`.

See Also `denormalmax`, `eps`, `quantizer`

Purpose Return diagonal matrices or the diagonals of a matrix

Description Refer to the MATLAB `diag` reference page for more information.

disp

Purpose Display an object

Description Refer to the MATLAB `disp` reference page for more information.

Purpose Divide two objects using a numeric type object

Syntax `c = T.divide(a,b)`

Description `c = T.divide(a,b)` performs division on the elements of `a` by the elements of `b` using numeric type object `T`.

`a` and `b` must have the same dimensions unless one is a scalar. If either `a` or `b` is scalar, then `c` has the dimensions of the nonscalar object.

If either `a` or `b` is a `fi` object, and the other is a MATLAB built-in numeric type, then the built-in object is cast to the word length of the `fi` object, preserving best-precision fraction length.

If `a` and `b` are both MATLAB built-in doubles, then `c` is the double-precision quotient `a./b`, and numeric type `T` is ignored.

Examples This example highlights the precision of the `fi divide` function.

First, create an unsigned `fi` object with an 80-bit word length and 2^{83} scaling, which puts the leading 1 of the representation into the most significant bit. Initialize the object with double-precision floating-point value 0.1, and examine the binary representation:

```
P =
fipref('NumberDisplay','bin','NumericTypeDisplay','short',...
'FimathDisplay','none');
a = fi(0.1, false, 80, 83)

a =

11001100110011001100110011001100110011001100110011001100110011001100110011010000000000000
000000000000000000
(bin)
u80,83
11001100110011001100110011001100110011001100110011001100110011001100110011001100110011001100
1100110011001100
```

divide

Notice that the infinite repeating representation is truncated after 52 bits, because the mantissa of an IEEE standard double-precision floating-point number has 52 bits.

Contrast the above to calculating $1/10$ in fixed-point arithmetic with the quotient set to the same numeric type as before:

```
T = numerictype('Signed',false,'WordLength',80,...
               'FractionLength',83);
```

```
a = fi(1);
b = fi(10);
c = T.divide(a,b);
c.bin
```

```
ans =
```

```
1100110011001100110011001100110011001100110011001100110011001100110011001100110011001100
1100110011001100
```

Notice that when you use the `divide` function, the quotient is calculated to the full 80 bits, regardless of the precision of `a` and `b`. Thus, the `fi` object `c` represents $1/10$ more precisely than IEEE standard double-precision floating-point number can.

With 1000 bits of precision,

```
T = numerictype('Signed',false,'WordLength',1000,...
               'FractionLength',1003);
```

```
a = fi(1);
b = fi(10);
c = T.divide(a,b);
c.bin
```

```
ans =
```

```
1100110011001100110011001100110011001100110011001100110011001100110011001100110011001100
1100110011001100110011001100110011001100110011001100110011001100110011001100110011001100
1100110011001100110011001100110011001100110011001100110011001100110011001100110011001100
```


double

Purpose Return the double-precision floating-point real-world value of a `fi` object

Syntax `double(a)`

Description Fixed-point numbers can be represented as

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

or, equivalently,

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

`double(a)` returns the real-world value of a `fi` object in double-precision floating point.

See Also `single`

Purpose Indicate last index of array

Description Refer to the MATLAB end reference page for more information.

eps

Purpose Return the quantized relative accuracy for `fi` objects or quantizer objects

Syntax `eps(obj)`

Description `eps(obj)` returns the value of the least significant bit of the value of the `fi` object or quantizer object `obj`. The result of this function is equivalent to that given by the Fixed-Point Toolbox `lsb` function.

See Also `lsb`

Purpose	Determine whether the real-world values of two <code>fi</code> objects are equal
Syntax	<code>c = eq(a,b)</code> <code>a == b</code>
Description	<code>c = eq(a,b)</code> is called for the syntax <code>'a == b'</code> when <code>a</code> or <code>b</code> is a <code>fi</code> object. <code>a</code> and <code>b</code> must have the same dimensions unless one is a scalar. A scalar can be compared with another object of any size. <code>a == b</code> does an element-by-element comparison between <code>a</code> and <code>b</code> and returns a matrix of the same size with elements set to 1 where the relation is true, and 0 where the relation is false.
See Also	<code>ge</code> , <code>gt</code> , <code>isequal</code> , <code>le</code> , <code>lt</code> , <code>ne</code>

errorbar

Purpose Plot error bars along a curve

Description Refer to the MATLAB errorbar reference page for more information.

Purpose Plot elimination tree

Description Refer to the MATLAB `etreeplot` reference page for more information.

exponentbias

Purpose Return the exponent bias for a quantizer object

Syntax `b = exponentbias(q)`

Description `b = exponentbias(q)` returns the exponent bias of the quantizer object `q`. For fixed-point quantizer objects, `exponentbias(q)` returns 0.

Examples

```
q = quantizer('double');
b = exponentbias(q)

b =

    1023
```

Algorithm For floating-point quantizer objects,

$$b = 2^{e-1} - 1$$

where $e = \text{eps}(q)$, and `exponentbias` is the same as the exponent maximum.

For fixed-point quantizer objects, $b = 0$ by definition.

See Also `eps`, `exponentlength`, `exponentmax`, `exponentmin`

Purpose Return the exponent length of a quantizer object

Syntax `e = exponentlength(q)`

Description `e = exponentlength(q)` returns the exponent length of quantizer object `q`. When `q` is a fixed-point quantizer object, `exponentlength(q)` returns 0. This is useful because exponent length is valid whether the quantizer object mode is floating point or fixed point.

Examples

```
q = quantizer('double');  
e = exponentlength(q)  
  
e =  
  
    11
```

Algorithm The exponent length is part of the format of a floating-point quantizer object `[w e]`. For fixed-point quantizer objects, `e = 0` by definition.

See Also `eps`, `exponentbias`, `exponentmax`, `exponentmin`

exponentmax

Purpose Return the maximum exponent for a quantizer object

Syntax `exponentmax(q)`

Description `exponentmax(q)` returns the maximum exponent for quantizer object `q`. When `q` is a fixed-point quantizer object, it returns 0.

Examples

```
q = quantizer('double');
exponentmax(q)

ans =

    1023
```

Algorithm For floating-point quantizer objects,

$$E_{max} = 2^{e-1} - 1$$

For fixed-point quantizer objects, $E_{max} = 0$ by definition.

See Also `eps`, `exponentbias`, `exponentlength`, `exponentmin`

Purpose Return the minimum exponent for a quantizer object

Syntax `emin = exponentmin(q)`

Description `emin = exponentmin(q)` returns the minimum exponent for quantizer object `q`. If `q` is a fixed-point quantizer object, `exponentmin` returns 0.

Examples

```
q = quantizer('double');
emin = exponentmin(q)

emin =

    -1022
```

Algorithm For floating-point quantizer objects,

$$E_{min} = -2^{e-1} + 2$$

For fixed-point quantizer objects, $E_{min} = 0$.

See Also `eps`, `exponentbias`, `exponentlength`, `exponentmax`

ezcontour

Purpose Easy-to-use contour plotter

Description Refer to the MATLAB ezcontour reference page for more information.

Purpose Easy-to-use filled contour plotter

Description Refer to the MATLAB ezcontourf reference page for more information.

ezmesh

Purpose Easy-to-use 3-D mesh plotter

Description Refer to the MATLAB ezmesh reference page for more information.

Purpose Easy-to-use function plotter

Description Refer to the MATLAB ezplot reference page for more information.

ezplot3

Purpose Easy-to-use 3-D parametric curve plotter

Description Refer to the MATLAB ezplot3 reference page for more information.

Purpose Easy-to-use polar coordinate plotter

Description Refer to the MATLAB `ezpolar` reference page for more information.

ezsurf

Purpose Easy-to-use 3-D colored surface plotter

Description Refer to the MATLAB ezsurf reference page for more information.

Purpose Easy-to-use combination surface/contour plotter

Description Refer to the MATLAB ezsurf reference page for more information.

feather

Purpose Plot velocity vectors

Description Refer to the MATLAB feather reference page for more information.

Purpose

Construct a `fi` object

Syntax

```
a = fi(v)
a = fi(v,s)
a = fi(v,s,w)
a = fi(v,s,w,f)
a = fi(v,s,w,slope,bias)
a = fi(v,s,w,slopeadjustmentfactor,fixedexponent,bias)
a = fi(v,T)
a = fi(v,T,F)
a = fi(...'PropertyName',PropertyValue...)
fi('PropertyName',PropertyValue...)
```

Description

You can use the `fi` constructor function in the following ways.

- `a = fi(v)` returns a signed fixed-point object with value `v`, 16-bit word length, and best-precision fraction length.
- `a = fi(v,s)` returns a fixed-point object with value `v`, signedness `s`, 16-bit word length, and best-precision fraction length. `s` can be 0 (false) for unsigned or 1 (true) for signed.
- `a = fi(v,s,w)` returns a fixed-point object with value `v`, signedness `s`, word length `w`, and best-precision fraction length.
- `a = fi(v,s,w,f)` returns a fixed-point object with value `v`, signedness `s`, word length `w`, and fraction length `f`.
- `a = fi(v,s,w,slope,bias)` returns a fixed-point object with value `v`, signedness `s`, word length `w`, slope, and bias.
- `a = fi(v,s,w,slopeadjustmentfactor,fixedexponent,bias)` returns a fixed-point object with value `v`, signedness `s`, word length `w`, slopeadjustmentfactor, fixedexponent, and bias.
- `a = fi(v,T)` returns a fixed-point object with value `v` and embedded.numericity `T`. Refer to for more information on numericity objects.

- `a = fi(v,T,F)` returns a fixed-point object with value `v`, embedded `numericType T`, and embedded `fimath F`. Refer to for more information on `fimath` objects.
- `a = fi(...'PropertyName',PropertyValue...)` and `fi('PropertyName',PropertyValue...)` allow you to set fixed-point objects for a `fi` object by property name/property value pairs.

The `fi` object has the following three general types of properties.

Note These properties are described in detail in “`fi` Object Properties” on page 9-2 in the Properties Reference.

- “Data Properties” on page 11-66
- “Fimath Properties” on page 11-67
- “NumericType Properties” on page 11-68

Data Properties

The data properties of a `fi` object are always writable.

- `bin` – Stored integer value of a `fi` object in binary
- `data` – Numerical real-world value of a `fi` object
- `dec` – Stored integer value of a `fi` object in decimal
- `double` – Real-world value of a `fi` object, stored as a MATLAB `double`
- `hex` – Stored integer value of a `fi` object in hexadecimal
- `int` – Stored integer value of a `fi` object, stored in a built-in MATLAB integer data type. You can also use `int8`, `int16`, `int32`, `uint8`, `uint16`, and `uint32` to get the stored integer value of a `fi` object in these formats
- `oct` – Stored integer value of a `fi` object in octal

These properties are described in detail in “fi Object Properties” on page 9-2 in the Properties Reference.

Fimath Properties

When you create a `fi` object, a `fimath` object is also automatically created as a property of the `fi` object.

- `fimath` – `fimath` object associated with a `fi` object

The following `fimath` properties are, by transitivity, also properties of a `fi` object. The properties of the `fimath` object listed below are always writable.

- `CastBeforeSum` – Whether both operands are cast to the sum data type before addition
- `MaxProductWordLength` – Maximum allowable word length for the product data type
- `MaxSumWordLength` – Maximum allowable word length for the sum data type
- `ProductFractionLength` – Fraction length, in bits, of the product data type
- `ProductMode` – Defines how the product data type is determined
- `ProductWordLength` – Word length, in bits, of the product data type
- `RoundMode` – Rounding mode
- `SumFractionLength` – Fraction length, in bits, of the sum data type
- `SumMode` – Defines how the sum data type is determined
- `SumWordLength` – Word length, in bits, of the sum data type

These properties are described in detail in “fi Object Properties” on page 9-2 in the Properties Reference.

Numerictype Properties

When you create a `fi` object, a `numerictype` object is also automatically created as a property of the `fi` object.

- `numerictype` – Object containing all the numeric type attributes of a `fi` object

The following `numerictype` properties are, by transitivity, also properties of a `fi` object. The properties of the `numerictype` object listed below are not writable once the `fi` object has been created. However, you can create a copy of a `fi` object with new values specified for the `numerictype` properties.

- `Bias` – Bias of a `fi` object
- `DataType` – Data type category associated with a `fi` object
- `DataTypeMode` – Data type and scaling mode of a `fi` object
- `FixedExponent` – Fixed-point exponent associated with a `fi` object
- `SlopeAdjustmentFactor` – Slope adjustment associated with a `fi` object
- `FractionLength` – Fraction length of the stored integer value of a `fi` object in bits
- `Scaling` – Fixed-point scaling mode of a `fi` object
- `Signed` – Whether a `fi` object is signed or unsigned
- `Slope` – Slope associated with a `fi` object
- `WordLength` – Word length of the stored integer value of a `fi` object in bits

These properties are described in detail in “`fi` Object Properties” on page 9-2 in the Properties Reference.

Examples

Note For information on the display format of fi objects, refer to “Display Settings” on page 1-5.

Example 1

For example, the following creates a fi object with a value of pi, a word length of 8 bits, and a fraction length of 3 bits.

```
a = fi(pi, 1, 8, 3)
```

```
a =
```

```
3.1250
```

```
DataTypeMode: Fixed-point: binary point scaling
Signed: true
WordLength: 8
FractionLength: 3
```

Example 2

The value v can also be an array.

```
a = fi((magic(3)/10), 1, 16, 12)
```

```
a =
```

```
0.8000    0.1001    0.6001
0.3000    0.5000    0.7000
0.3999    0.8999    0.2000
```

```
DataTypeMode: Fixed-point: binary point scaling
Signed: true
WordLength: 16
```

FractionLength: 12

Example 3

If you omit the argument *f*, it is set automatically to the best precision possible.

```
a = fi(pi, 1, 8)
```

```
a =
```

```
3.1563
```

```
DataTypeMode: Fixed-point: binary point scaling  
Signed: true  
WordLength: 8  
FractionLength: 5
```

Example 4

If you omit *w* and *f*, they are set automatically to 16 bits and the best precision possible, respectively.

```
a = fi(pi, 1)
```

```
a =
```

```
3.1416
```

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


Example 5

You can use property name/property value pairs to set `fi` properties when you create the object.

```
a = fi(pi, 'roundmode', 'floor', 'overflowmode', 'wrap')
```

```
a =
```

```
3.1415
```

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

See Also

`fimath`, `fipref`, `numerictype`, `quantizer`, “`fi` Object Properties” on page 9-2

fimath

Purpose Construct a fimath object

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

Description You can use the `fimath` constructor function in the following ways:

- `F = fimath` creates a default `fimath` object.
- `F = fimath(...'PropertyName',PropertyValue...)` allows you to set the attributes of a `fimath` object using property name/property value pairs.

The properties of the `fimath` object are listed below. These properties are described in detail in “`fimath` Object Properties” on page 9-5 in the Properties Reference.

- `CastBeforeSum` – Whether both operands are cast to the sum data type before addition
- `MaxProductWordLength` – Maximum allowable word length for the product data type
- `MaxSumWordLength` – Maximum allowable word length for the sum data type
- `OverflowMode` – Overflow-handling mode
- `ProductFractionLength` – Fraction length, in bits, of the product data type
- `ProductMode` – Defines how the product data type is determined
- `ProductWordLength` – Word length, in bits, of the product data type
- `RoundMode` – Rounding mode
- `SumFractionLength` – Fraction length, in bits, of the sum data type
- `SumMode` – Defines how the sum data type is determined
- `SumWordLength` – Word length, in bits, of the sum data type

Examples

Example 1

Type

```
F = fmath
```

to create a default fmath object.

```
F = fmath
```

```
F =
```

```
          RoundMode: round
    OverflowMode: saturate
      ProductMode: FullPrecision
MaxProductWordLength: 128
          SumMode: FullPrecision
    MaxSumWordLength: 128
      CastBeforeSum: true
```

Example 2

You can set properties of fmath objects at the time of object creation by including properties after the arguments of the fmath constructor function. For example, to set the overflow mode to saturate and the rounding mode to convergent,

```
F = fmath('OverflowMode','saturate','RoundMode','convergent')
```

```
F =
```

```
          RoundMode: convergent
    OverflowMode: saturate
      ProductMode: FullPrecision
MaxProductWordLength: 128
          SumMode: FullPrecision
    MaxSumWordLength: 128
```

CastBeforeSum: true

See Also

fi, fipref, numerictype, quantizer, “fimath Object Properties” on page 9-5

Purpose Find indices and values of nonzero elements

Description Refer to the MATLAB `find` reference page for more information.

fipref

Purpose Construct a fipref object

Syntax
`P = fipref`
`P = fipref(...'PropertyName',PropertyValue...)`

Description You can use the fipref constructor function in the following ways:

- `P = fipref` creates a default fipref object.
- `P = fipref(...'PropertyName',PropertyValue...)` allows you to set the attributes of a object using property name/property value pairs.

The properties of the fipref object are listed below. These properties are described in detail in “fipref Object Properties” on page 9-10.

- `FimathDisplay` – Display options for the fimath attributes of a fi object
- `NumericTypeDisplay` – Display options for the numeric type attributes of a fi object
- `NumberDisplay` – Display options for the value of a fi object
- `LoggingMode` – Logging options for operations performed on fi objects

Your fipref settings persist throughout your MATLAB session. Use `reset(fipref)` to return to the default settings during your session. Use `savefipref` to save your display preferences for subsequent MATLAB sessions.

Examples

Example 1

Type

```
P = fipref
```

to create a default fipref object.

```
P =  
  
    NumberDisplay: 'RealWorldValue'  
    NumericTypeDisplay: 'full'  
    FimathDisplay: 'full'  
    LoggingMode: 'Off'
```

Example 2

You can set properties of fipref objects at the time of object creation by including properties after the arguments of the fipref constructor function. For example, to set NumberDisplay to bin and AttributesDisplay to short,

```
P = fipref('NumberDisplay', 'bin', 'NumericType', 'short')
```

```
P =  
  
    NumberDisplay: 'bin'  
    NumericTypeDisplay: 'short'  
    FimathDisplay: 'full'  
    LoggingMode: 'Off'
```

See Also

fi, fimath, numerictype, quantizer, savefipref, “fipref Object Properties” on page 9-10

fplot

Purpose Plot a function between specified limits

Description Refer to the MATLAB `fplot` reference page for more information.

Purpose Return the fraction length of a quantizer object

Syntax `fractionlength(q)`

Description `fractionlength(q)` returns the fraction length of quantizer object `q`.

Examples For a floating-point quantizer object,

```
q = quantizer('float',[32 8]);  
f = fractionlength(q)
```

```
f =
```

```
23
```

where $f = 23 = 32 - 8 - 1$.

For a fixed-point quantizer object,

```
q = quantizer('fixed',[6 4])  
f = fractionlength(q)
```

```
q =
```

```
    DataMode = fixed  
    RoundMode = floor  
    OverflowMode = saturate  
    Format = [6 4]
```

```
    Max = reset  
    Min = reset  
    NOverflows = 0  
    NUnderflows = 0  
    NOperations = 0
```

```
f =
```

fractionlength

4

Algorithm For floating-point quantizer objects, $f = w - e - 1$, where w is the word length and e is the exponent length.

For fixed-point quantizer objects, f is part of the format $[w f]$.

See Also `fi`, `numerictype`, `quantizer`, `wordlength`

Purpose	Determine whether the real-world value of one <code>fi</code> object is greater than or equal to another
Syntax	<code>c = ge(a,b)</code> <code>a >= b</code>
Description	<code>c = ge(a,b)</code> is called for the syntax ' <code>a >= b</code> ' when <code>a</code> or <code>b</code> is a <code>fi</code> object. <code>a</code> and <code>b</code> must have the same dimensions unless one is a scalar. A scalar can be compared with another object of any size. <code>a >= b</code> does an element-by-element comparison between <code>a</code> and <code>b</code> and returns a matrix of the same size with elements set to 1 where the relation is true, and 0 where the relation is false.
See Also	<code>eq</code> , <code>gt</code> , <code>le</code> , <code>lt</code> , <code>ne</code>

get

Purpose Return the property values of a quantizer object

Syntax

```
get(q,pn,pv)
value = get(q, 'propertyname')
structure = get(q)
```

Description `get(q,pn,pv)` displays the property names and property values associated with quantizer object `q`.

`pn` is the name of a property of the object `obj`, and `pv` is the value associated with `pn`.

`value = get(q, 'propertyname')` returns the property value associated with the property named in the string `'propertyname'` for the quantizer object `q`. If you replace the string `'propertyname'` by a cell array of a vector of strings containing property names, `get` returns a cell array of a vector of corresponding values.

`structure = get(q)` returns a structure containing the properties and states of quantizer object `q`.

See Also `quantizer`, `set`

Purpose Plot set of nodes using an adjacency matrix

Description Refer to the MATLAB `gplot` reference page for more information.

Purpose Determine whether the real-world value of one `fi` object is greater than another

Syntax
`c = gt(a,b)`
`a > b`

Description `c = gt(a,b)` is called for the syntax '`a > b`' when `a` or `b` is a `fi` object. `a` and `b` must have the same dimensions unless one is a scalar. A scalar can be compared with another object of any size.

`a > b` does an element-by-element comparison between `a` and `b` and returns a matrix of the same size with elements set to 1 where the relation is true, and 0 where the relation is false.

See Also `eq`, `ge`, `le`, `lt`, `ne`

Purpose Return a Hankel matrix

Description Refer to the MATLAB `hankel` reference page for more information.

hex

Purpose Return the hexadecimal representation of the stored integer of a `fi` object as a string

Syntax `hexadecimal(a)`

Description Fixed-point numbers can be represented as

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

or, equivalently,

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

The stored integer is the raw binary number, in which the binary point is assumed to be at the far right of the word.

`hexadecimal(a)` returns the stored integer of `fi` object `a` in hexadecimal format as a string.

Examples

Example 1

The following code

```
a = fi([-1 1],1,8,7);  
hex(a)
```

returns

```
80 7f
```

See Also

`bin`, `dec`, `int`, `oct`

Purpose Convert a hexadecimal string to a number using a quantizer object

Syntax

```
x = hex2num(q,h)
[x1,x2,...] = hex2num(q,h1,h2,...)
```

Description `x = hex2num(q,h)` converts hexadecimal string `h` to numeric matrix `x`. The attributes of the numbers in `x` are specified by quantizer object `q`. When `h` is a cell array containing hexadecimal strings, `hex2num` returns `x` as a cell array of the same dimension containing numbers. For fixed-point hexadecimal strings, `hex2num` uses two's complement representation. For floating-point strings, the representation is IEEE Standard 754 style.

When there are fewer hexadecimal digits than needed to represent the number, the fixed-point conversion zero-fills on the left. Floating-point conversion zero-fills on the right.

`[x1,x2,...] = hex2num(q,h1,h2,...)` converts hexadecimal strings `h1, h2,...` to numeric matrices `x1, x2,...`

`hex2num` and `num2hex` are inverses of one another, with the distinction that `num2hex` returns the hexadecimal strings in a column.

Examples To create all the 4-bit fixed-point two's complement numbers in fractional form, use the following code.

```
q = quantizer([4 3]);
h = ['7 3 F B'; '6 2 E A'; '5 1 D 9'; '4 0 C 8'];
x = hex2num(q,h)
```

```
x =

    0.8750    0.3750   -0.1250   -0.6250
    0.7500    0.2500   -0.2500   -0.7500
    0.6250    0.1250   -0.3750   -0.8750
    0.5000         0   -0.5000   -1.0000
```

See Also `bin2num`, `num2bin`, `num2hex`, `num2int`

hist

Purpose Create histogram plot

Description Refer to the MATLAB `hist` reference page for more information.

Purpose Return histogram count

Description Refer to the MATLAB `histc` reference page for more information.

horzcat

Purpose Horizontally concatenate two or more `fi` objects

Syntax `c = horzcat(a,b,...)`
`[a, b, ...]`

Description `c = horzcat(a,b,...)` is called for the syntax `[a, b, ...]` when any of `a, b, ...`, is a `fi` object.

`[a b, ...]` or `[a,b, ...]` is the horizontal concatenation of matrices `a` and `b`. `a` and `b` must have the same number of rows. Any number of matrices can be concatenated within one pair of brackets. N-D arrays are horizontally concatenated along the second dimension. The first and remaining dimensions must match.

Horizontal and vertical concatenation can be combined together as in `[1 2;3 4]`.

`[a b; c]` is allowed if the number of rows of `a` equals the number of rows of `b`, and if the number of columns of `a` plus the number of columns of `b` equals the number of columns of `c`.

The matrices in a concatenation expression can themselves be formed via a concatenation as in `[a b;[c d]]`.

Note The `fi`math and `numeric`type objects of a concatenated matrix of `fi` objects `c` are taken from the leftmost `fi` object in the list `(a,b,...)`

See Also `vertcat`

Purpose Return the imaginary part

Description Refer to the MATLAB `imag` reference page for more information.

innerprodintbits

Purpose Return the number of integer bits needed for a fixed-point inner product

Syntax `innerprodintbits(a,b)`

Description `innerprodintbits(a,b)` computes the minimum number of integer bits necessary in the inner product of $a' * b$ to guarantee that no overflows occur and to preserve best precision.

- a and b are `fi` vectors
- The values of a are known
- Only the numeric type of b is relevant. The values of b are ignored

Examples The primary use of this function is to determine the number of integer bits necessary in the output Y of an FIR filter that computes the inner product between constant coefficient row vector B and state column vector Z . For example,

```
for k=1:length(X);
    Z = [X(k);Z(1:end-1)];
    Y(k) = B * Z;
end
```

Algorithm In general, an inner product grows $\log_2(n)$ bits for vectors of length n . However, in the case of this function the vector a is known and its values do not change. This knowledge is used to compute the smallest number of integer bits that are necessary in the output to guarantee that no overflow will occur.

The largest gain occurs when the vector b has the same sign as the constant vector a . Therefore, the largest gain due to the vector a is $a * \text{sign}(a')$, which is equal to $\text{sum}(\text{abs}(a))$.

The overall number of integer bits necessary to guarantee that no overflow occurs in the inner product is computed by:

$$\log_2(\text{sum}(\text{abs}(a)) + \text{number of integer bits in } b + 1 \text{ sign bit})$$

Purpose Display Property Inspector

Description Refer to the MATLAB `inspect` reference page for more information.

int

Purpose Return the smallest built-in integer in which the stored integer value of a `fi` object will fit

Syntax `int(a)`

Description Fixed-point numbers can be represented as

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

or, equivalently,

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

The stored integer is the raw binary number, in which the binary point is assumed to be at the far right of the word.

`int(a)` returns the smallest built-in integer of the data type in which the stored integer value of `fi` object `a` will fit.

The following table gives the return type of the `int` function.

Word Length	Return Type for Signed <code>fi</code>	Return Type for Unsigned <code>fi</code>
word length <= 8 bits	<code>int8</code>	<code>uint8</code>
8 bits < word length <= 16 bits	<code>int16</code>	<code>uint16</code>
16 bits < word length <= 32 bits	<code>int32</code>	<code>uint32</code>
32 < word length	<code>double</code>	<code>double</code>

Note When the word length is greater than 52 bits, the return value can have quantization error. For bit-true integer representation of very large word lengths, use `bin`, `oct`, `dec`, `hex`, or `sdec`.

See Also

`int8`, `int16`, `int32`, `uint8`, `uint16`, `uint32`

int8

Purpose Return the stored integer value of a `fi` object as a built-in `int8`

Syntax `int8(a)`

Description Fixed-point numbers can be represented as

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

or, equivalently,

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

The stored integer is the raw binary number, in which the binary point is assumed to be at the far right of the word.

`int8(a)` returns the stored integer value of `fi` object `a` as a built-in `int8`. If the stored integer word length is too big for an `int8`, or if the stored integer is unsigned, the returned value saturates to an `int8`.

See Also `int`, `int16`, `int32`, `uint8`, `uint16`, `uint32`

Purpose Return the stored integer value of a `fi` object as a built-in `int16`

Syntax `int16(a)`

Description Fixed-point numbers can be represented as

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

or, equivalently,

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

The stored integer is the raw binary number, in which the binary point is assumed to be at the far right of the word.

`int16(a)` returns the stored integer value of `fi` object `a` as a built-in `int16`. If the stored integer word length is too big for an `int16`, or if the stored integer is unsigned, the returned value saturates to an `int16`.

See Also `int`, `int8`, `int32`, `uint8`, `uint16`, `uint32`

int32

Purpose Return the stored integer value of a `fi` object as a built-in `int32`

Syntax `int32(a)`

Description Fixed-point numbers can be represented as

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

or, equivalently,

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

The stored integer is the raw binary number, in which the binary point is assumed to be at the far right of the word.

`int32(a)` returns the stored integer value of `fi` object `a` as a built-in `int32`. If the stored integer word length is too big for an `int32`, or if the stored integer is unsigned, the returned value saturates to an `int32`.

See Also `int`, `int8`, `int16`, `uint8`, `uint16`, `uint32`

Purpose Return the largest positive stored integer value representable by the `numerictype` of a `fi` object

Syntax `x = intmax(a)`

Description `x = intmax(a)` returns the largest positive stored integer value representable by the `numerictype` of `a`.

See Also `intmin`, `lsb`, `stripscaling`

intmin

Purpose Return smallest stored integer value representable by numeric type of `fi` object

Syntax `x = intmin(a)`

Description `x = intmin(a)` returns the smallest stored integer value representable by the numeric type of `a`.

Examples

```
a = fi(pi, true, 16, 12);
x = intmin(a)
```

```
x =
```

```
-32768
```

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

See Also `intmax`, `lsb`, `stripScaling`

Purpose Inverse permute the dimensions of a multidimensional array

Description Refer to the MATLAB `ipermute` reference page for more information.

iscolumn

Purpose	Determine whether a <code>fi</code> object is a column vector
Syntax	<code>iscolumn(a)</code>
Description	<code>iscolumn(a)</code> returns 1 if the <code>fi</code> object <code>a</code> is a column vector, and 0 otherwise.
See Also	<code>isrow</code>

Purpose Determine if array is empty

Description Refer to the MATLAB `isempty` reference page for more information.

isequal

Purpose Determine whether the real-world values of two `fi` objects are equal, or determine whether the properties of two `fimath`, `numericType`, or `quantizer` objects are equal

Syntax

```
isequal(a,b,...)  
isequal(F,G,...)  
isequal(T,U,...)  
isequal(q,r,...)
```

Description

`isequal(a,b,...)` returns 1 if all the `fi` object inputs have the same real-world value. Otherwise, the function returns 0.

`isequal(F,G,...)` returns 1 if all the `fimath` object inputs have the same properties. Otherwise, the function returns 0.

`isequal(T,U,...)` returns 1 if all the `numericType` object inputs have the same properties. Otherwise, the function returns 0.

`isequal(q,r,...)` returns 1 if all the `quantizer` object inputs have the same properties. Otherwise, the function returns 0.

See Also `eq`, `ispropequal`

Purpose	Determine whether a variable is a <code>fi</code> object
Syntax	<code>isfi(a)</code>
Description	<code>isfi(a)</code> returns 1 if <code>a</code> is a <code>fi</code> object, and 0 otherwise.
See Also	<code>fi</code> , <code>isfimath</code> , <code>isnumericitype</code>

isfimath

Purpose	Determine whether a variable is a fimath object
Syntax	<code>isfimath(F)</code>
Description	<code>isfimath(F)</code> returns 1 if F is a fimath object, and 0 otherwise.
See Also	<code>fimath</code> , <code>isfi</code> , <code>isnumericitype</code>

Purpose Determine if input is numeric array

Description Refer to the MATLAB `isnumeric` reference page for more information.

isnumerictype

Purpose	Determine whether a variable is a numerictype object
Syntax	<code>isnumerictype(T)</code>
Description	<code>isnumerictype(T)</code> returns 1 if a is a numerictype object, and 0 otherwise.
See Also	<code>isfi</code> , <code>isfimath</code> , <code>numerictype</code>

Purpose Determine if input is MATLAB OOPS object

Description Refer to the MATLAB `isobject` reference page for more information.

ispropequal

Purpose	Determine whether the properties of two <code>fi</code> objects are equal
Syntax	<code>ispropequal(a,b,...)</code>
Description	<code>ispropequal(a,b,...)</code> returns 1 if all the inputs are <code>fi</code> objects and all the inputs have the same properties. Otherwise, the function returns 0. To compare the real-world values of two <code>fi</code> objects <code>a</code> and <code>b</code> , use <code>a == b</code> or <code>isequal(a,b)</code> .
See Also	<code>fi</code> , <code>isequal</code>

Purpose Determine if all array elements are real numbers

Description Refer to the MATLAB `isreal` reference page for more information.

isrow

Purpose Determine whether a `fi` object is a row vector

Syntax `isrow(a)`

Description `isrow(a)` returns 1 if the `fi` object `a` is a row vector, and 0 otherwise.

See Also `iscolumn`

Purpose Determine if input is scalar

Description Refer to the MATLAB `isscalar` reference page for more information.

issigned

Purpose	Determine whether a <code>fi</code> object is signed
Syntax	<code>issigned(a)</code>
Description	<code>issigned(a)</code> returns 1 if the <code>fi</code> object <code>a</code> is signed, and 0 if it is unsigned.

Purpose Determine if input is vector

Description Refer to the MATLAB `isvector` reference page for more information.

le

Purpose Determine whether the real-world value of a `fi` object is less than or equal to another

Syntax
`c = le(a,b)`
`a <= b`

Description `c = le(a,b)` is called for the syntax '`a <= b`' when `a` or `b` is a `fi` object. `a` and `b` must have the same dimensions unless one is a scalar. A scalar can be compared with another object of any size.

`a <= b` does an element-by-element comparison between `a` and `b` and returns a matrix of the same size with elements set to 1 where the relation is true, and 0 where the relation is false.

See Also `eq`, `ge`, `gt`, `lt`, `ne`

Purpose Return the length of a vector

Description Refer to the MATLAB length reference page for more information.

line

Purpose Create line object

Description Refer to the MATLAB `line` reference page for more information.

Purpose Convert numeric values to logical

Description Refer to the MATLAB `logical` reference page for more information.

loglog

Purpose Create log-log scale plot

Description Refer to the MATLAB loglog reference page for more information.

Purpose Return lower bound of range of f i object

Syntax lowerbound(a)

Description lowerbound(a) returns the lower bound of the range of f i object a. If $L = \text{lowerbound}(a)$ and $U = \text{upperbound}(a)$, then $[L, U] = \text{range}(a)$.

See Also range, upperbound

lsb

Purpose	Return the scaling of the least significant bit of a <code>fi</code> object
Syntax	<code>lsb(a)</code>
Description	<code>lsb(a)</code> returns the scaling of the least significant bit of <code>fi</code> object <code>a</code> . The result is equivalent to the result given by the <code>eps</code> function.
See Also	<code>eps</code>

Purpose	Determine whether the real-world value of a <code>fi</code> object is less than another
Syntax	<code>c = lt(a,b)</code> <code>a < b</code>
Description	<code>c = lt(a,b)</code> is called for the syntax ' <code>a < b</code> ' when <code>a</code> or <code>b</code> is a <code>fi</code> object. <code>a</code> and <code>b</code> must have the same dimensions unless one is a scalar. A scalar can be compared with another object of any size. <code>a < b</code> does an element-by-element comparison between <code>a</code> and <code>b</code> and returns a matrix of the same size with elements set to 1 where the relation is true, and 0 where the relation is false.
See Also	<code>eq</code> , <code>ge</code> , <code>gt</code> , <code>le</code> , <code>ne</code>

max

Purpose Return the largest element in an array of `fi` objects or the maximum value of a quantizer object before quantization

Syntax

```
max(a)
max(a,b)
[y,v] = max(a)
[y,v] = max(a,[],dim)
max(q)
```

Description

- For vectors, `max(a)` is the largest element in `a`.
- For matrices, `max(a)` is a row vector containing the maximum element from each column.
- For N-D arrays, `max(a)` operates along the first nonsingleton dimension.

`max(a,b)` returns an array the same size as `a` and `b` with the largest elements taken from `a` or `b`. Either one can be a scalar.

`[y,v] = max(a)` returns the indices of the maximum values in vector `v`. If the values along the first nonsingleton dimension contain more than one maximal element, the index of the first one is returned.

`[y,v] = max(a,[],dim)` operates along the dimension `dim`.

When complex, the magnitude `max(abs(a))` is used, and the angle `angle(a)` is ignored. NaNs are ignored when computing the maximum.

`max(q)` is the maximum value before quantization during a call to `quantize(q,...)` for quantizer object `q`. This value is the maximum value encountered over successive calls to `quantize` and is reset with `reset(q)`. `max(q)` is equivalent to `get(q,'max')` and `q.max`.

Examples

```
q = quantizer;
warning on
y = quantize(q,-20:10);
max(q)
Warning: 29 overflows.
ans =
```

10

See Also

min, quantize

mesh

Purpose Create mesh plot

Description Refer to the MATLAB mesh reference page for more information.

Purpose Create mesh plot with contour plot

Description Refer to the MATLAB `meshc` reference page for more information.

meshz

Purpose Create mesh plot with curtain plot

Description Refer to the MATLAB `meshz` reference page for more information.

Purpose Return the smallest element in an array of `fi` objects or the minimum value of a quantizer object before quantization

Syntax

```
min(a)
min(a,b)
[y,v] = min(a)
[y,v] = min(a,[],dim)
min(q)
```

Description

- For vectors, `min(a)` is the smallest element in `a`.
- For matrices, `min(a)` is a row vector containing the minimum element from each column.
- For N-D arrays, `min(a)` operates along the first nonsingleton dimension.

`min(a,b)` returns an array the same size as `a` and `b` with the smallest elements taken from `a` or `b`. Either one can be a scalar.

`[y,v] = min(a)` returns the indices of the minimum values in vector `v`. If the values along the first nonsingleton dimension contain more than one minimal element, the index of the first one is returned.

`[y,v] = min(a,[],dim)` operates along the dimension `dim`.

When complex, the magnitude `min(abs(a))` is used, and the angle `angle(a)` is ignored. NaNs are ignored when computing the minimum.

`min(q)` is the minimum value before quantization during a call to `quantize(q,...)` for quantizer object `q`. This value is the minimum value encountered over successive calls to `quantize` and is reset with `reset(q)`. `min(q)` is equivalent to `get(q,'min')` and `q.min`.

See Also `max`, `quantize`

minus

Purpose Return the matrix difference between `fi` objects

Syntax `minus(a,b)`

Description `minus(a,b)` is called for the syntax '`a - b`' when `a` or `b` is an object. `a - b` subtracts matrix `b` from matrix `a`. `a` and `b` must have the same dimensions unless one is a scalar (a 1-by-1 matrix). A scalar can be subtracted from anything.

See Also `mtimes`, `plus`, `times`, `uminus`

Purpose Multiply two objects using a `fimath` object

Syntax `c = F.mpy(a,b)`

Description `c = F.mpy(a,b)` performs elementwise multiplication on `a` and `b` using `fimath` object `F`. This is helpful in cases when you want to override the `fimath` objects of `a` and `b`, or if the `fimath` objects of `a` and `b` are different. `a` and `b` must have the same dimensions unless one is a scalar. If either `a` or `b` is scalar, then `c` has the dimensions of the nonscalar object. If either `a` or `b` is a `fi` object, and the other is a MATLAB built-in numeric type, then the built-in object is cast to the word length of the `fi` object, preserving best-precision fraction length.

Examples In this example, `c` is the 40-bit product of `a` and `b` with fraction length 30.

```
a = fi(pi);
b = fi(exp(1));
F = fimath('ProductMode','SpecifyPrecision',...
'ProductWordLength',40,'ProductFractionLength',30);
c = F.mpy(a, b)
```

```
c =
```

```
8.5397
```

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

```
RoundMode: round
OverflowMode: saturate
ProductMode: SpecifyPrecision
ProductWordLength: 40
ProductFractionLength: 30
```

```
SumMode: FullPrecision
MaxSumWordLength: 128
CastBeforeSum: true
```

Algorithm

`c = F.mpy(a,b)` is equivalent to

```
a.fimath = F;
b.fimath = F;
c = a .* b;
```

except that the `fimath` properties of `a` and `b` are not modified when you use the functional form.

See Also

`add`, `divide`, `fi`, `fimath`, `numericType`, `sub`, `sum`

Purpose Return the matrix product of `fi` objects

Syntax `mtimes(a,b)`

Description `mtimes(a,b)` is called for the syntax '`a * b`' when `a` or `b` is an object. `a * b` is the matrix product of `a` and `b`. Any scalar (a 1-by-1 matrix) can multiply anything. Otherwise, the number of columns of `a` must equal the number of rows of `b`.

See Also `plus`, `minus`, `times`, `uminus`

ndims

Purpose Return number of array dimensions

Description Refer to the MATLAB `ndims` reference page for more information.

Purpose	Determine whether the real-world values of two <code>fi</code> objects are not equal
Syntax	<code>c = ne(a,b)</code> <code>a ~= b</code>
Description	<code>c = ne(a,b)</code> is called for the syntax ' <code>a ~= b</code> ' when <code>a</code> or <code>b</code> is a <code>fi</code> object. <code>a</code> and <code>b</code> must have the same dimensions unless one is a scalar. A scalar can be compared with another object of any size. <code>a ~= b</code> does an element-by-element comparison between <code>a</code> and <code>b</code> and returns a matrix of the same size with elements set to 1 where the relation is true, and 0 where the relation is false.
See Also	<code>eq</code> , <code>ge</code> , <code>gt</code> , <code>le</code> , <code>lt</code>

not

Purpose Find logical NOT of array or scalar input

Description Refer to the MATLAB not reference page for more information.

Purpose	Return the number of quantization operations performed by a quantizer object
Syntax	<code>noperations(q)</code>
Description	<p><code>noperations(q)</code> is the number of quantization operations during a call to <code>quantize(q, ...)</code> for quantizer object <code>q</code>. This value accumulates over successive calls to <code>quantize</code>. You reset the value of <code>noperations</code> to zero by issuing the command <code>reset(q)</code>.</p> <p>Each time any data element is quantized, <code>noperations</code> is incremented by one. The real and complex parts are counted separately. For example, <code>(complex * complex)</code> counts four quantization operations for products and two for sum, since $(a+bi)(c+di) = (a*c - b*d) + (a*d + b*c)$. In contrast, <code>(real*real)</code> counts one quantization operation.</p> <p>In addition, the real and complex parts of the inputs are quantized individually. As a result, for a complex input of length 204 elements, <code>noperations</code> counts 408 quantizations: 204 for the real part of the input and 204 for the complex part.</p> <p>If any inputs, states, or coefficients are complex-valued, they are all expanded from real values to complex values, with a corresponding increase in the number of quantization operations recorded by <code>noperations</code>. In concrete terms, <code>(real*real)</code> requires fewer quantizations than <code>(real*complex)</code> and <code>(complex*complex)</code>. Changing all the values to complex because one is complex, such as the coefficient, makes the <code>(real*real)</code> into <code>(real*complex)</code>, raising <code>noperations</code> count.</p>
See Also	<code>get</code> , <code>quantizer</code> , <code>reset</code>

noverflows

Purpose	Return the number of overflows from quantization operations performed by a quantizer object
Syntax	<code>noverflows(q)</code>
Description	<code>noverflows(q)</code> returns the accumulated number of overflows resulting from quantization operations performed by a quantizer object <code>q</code> .
See Also	<code>get</code> , <code>max</code> , <code>range</code> , <code>reset</code>

Purpose Convert a number to a binary string using a quantizer object

Syntax `y = num2bin(q,x)`

Description `y = num2bin(q,x)` converts numeric array `x` into binary strings returned in `y`. When `x` is a cell array, each numeric element of `x` is converted to binary. If `x` is a structure, each numeric field of `x` is converted to binary.

`num2bin` and `bin2num` are inverses of one another, differing in that `num2bin` returns the binary strings in a column.

Examples

```
x = magic(3)/9;
q = quantizer([4,3]);
y = num2bin(q,x)
Warning: 1 overflow.
y =

0111
0010
0011
0000
0100
0111
0101
0110
0001
```

See Also `bin2num`, `hex2num`, `num2hex`, `num2int`

num2hex

Purpose Convert a number to its hexadecimal equivalent using a quantizer object

Syntax `y = num2hex(q,x)`

Description `y = num2hex(q,x)` converts numeric array `x` into hexadecimal strings returned in `y`. When `x` is a cell array, each numeric element of `x` is converted to hexadecimal. If `x` is a structure, each numeric field of `x` is converted to hexadecimal.

For fixed-point quantizer objects, the representation is two's complement. For floating-point quantizer objects, the representation is IEEE Standard 754 style.

For example, for `q = quantizer('double')`

```
num2hex(q,nan)
ans =
fff8000000000000
```

The leading fraction bit is 1, all other fraction bits are 0. Sign bit is 1, exponent bits are all 1.

```
num2hex(q,inf)
ans =
7ff0000000000000
```

Sign bit is 0, exponent bits are all 1, all fraction bits are 0.

```
num2hex(q,-inf)
ans =
fff0000000000000
```

Sign bit is 1, exponent bits are all 1, all fraction bits are 0.

num2hex and hex2num are inverses of each other, except that num2hex returns the hexadecimal strings in a column.

Examples

This is a floating-point example using a quantizer object q that has 6-bit word length and 3-bit exponent length.

```
x = magic(3);  
q = quantizer('float',[6 3]);  
y = num2hex(q,x)
```

```
y =
```

```
18  
12  
14  
0c  
15  
18  
16  
17  
10
```

See Also

bin2num, hex2num, num2bin, num2int

num2int

Purpose Convert a number to a signed integer

Syntax
`y = num2int(q,x)`
`[y1,y,...] = num2int(q,x1,x,...)`

Description `y = num2int(q,x)` uses `q.format` to convert numeric `x` to an integer.
`[y1,y,...] = num2int(q,x1,x,...)` uses `q.format` to convert numeric values `x1, x2,...` to integers `y1,y2,...`.

Examples All the two's complement 4-bit numbers in fractional form are given by

```
x = [0.875 0.375 -0.125 -0.625  
      0.750 0.250 -0.250 -0.750  
      0.625 0.125 -0.375 -0.875  
      0.500 0.000 -0.500 -1.000];
```

```
q=quantizer([4 3]);
```

```
y = num2int(q,x)  
y =
```

```
 7     3     -1     -5  
 6     2     -2     -6  
 5     1     -3     -7  
 4     0     -4     -8
```

Algorithm When `q` is a fixed-point quantizer object, `f` is equal to `fractionlength(q)`, and `x` is numeric

$$y = x \times 2^f$$

When `q` is a floating-point quantizer object, `y = x`. `num2int` is meaningful only for fixed-point quantizer objects.

See Also `bin2num`, `hex2num`, `num2bin`, `num2hex`

Purpose Return number of data elements in `fi` array

Syntax `numberofelements(a)`

Description `numberofelements(a)` returns the number of data elements in a `fi` array. `numberofelements(a) == prod(size(a))`.

Note that `fi` is a MATLAB object, and therefore `numel(a)` returns 1 when `a` is a `fi` object. Refer to the information about classes in the MATLAB `numel` reference page.

See Also `max`, `min`, `numel`

numerictype

Purpose Construct a numerictype object

Syntax

```
T = numerictype
T = numerictype(s)
T = numerictype(s,w)
T = numerictype(s,w,f)
T = numerictype(s,w,slope,bias)
T = numerictype(s,w,slopeadjustmentfactor,fixedexponent,bias)
T = numerictype(property1,value1, ...)
T = numerictype(T1, property1, value1, ...)
```

Description You can use the numerictype constructor function in the following ways:

- `T = numerictype` creates a default numerictype object.
- `T = numerictype(s)` creates a numerictype object with Fixed-point: binary point scaling, signedness `s`, 16-bit word length and 15-bit fraction length.
- `T = numerictype(s,w)` creates a numerictype object with Fixed-point: binary point scaling, signedness `s`, word length `w` and 15-bit fraction length.
- `T = numerictype(s,w,f)` creates a numerictype object with Fixed-point: binary point scaling, signedness `s`, word length `w` and fraction length `f`.
- `T = numerictype(s,w,slope,bias)` creates a numerictype object with Fixed-point: slope and bias scaling, signedness `s`, word length `w`, slope, and bias.
- `T = numerictype(s,w,slopeadjustmentfactor,fixedexponent,bias)` creates a numerictype object with Fixed-point: slope and bias scaling, signedness `s`, word length `w`, slopeadjustmentfactor, fixedexponent, and bias.
- `T = numerictype(property1,value1, ...)` allows you to set properties for a numerictype object using property name/property value pairs.

- `T = numerictype(T1, property1, value1, ...)` allows you to make a copy of an existing `numerictype` object, while modifying any or all of the property values.

The properties of the `numerictype` object are listed below. These properties are described in detail in “`numerictype` Object Properties” on page 9-12.

- `Bias` – Bias
- `DataType` – Data type category
- `DataTypeMode` – Data type and scaling mode
- `FixedExponent` – Fixed-point exponent
- `SlopeAdjustmentFactor` – Slope adjustment
- `FractionLength` – Fraction length of the stored integer value, in bits
- `Scaling` – Fixed-point scaling mode
- `Signed` – Signed or unsigned
- `Slope` – Slope
- `WordLength` – Word length of the stored integer value, in bits

Examples

Example 1

Type

```
T = numerictype
```

to create a default `numerictype` object.

```
T =
```

```
DataType: Fixed  
Scaling: BinaryPoint  
Signed: true
```

numerictype

```
WordLength: 16
FractionLength: 15
```

Example 2

The following creates a signed numerictype object with a 32-bit word length and 30-bit fraction length.

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

```
T =
```

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

Example 3

If you omit the argument *f*, it is automatically set to the best precision possible.

```
T = numerictype(1, 32)
```

```
T =
```

```
DataTypeMode: Fixed-point: binary point scaling
Signed: true
WordLength: 32
FractionLength: 15
```

Example 4

```
T = numerictype(1)
```

```
T =
```

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

Example 5

```
T = numerictype('Signed', true, 'DataTypeMode', ...
'Fixed-point: slope and bias', 'WordLength', 32, 'Slope', ...
2^-2, 'Bias', 4)
```

```
T =
```

```
DataTypeMode: Fixed-point: slope and bias scaling
Signed: true
WordLength: 32
Slope: 0.25
Bias: 4
```

Example 6

To copy a numerictype object, use the numerictype constructor function:

```
T = numerictype
```

```
T =
```

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

```
U = numerictype(T)
```

numerictype

U =

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

See Also

fi, fimath, fipref, quantizer, “numerictype Object Properties” on page 9-12

Purpose	Return the number of underflows from quantization operations performed by a quantizer object
Syntax	<code>nunderflows(q)</code>
Description	<code>nunderflows(q)</code> returns the accumulated number of underflows resulting from quantization operations performed by a quantizer object <code>q</code> . An underflow is defined as a number that is nonzero before it is quantized, and zero after it is quantized.
See Also	<code>denormalmin</code> , <code>eps</code> , <code>quantize</code> , <code>quantizer</code> , <code>reset</code>

oct

Purpose Return the octal representation of the stored integer of a `fi` object as a string

Syntax `oct(a)`

Description Fixed-point numbers can be represented as

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

or, equivalently,

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

The stored integer is the raw binary number, in which the binary point is assumed to be at the far right of the word.

`oct(a)` returns the stored integer of `fi` object `a` in octal format as a string.

Examples

Example 1

The following code

```
a = fi([-1 1],1,8,7);  
oct(a)
```

returns

```
200 177
```

See Also

`bin`, `dec`, `hex`, `int`

Purpose Find logical OR of array or scalar inputs

Description Refer to the MATLAB or reference page for more information.

patch

Purpose Create patch graphics object

Description Refer to the MATLAB patch reference page for more information.

Purpose Create pseudocolor plot

Description Refer to the MATLAB `pcolor` reference page for more information.

permute

Purpose Rearrange the dimensions of a multidimensional array

Description Refer to the MATLAB permute reference page for more information.

Purpose Create linear 2-D plot

Description Refer to the MATLAB `plot` reference page for more information.

plot3

Purpose Create 3-D line plot

Description Refer to the MATLAB `plot3` reference page for more information.

Purpose Draw scatter plots

Description Refer to the MATLAB `plotmatrix` reference page for more information.

plotyy

Purpose Create graph with y-axes on both right and left sides

Description Refer to the MATLAB `plotyy` reference page for more information.

Purpose Return the matrix sum of `fi` objects

Syntax `plus(a,b)`

Description `plus(a,b)` is called for the syntax '`a + b`' when `a` or `b` is an object. `a + b` adds matrices `a` and `b`. `a` and `b` must have the same dimensions unless one is a scalar (a 1-by-1 matrix). A scalar can be added to anything.

See Also `minus`, `mtimes`, `times`, `uminus`

polar

Purpose Plot polar coordinates

Description Refer to the MATLAB `polar` reference page for more information.

Purpose Multiply by a power of 2

Syntax `b = pow2(a, K)`

Description `b = pow2(a, K)` returns

$$b = a \times 2^K$$

where K is an integer and a and b are `fi` objects. If K is a non-integer, it will be rounded to `floor` before the calculation is performed. The scaling of a must be equivalent to binary point-only scaling; in other words, it must have a fractional slope of 1 and a bias of 0.

The syntax `b = pow2(a)` is not supported when a is a `fi` object.

a can be real or complex. If a is complex, `pow2` operates on both the real and complex portions of a .

Examples The following example shows the use of `pow2` with a complex `fi` object:

```
format long g
P = fipref('NumericTypeDisplay', 'short', 'FimathDisplay',...
'none');
a = fi(57 - 2i, 1, 16, 8)

a =

                    57 -                    2i

                    s16,8
pow2(a, 2)

ans =

                    127.99609375 -                    8i

                    s16,8
```

pow2

See Also

`bitshift`

Purpose Apply a quantizer object to data

Syntax `y = quantize(q, x)`
`[y1,y2,...]quantize(q,x1,x2,...)`

Description `y = quantize(q, x)` uses the quantizer object `q` to quantize `x`. When `x` is a numeric array, each element of `x` is quantized. When `x` is a cell array, each numeric element of the cell array is quantized. When `x` is a structure, each numeric field of `x` is quantized. Nonnumeric elements or fields of `x` are left unchanged and `quantize` does not issue warnings for nonnumeric values.

`[y1,y2,...]quantize(q,x1,x2,...)` is equivalent to

`y1 = quantize(q,x1), y2 = quantize(q,x2),...`

The quantizer object states

- `max` – Maximum value before quantizing
- `min` – Minimum value before quantizing
- `noverflows` – Number of overflows
- `nunderflows` – Number of underflows
- `noperations` – Number of quantization operations

are updated during the call to `quantize`, and running totals are kept until a call to `reset` is made.

Examples The following examples demonstrate using `quantize` to quantize data.

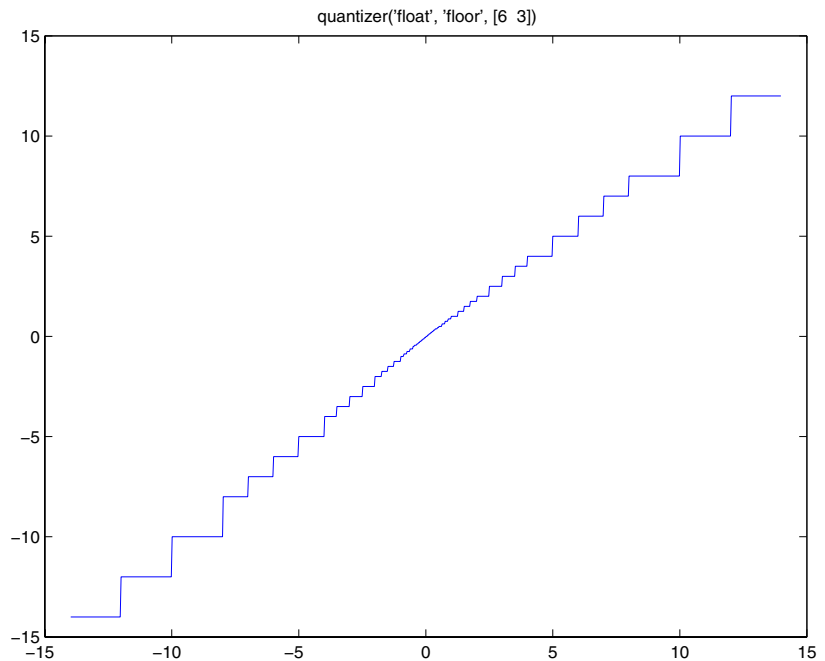
Example 1 - Custom Precision Floating-Point

The code listed here produces the plot shown in the following figure.

```
u=linspace(-15,15,1000);  
q=quantizer([6 3], 'float');  
range(q)
```

quantize

```
ans =  
-14    14  
y=quantize(q,u);  
plot(u,y);title(tostring(q))  
Warning: 68 overflows.
```

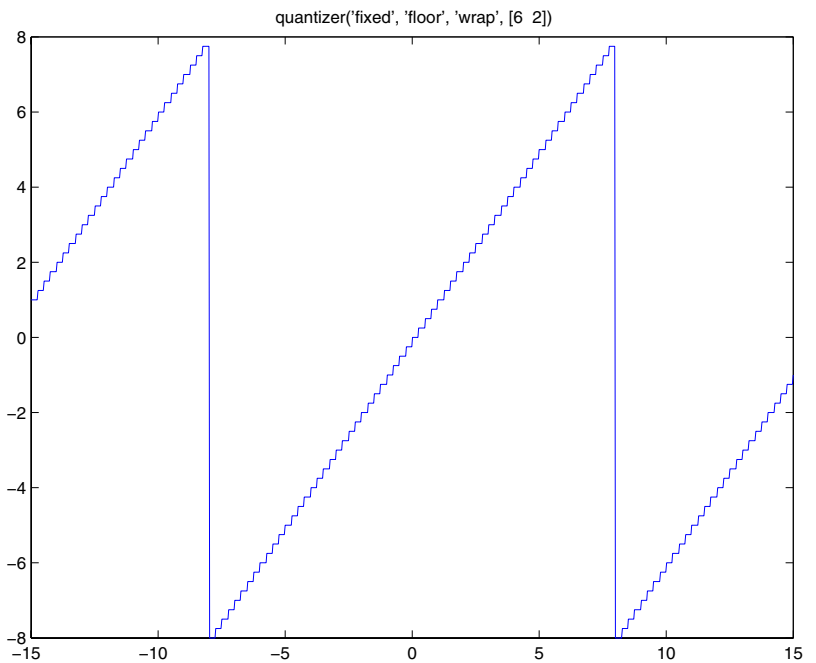


Example 2 - Fixed-Point

The code listed here produces the plot shown in the following figure.

```
u=linspace(-15,15,1000);  
q=quantizer([6 2], 'wrap');
```

```
range(q)  
  
ans =  
  
    -8.0000    7.7500  
y=quantize(q,u);  
plot(u,y);title(tostring(q))  
Warning: 468 overflows.
```



See Also quantizer, set

quantizer

Purpose Construct a quantizer object

Syntax

```
q = quantizer
q = quantizer('PropertyName',PropertyValue1,...)
q = quantizer(PropertyValue1,PropertyValue2,...)
q = quantizer(struct)
q = quantizer(pn,pv)
```

Description `q = quantizer` creates a quantizer object with properties set to their default values.

`q = quantizer('PropertyName',PropertyValue1,...)` uses property name/ property value pairs.

`q = quantizer(PropertyValue1,PropertyValue2,...)` creates a quantizer object with the listed property values. When two values conflict, `quantizer` sets the last property value in the list. Property values are unique; you can set the property names by specifying just the property values in the command.

`q = quantizer(struct)`, where `struct` is a structure whose field names are property names, sets the properties named in each field name with the values contained in the structure.

`q = quantizer(pn,pv)` sets the named properties specified in the cell array of strings `pn` to the corresponding values in the cell array `pv`.

The quantizer object property values are listed below. These properties are described in detail in “quantizer Object Properties” on page 9-16.

Property Name	Property Value	Description
mode	'double'	Double-precision mode. Override all other parameters.
	'float'	Custom-precision floating-point mode.

Property Name	Property Value	Description
	'fixed'	Signed fixed-point mode.
	'single'	Single-precision mode. Override all other parameters.
	'ufixed'	Unsigned fixed-point mode.
roundmode	'ceil'	Round toward positive infinity.
	'convergent'	Convergent rounding.
	'fix'	Round toward zero.
	'floor'	Round toward negative infinity.
	'round'	Round toward nearest.
overflowmode (fixed-point only)	'saturate'	Saturate on overflow.
	'wrap'	Wrap on overflow.
format	[wordlength exponentlength]	Format for fixed or ufixed mode.
	[wordlength exponentlength]	Format for float mode.

The default property values for a quantizer object are

```

mode = 'fixed';
roundmode = 'floor';
overflowmode = 'saturate';
format = [16 15];

```

quantizer

Along with the preceding properties, quantizer objects have read-only properties: 'max', 'min', 'noverflows', 'nunderflows', and 'noperations'. They can be accessed through `quantizer/get` or `q.max`, `q.min`, `q.noverflows`, `q.nunderflows`, and `q.noperations`, but they cannot be set. They are updated during the `quantizer/quantize` method, and are reset by the `quantizer/reset` method.

The following table lists the read-only quantizer object properties:

Property Name	Description
'max'	Maximum value before quantizing
'min'	Minimum value before quantizing
'noverflows'	Number of overflows
'nunderflows'	Number of underflows
'noperations'	Number of data points quantized

Examples

The following example operations are equivalent.

Setting quantizer object properties by listing property values only in the command,

```
q = quantizer('fixed', 'ceil', 'saturate', [5 4])
```

Using a structure struct to set quantizer object properties,

```
struct.mode = 'fixed';  
struct.roundmode = 'ceil';  
struct.overflowmode = 'saturate';  
struct.format = [5 4];  
q = quantizer(struct);
```

Using property name and property value cell arrays pn and pv to set quantizer object properties,

```
pn = {'mode', 'roundmode', 'overflowmode', 'format'};  
pv = {'fixed', 'ceil', 'saturate', [5 4]};  
q = quantizer(pn, pv)
```

Using property name/property value pairs to configure a quantizer object,

```
q = quantizer('mode', 'fixed', 'roundmode', 'ceil', ...  
             'overflowmode', 'saturate', 'format', [5 4]);
```

See Also

`fi`, `fimath`, `fipref`, `numericType`, `quantize`, `set`, “[quantizer Object Properties](#)” on page 9-16

quiver

Purpose Create quiver or velocity plot

Description Refer to the MATLAB `quiver` reference page for more information.

Purpose Create 3-D quiver or velocity plot

Description Refer to the MATLAB `quiver3` reference page for more information.

randquant

Purpose Generate a uniformly distributed, quantized random number using a quantizer object

Syntax

```
randquant(q,n)
randquant(q,m,n)
randquant(q,m,n,p,...)
randquant(q,[m,n])
randquant(q,[m,n,p,...])
```

Description `randquant(q,n)` uses quantizer object `q` to generate an `n`-by-`n` matrix with random entries whose values cover the range of `q` when `q` is a fixed-point quantizer object. When `q` is a floating-point quantizer object, `randquant` populates the `n`-by-`n` array with values covering the range

-[square root of `realmax(q)`] to [square root of `realmax(q)`]

`randquant(q,m,n)` uses quantizer object `q` to generate an `m`-by-`n` matrix with random entries whose values cover the range of `q` when `q` is a fixed-point quantizer object. When `q` is a floating-point quantizer object, `randquant` populates the `m`-by-`n` array with values covering the range

-[square root of `realmax(q)`] to [square root of `realmax(q)`]

`randquant(q,m,n,p,...)` uses quantizer object `q` to generate an `m`-by-`n`-by-`p`-by ... matrix with random entries whose values cover the range of `q` when `q` is fixed-point quantizer object. When `q` is a floating-point quantizer object, `randquant` populates the matrix with values covering the range

-[square root of `realmax(q)`] to [square root of `realmax(q)`]

`randquant(q,[m,n])` uses quantizer object `q` to generate an `m`-by-`n` matrix with random entries whose values cover the range of `q` when `q` is a fixed-point quantizer object. When `q` is a floating-point quantizer object, `randquant` populates the `m`-by-`n` array with values covering the range

-[square root of `realmax(q)`] to [square root of `realmax(q)`]

`randquant(q,[m,n,p,...])` uses quantizer object `q` to generate `p` `m`-by-`n` matrices containing random entries whose values cover the range of `q` when `q` is a fixed-point quantizer object. When `q` is a floating-point quantizer object, `randquant` populates the `m`-by-`n` arrays with values covering the range

-[square root of `realmax(q)`] to [square root of `realmax(q)`]

`randquant` produces pseudorandom numbers. The number sequence `randquant` generates during each call is determined by the state of the generator. Because MATLAB resets the random number generator state at startup, the sequence of random numbers generated by the function remains the same unless you change the state.

`randquant` works like `rand` in most respects, including the generator used, but it does not support the 'state' and 'seed' options available in `rand`.

Examples

```
q=quantizer([4 3]);  
rand('state',0)  
randquant(q,3)
```

ans =

```
    0.7500    -0.1250    -0.2500  
   -0.6250     0.6250    -1.0000  
    0.1250     0.3750     0.5000
```

See Also

`quantizer`, `range`, `realmax`

range

Purpose Return the numerical range of a `fi` object or quantizer object

Syntax

```
range(a)
[min, max] = range(a)
r = range(q)
[min, max] = range(q)
```

Description `range(a)` returns the minimum and maximum possible values of `fi` object `a` in two-vector format. All possible quantized real-world values of `a` are in the range returned. If `a` is a complex number, then all possible values of `real(a)` and `imag(a)` are in the range returned.

`[min, max] = range(a)` returns the minimum and maximum values of `fi` object `a` in separate output variables.

`r = range(q)` returns the two-element row vector $r = [a \ b]$ such that for all real x , $y = \text{quantize}(q, x)$ returns y in the range $a \leq y \leq b$.

`[min, max] = range(q)` returns the minimum and maximum values of the range in separate output variables.

Examples

```
q = quantizer('float',[6 3]);
r = range(q)

r =

    -14     14
q = quantizer('fixed',[4 2],'floor');
[min,max] = range(q)

min =

    -2

max =

    1.7500
```


Algorithm

If q is a floating-point quantizer object, $a = -\text{realmax}(q)$, $b = \text{realmax}(q)$.

If q is a signed fixed-point quantizer object (`datamode = 'fixed'`),

$$a = -\text{realmax}(q) - \text{eps}(q) = \frac{-2^{w-1}}{2^f}$$

$$b = \text{realmax}(q) = \frac{2^{w-1} - 1}{2^f}$$

If q is an unsigned fixed-point quantizer object (`datamode = 'ufixed'`),

$$a = 0$$

$$b = \text{realmax}(q) = \frac{2^w - 1}{2^f}$$

See `realmax` for more information.

See Also

`exponentmin`, `fractionlength`, `max`, `min`, `realmax`, `realmin`

real

Purpose Return real part of complex number

Description Refer to the MATLAB `real` reference page for more information.

Purpose Return the largest positive fixed-point value or quantized number

Syntax `realmax(a)`
`realmax(q)`

Description `realmax(a)` is the largest real-world value that can be represented in the data type of fi object `a`. Anything larger overflows.

`realmax(q)` is the largest quantized number that can be represented where `q` is a quantizer object. Anything larger overflows.

Examples

```
q = quantizer('float',[6 3]);
x = realmax(q)

x =

    14
```

Algorithm If `q` is a floating-point quantizer object, the largest positive number, x , is

$$x = 2^{E_{max}} \cdot (2 - eps(q))$$

If `q` is a signed fixed-point quantizer object, the largest positive number, x , is

$$x = \frac{2^{w-1} - 1}{2^f}$$

If `q` is an unsigned fixed-point quantizer object (`datamode = 'ufixed'`), the largest positive number, x , is

$$x = \frac{2^w - 1}{2^f}$$

See Also `quantizer`, `realmin`, `exponentmin`, `fractionlength`

realmin

Purpose Return the smallest positive normalized fixed-point value or quantized number

Syntax `realmin(a)`
`realmin(q)`

Description `realmin(a)` is the smallest real-world value that can be represented in the data type of fi object `a`. Anything smaller underflows.

`realmin(q)` is the smallest positive normal quantized number where `q` is a quantizer object. Anything smaller than `x` underflows or is an IEEE "denormal" number.

Examples

```
q = quantizer('float',[6 3]);  
realmin(q)
```

```
ans =  
  
    0.2500
```

Algorithm If `q` is a floating-point quantizer object, $x = 2^{E_{min}}$ where $E_{min} = \text{exponentmin}(q)$ is the minimum exponent.

If `q` is a signed or unsigned fixed-point quantizer object, $x = 2^{-f} = \epsilon$ where f is the fraction length.

See Also `exponentmin`, `fractionlength`, `realmax`

Purpose Replicate and tile an array

Description Refer to the MATLAB repmat reference page for more information.

rescale

Purpose Change the scaling of a `fi` object

Syntax

```
b = rescale(a, fractionlength)
b = rescale(a, slope, bias)
b = rescale(a, slopeadjustmentfactor, fixedexponent, bias)
b = rescale(a, ..., PropertyName, PropertyValue, ...)
```

Description The `rescale` function acts similarly to the `fi` copy function with the following exceptions:

- The `fi` copy constructor preserves the real-world value, while `rescale` preserves the stored integer value.
- `rescale` does not allow the `Signed` and `WordLength` properties to be changed.

Examples In the following example, `fi` object `a` is rescaled to create `fi` object `b`. The real-world values of `a` and `b` are different, while their stored integer values are the same:

```
p = fipref('FimathDisplay','none',...
          'NumericTypeDisplay','short');
a = fi(10, 1, 8, 3)

a =

    10

    s8,3

b = rescale(a, 1)

b =

    40
```

```
s8,1  
  
stored_integer_a = a.int;  
stored_integer_b = b.int;  
isequal(stored_integer_a, stored_integer_b)  
  
ans =  
  
1
```

See Also

fi

reset

Purpose	Reset one or more objects to their initial conditions
Syntax	<code>reset(obj)</code> <code>reset(q1, q2, ...)</code>
Description	<p><code>reset(obj)</code> resets <code>fi</code>, <code>fimath</code>, <code>fipref</code>, or quantizer object <code>obj</code> to its initial conditions.</p> <p><code>reset(q1, q2, ...)</code> resets the states of the quantizer objects <code>q1</code>, <code>q2, ...</code> to their initial conditions.</p> <p>The states of a quantizer object are</p> <ul style="list-style-type: none">• <code>max</code> – Maximum value before quantizing• <code>min</code> – Minimum value before quantizing• <code>noverflows</code> – Number of overflows• <code>nunderflows</code> – Number of underflows• <code>noperations</code> – Number of quantization operations performed
See Also	<code>quantizer</code> , <code>set</code>

Purpose Reshape array

Description Refer to the MATLAB reshape reference page for more information.

rgbplot

Purpose Plot colormap

Description Refer to the MATLAB `rgbplot` reference page for more information.

Purpose Create ribbon plot

Description Refer to the MATLAB ribbon reference page for more information.

rose

Purpose Create angle histogram

Description Refer to the MATLAB rose reference page for more information.

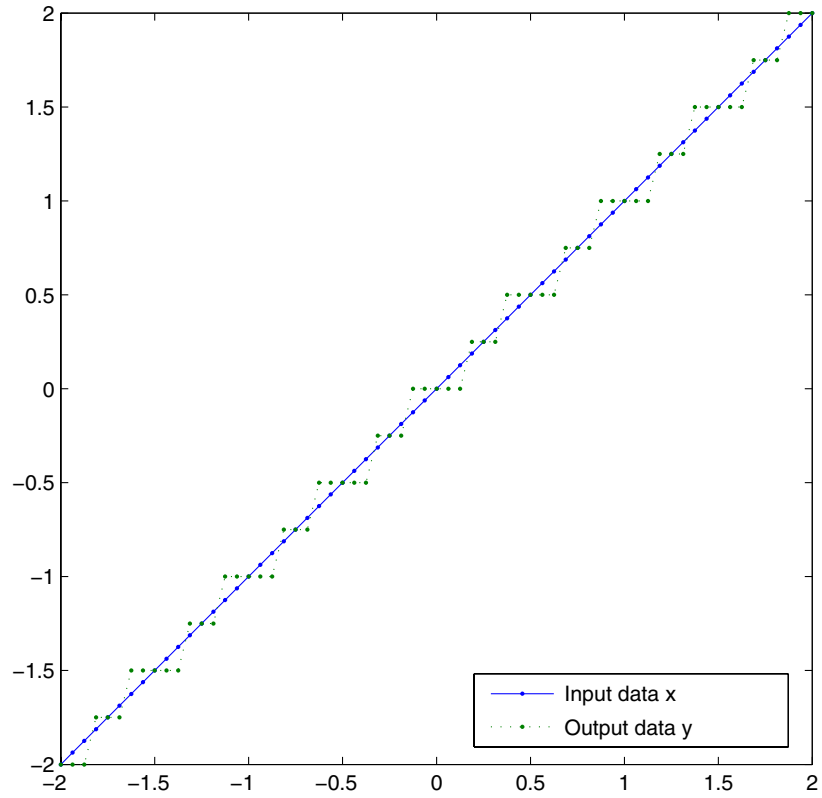
- Purpose** Round input data using a quantizer object without checking for overflow
- Syntax** `round(q,x)`
- Description** `round(q,x)` uses the `RoundMode` and `FractionLength` settings of `q` to round the numeric data `x`, but does not check for overflows during the operation. Compare to `quantize`.

Example Create a quantizer object and use it to quantize input data. The quantizer object applies its properties to the input data to return quantized output.

```
q = quantizer('fixed', 'convergent', 'wrap', [3 2]);  
x = (-2:eps(q)/4:2)';  
y = round(q,x);  
plot(x,[x,y],'.-'); axis square;
```

Applying quantizer object `q` to the data results in the staircase shape output plot shown here. Where the input data is linear, output `y` shows distinct quantization levels.

round



See Also

quantize, quantizer

Purpose	Save <code>fi</code> preferences for the next MATLAB session
Syntax	<code>savefipref</code>
Description	<code>savefipref</code> saves the settings of the current <code>fipref</code> object for the next MATLAB session.
See Also	<code>fipref</code>

scatter

Purpose Create a scatter or bubble plot

Description Refer to the MATLAB scatter reference page for more information.

Purpose Create a 3-D scatter or bubble plot

Description Refer to the MATLAB `scatter3` reference page for more information.

sdec

Purpose Return signed decimal representation of stored integer of `fi` object as string

Syntax `sdec(a)`

Description Fixed-point numbers can be represented as

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

or, equivalently,

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

The stored integer is the raw binary number, in which the binary point is assumed to be at the far right of the word.

`sdec(a)` returns the stored integer of `fi` object `a` in signed decimal format as a string.

Examples The code

```
a = fi([-1 1],1,8,7);  
sdec(a)
```

returns

```
-128    127
```

See Also `bin`, `dec`, `hex`, `int`, `oct`

Purpose Create semilogarithmic plot with logarithmic x-axis

Description Refer to the MATLAB `semilogx` reference page for more information.

semilogy

Purpose Create semilogarithmic plot with logarithmic y-axis

Description Refer to the MATLAB `semilogy` reference page for more information.

Purpose

Set or display property values for quantizer objects

Syntax

```
set(q, PropertyValue1, PropertyValue2,...)
set(q,s)
set(q,pn,pv)
set(q,'PropertyName1',PropertyValue1,'PropertyName2',
    PropertyValue2,...)
q.PropertyName = Value
s = set(q)
```

Description

`set(q, PropertyValue1, PropertyValue2, ...)` sets the properties of quantizer object `q`. If two property values conflict, the last value in the list is the one that is set.

`set(q, s)`, where `s` is a structure whose field names are object property names, sets the properties named in each field name with the values contained in the structure.

`set(q, pn, pv)` sets the named properties specified in the cell array of strings `pn` to the corresponding values in the cell array `pv`.

`set(q, 'PropertyName1', PropertyValue1, 'PropertyName2', PropertyValue2, ...)` sets multiple property values with a single statement. Note that you can use property name/property value string pairs, structures, and property name/property value cell array pairs in the same call to `set`.

`q.PropertyName = Value` uses dot notation to set property `PropertyName` to `Value`.

`set(q)` displays the possible values for all properties of quantizer object `q`.

`s = set(q)` returns a structure containing the possible values for the properties of quantizer object `q`.

The states are cleared when you set any value other than `WarnIfOverflow`.

set

See Also

[get](#)

Purpose Perform signum function on array

Syntax `c = sign(a)`

Description `c = sign(a)` returns an array `c` the same size as `a`, where each element of `c` is

- 1 if the corresponding element of `a` is greater than zero
- 0 if the corresponding element of `a` is zero
- -1 if the corresponding element of `a` is less than zero

The elements of `c` are of data type `int8`.

`sign` does not support complex `fi` inputs.

single

Purpose Return the single-precision floating-point real-world value of a `fi` object

Syntax `single(a)`

Description Fixed-point numbers can be represented as

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

or, equivalently,

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

`single(a)` returns the real-world value of a `fi` object in single-precision floating point.

See Also `double`

Purpose Return array dimensions

Description Refer to the MATLAB `size` reference page for more information.

slice

Purpose Create volumetric slice plot

Description Refer to the MATLAB `slice` reference page for more information.

Purpose Visualize sparsity pattern

Description Refer to the MATLAB `spy` reference page for more information.

squeeze

Purpose Remove singleton dimensions

Description Refer to the MATLAB `squeeze` reference page for more information.

Purpose Create stairstep graph

Description Refer to the MATLAB stairs reference page for more information.

stem

Purpose Plot discrete sequence data

Description Refer to the MATLAB `stem` reference page for more information.

Purpose Plot 3-D discrete sequence data

Description Refer to the MATLAB `stem3` reference page for more information.

streamribbon

Purpose Create a 3-D stream ribbon plot

Description Refer to the MATLAB streamribbon reference page for more information.

Purpose Draw streamlines in slice planes

Description Refer to the MATLAB `streamslice` reference page for more information.

streamtube

Purpose Create a 3-D stream tube plot

Description Refer to the MATLAB streamtube reference page for more information.

Purpose Return the stored integer of a `fi` object

Syntax `I = stripscaling(a)`

Description `I = stripscaling(a)` returns the stored integer of `a` as a `fi` object with zero bias and the same word length and sign as `a`.

sub

Purpose Subtract two objects using a `fimath` object

Syntax `c = F.sub(a,b)`

Description `c = F.sub(a,b)` subtracts objects `a` and `b` using `fimath` object `F`. This is helpful in cases when you want to override the `fimath` objects of `a` and `b`, or if the `fimath` objects of `a` and `b` are different.

`a` and `b` must have the same dimensions unless one is a scalar. If either `a` or `b` is scalar, then `c` has the dimensions of the nonscalar object.

If either `a` or `b` is a `fi` object, and the other is a MATLAB built-in numeric type, then the built-in object is cast to the word length of the `fi` object, preserving best-precision fraction length.

Examples In this example, `c` is the 32-bit difference of `a` and `b` with fraction length 16.

```
a = fi(pi);
b = fi(exp(1));
F = fimath('SumMode','SpecifyPrecision',...
'SumWordLength',32,'SumFractionLength',16);
c = F.sub(a, b)
```

```
c =
```

```
0.4233
```

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

```
RoundMode: round
OverflowMode: saturate
ProductMode: FullPrecision
MaxProductWordLength: 128
```

```
SumMode: SpecifyPrecision
SumWordLength: 32
SumFractionLength: 16
CastBeforeSum: true
```

Algorithm

`c = F.sub(a,b)` is equivalent to

```
a.fimath = F;
b.fimath = F;
c = a - b;
```

except that the `fimath` properties of `a` and `b` are not modified when you use the functional form.

See Also

`add`, `divide`, `fi`, `fimath`, `mpy`, `numericType`

subsasgn

Purpose Subscripted assignment

Syntax

```
a(I) = b
a(I,J) = b
a(I,:) = b
a(:,I) = b
a(I,J,K,...) = b
a = subsasgn(a,S,b)
```

Description `a(I) = b` assigns the values of `b` into the elements of `a` specified by the subscript vector `I`. `b` must have the same number of elements as `I` or be a scalar.

`a(I,J) = b` assigns the values of `b` into the elements of the rectangular submatrix of `a` specified by the subscript vectors `I` and `J`. `b` must have `LENGTH(I)` rows and `LENGTH(J)` columns.

A colon used as a subscript, as in `a(I,:) = b` or `a(:,I) = b` indicates the entire column or row.

For multidimensional arrays, `a(I,J,K,...) = b` assigns `b` to the specified elements of `a`. `b` must be `length(I)-by-length(J)-by-length(K)-...` or be shiftable to that size by adding or removing singleton dimensions.

`a = subsasgn(a,S,b)` is called for the syntax `a(i)=b`, `a{i}=b`, or `a.i=b` when `a` is an object. `S` is a structure array with the fields

- `type` – String containing `'()'`, `'{}'`, or `'.'` specifying the subscript type
- `subs` – Cell array or string containing the actual subscripts

For instance, the syntax `a(1:2,:)=b` calls `a=subsasgn(a,S,b)` where `S` is a 1-by-1 structure with `S.type='()'` and `S.subs = {1:2, ':'}`. A colon used as a subscript is passed as the string `':'`.

See Also `subref`

Purpose Subscripted reference

Description Refer to the MATLAB subsref reference page for more information.

sum

Purpose Return sum of array elements

Syntax
`b = sum(a)`
`b = sum(a, dim)`

Description `b = sum(a)` returns the sum along different dimensions of the `fi` array `a`.

If `a` is a vector, `sum(a)` returns the sum of the elements.

If `a` is a matrix, `sum(a)` treats the columns of `a` as vectors, returning a row vector of the sums of each column.

If `a` is a multidimensional array, `sum(a)` treats the values along the first nonsingleton dimension as vectors, returning an array of row vectors.

`b = sum(a, dim)` sums along the dimension `dim` of `a`.

The `fimath` object is used in the calculation of the sum. If `SumMode` is `FullPrecision`, `KeepLSB`, or `KeepMSB`, then the number of integer bits of growth for `sum(a)` is `ceil(log2(length(a)))`.

See Also `add`, `divide`, `fi`, `fimath`, `mpy`, `numerictype`, `sub`

Purpose Create 3-D shaded surface plot

Description Refer to the MATLAB `surf` reference page for more information.

surf

Purpose Create 3-D shaded surface plot with contour plot

Description Refer to the MATLAB `surf` reference page for more information.

Purpose Create a surface plot with colormap-based lighting

Description Refer to the MATLAB `surf1` reference page for more information.

surfnorm

Purpose Compute and display 3-D surface normals

Description Refer to the MATLAB `surfnorm` reference page for more information.

Purpose Create text object in current axes

Description Refer to the MATLAB text reference page for more information.

times

Purpose	Return the result of element-by-element multiplication of <code>fi</code> objects
Syntax	<code>times(a,b)</code>
Description	<code>times(a,b)</code> is called for the syntax ' <code>a .* b</code> ' when <code>a</code> or <code>b</code> is an object. <code>a .* b</code> denotes element-by-element multiplication. <code>a</code> and <code>b</code> must have the same dimensions unless one is a scalar. A scalar can be multiplied into anything.
See Also	<code>plus</code> , <code>minus</code> , <code>mtimes</code> , <code>uminus</code>

Purpose Create Toeplitz matrix

Syntax `t = toeplitz(a, b)`
`t = toeplitz(b)`

Description `t = toeplitz(a, b)` returns a nonsymmetric Toeplitz matrix having `a` as its first column and `b` as its first row. `b` is cast to the `numericType` of `a`.
`t = toeplitz(b)` returns the symmetric or Hermitian Toeplitz matrix formed from vector `b`, where `b` is the first row of the matrix.
The `numericType` and `fimath` objects of the leftmost input that is a `fi` object are applied to the output `t`.

tostring

Purpose Convert a quantizer object to a string

Syntax `s = tostring(q)`

Description `s = tostring(q)` converts quantizer object `q` to a string `s`. After converting `q` to a string, the function `eval(s)` can use `s` to create a quantizer object with the same properties as `q`.

Examples When you use `tostring` with a quantizer object you see the following response:

```
q = quantizer

q =

    DataMode = fixed
    RoundMode = floor
    OverflowMode = saturate
    Format = [16 15]

    Max = reset
    Min = reset
    NOverflows = 0
    NUnderflows = 0
    NOperations = 0

s = tostring(q)

s =

quantizer('fixed', 'floor', 'saturate', [16 15])

eval(s)

ans =

    DataMode = fixed
```



```
RoundMode = floor
OverflowMode = saturate
Format = [16 15]
```

```
Max = reset
Min = reset
NOverflows = 0
NUnderflows = 0
NOperations = 0
```

Note that s is the same as q .

See Also

quantizer

transpose

Purpose Return the transpose

Description Refer to the MATLAB arithmetic operators reference page for more information.

Purpose Plot picture of tree

Description Refer to the MATLAB `treeplot` reference page for more information.

tril

Purpose Return the lower triangular part of a matrix

Description Refer to the MATLAB `tril` reference page for more information.

Purpose Create triangular mesh plot

Description Refer to the MATLAB `trimesh` reference page for more information.

triplot

Purpose Create 2-D triangular plot

Description Refer to the MATLAB triplot reference page for more information.

Purpose Create triangular surface plot

Description Refer to the MATLAB `trisurf` reference page for more information.

triu

Purpose Return the upper triangular part of a matrix

Description Refer to the MATLAB `triu` reference page for more information.

Purpose Return the stored integer value of a `fi` object as a built-in `uint8`

Syntax `uint8(a)`

Description Fixed-point numbers can be represented as

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

or, equivalently,

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

The stored integer is the raw binary number, in which the binary point is assumed to be at the far right of the word.

`uint8(a)` returns the stored integer value of `fi` object `a` as a built-in `uint8`. If the stored integer word length is too big for a `uint8`, or if the stored integer is signed, the returned value saturates to a `uint8`.

See Also `int`, `int8`, `int16`, `int32`, `uint16`, `uint32`

uint16

Purpose Return the stored integer value of a `fi` object as a built-in `uint16`

Syntax `uint16(a)`

Description Fixed-point numbers can be represented as

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

or, equivalently,

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

The stored integer is the raw binary number, in which the binary point is assumed to be at the far right of the word.

`uint16(a)` returns the stored integer value of `fi` object `a` as a built-in `uint16`. If the stored integer word length is too big for a `uint16`, or if the stored integer is signed, the returned value saturates to a `uint16`.

See Also `int`, `int8`, `int16`, `int32`, `uint8`, `uint32`

Purpose Return the stored integer value of a `fi` object as a built-in `uint32`

Syntax `uint32(a)`

Description Fixed-point numbers can be represented as

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

or, equivalently,

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

The stored integer is the raw binary number, in which the binary point is assumed to be at the far right of the word.

`uint32(a)` returns the stored integer value of `fi` object `a` as a built-in `uint32`. If the stored integer word length is too big for a `uint32`, or if the stored integer is signed, the returned value saturates to a `uint32`.

See Also `int`, `int8`, `int16`, `int32`, `uint8`, `uint16`

uminus

Purpose	Negate the elements of a fi object array
Syntax	<code>uminus(a)</code>
Description	<code>uminus(a)</code> is called for the syntax <code>'-a'</code> when <code>a</code> is an object. <code>-a</code> negates the elements of <code>a</code> .
See Also	<code>plus</code> , <code>minus</code> , <code>mtimes</code> , <code>times</code>

Purpose

Unary plus

Description

Refer to the MATLAB arithmetic operators reference page for more information.

upperbound

Purpose Return upper bound of range of fi object

Syntax upperbound(a)

Description upperbound(a) returns the upper bound of the range of fi object a. If $L = \text{lowerbound}(a)$ and $U = \text{upperbound}(a)$, then $[L, U] = \text{range}(a)$.

See Also lowerbound, range

Purpose	Vertically concatenate two or more <code>fi</code> objects
Syntax	<code>c = vertcat(a,b,...)</code> <code>[a; b; ...]</code> <code>[a;b]</code>
Description	<p><code>c = vertcat(a,b,...)</code> is called for the syntax <code>[a; b; ...]</code> when any of <code>a</code>, <code>b</code>, ..., is a <code>fi</code> object.</p> <p><code>[a;b]</code> is the vertical concatenation of matrices <code>a</code> and <code>b</code>. <code>a</code> and <code>b</code> must have the same number of columns. Any number of matrices can be concatenated within one pair of brackets. N-D arrays are vertically concatenated along the first dimension. The remaining dimensions must match.</p> <p>Horizontal and vertical concatenation can be combined, as in <code>[1 2;3 4]</code>.</p> <p><code>[a b; c]</code> is allowed if the number of rows of <code>a</code> equals the number of rows of <code>b</code>, and if the number of columns of <code>a</code> plus the number of columns of <code>b</code> equals the number of columns of <code>c</code>.</p> <p>The matrices in a concatenation expression can themselves be formed via a concatenation, as in <code>[a b;[c d]]</code>.</p> <hr/> <p>Note The <code>fimath</code> and <code>numericType</code> objects of a concatenated matrix of <code>fi</code> objects <code>c</code> are taken from the leftmost <code>fi</code> object in the list <code>(a,b,...)</code></p> <hr/>
See Also	<code>horzcat</code>

voronoi

Purpose Create Voronoi diagram

Description Refer to the MATLAB voronoi reference page for more information.

Purpose Create n-dimensional Voronoi diagram

Description Refer to the MATLAB voronoin reference page for more information.

waterfall

Purpose Create waterfall plot

Description Refer to the MATLAB waterfall reference page for more information.

Purpose	Return the word length of a quantizer object
Syntax	<code>wordlength(q)</code>
Description	<code>wordlength(q)</code> returns the word length of the quantizer object <code>q</code> .
Examples	<pre>q = quantizer([16 15]); wordlength(q) ans = 16</pre>
See Also	<code>fi</code> , <code>fractionlength</code> , <code>exponentlength</code> , <code>numerictype</code> , <code>quantizer</code>

xlim

Purpose Set or query x-axis limits

Description Refer to the MATLAB `xlim` reference page for more information.

Purpose Set or query y-axis limits

Description Refer to the MATLAB `ylim` reference page for more information.

zlim

Purpose Set or query z-axis limits

Description Refer to the MATLAB `zlim` reference page for more information.

This glossary defines terms related to fixed-point data types and numbers. These terms may appear in some or all of the documents that describe products from The MathWorks that have fixed-point support.

arithmetic shift

Shift of the bits of a binary word for which the sign bit is recycled for each bit shift to the right. A zero is incorporated into the least significant bit of the word for each bit shift to the left. In the absence of overflows, each arithmetic shift to the right is equivalent to a division by 2, and each arithmetic shift to the left is equivalent to a multiplication by 2.

See also binary point, binary word, bit, logical shift, most significant bit

bias

Part of the numerical representation used to interpret a fixed-point number. Along with the slope, the bias forms the scaling of the number. 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

$$\text{slope} = \text{fractional slope} \times 2^{\text{exponent}}$$

See also fixed-point representation, fractional slope, integer, scaling, slope, [Slope Bias]

binary number

Value represented in a system of numbers that has two as its base and that uses 1's and 0's (bits) for its notation.

See also bit

binary point

Symbol in the shape of a period that separates the integer and fractional parts of a binary number. Bits to the left of the binary point are integer bits and/or sign bits, and bits to the right of the binary point are fractional bits.

See also binary number, bit, fraction, integer, radix point

binary point-only scaling

Scaling of a binary number that results from shifting the binary point of the number right or left, and which therefore can only occur by powers of two.

See also binary number, binary point, scaling

binary word

Fixed-length sequence of bits (1's and 0's). In digital hardware, numbers are stored in binary words. The way in which hardware components or software functions interpret this sequence of 1's and 0's is described by a data type.

See also bit, data type, word

bit

Smallest unit of information in computer software or hardware. A bit can have the value 0 or 1.

ceiling (round toward)

Rounding mode that rounds to the closest representable number in the direction of positive infinity. This is equivalent to the `ceil` mode in Fixed-Point Toolbox.

See also convergent rounding, floor (round toward), nearest (round toward), rounding, truncation, zero (round toward)

contiguous binary point

Binary point that occurs within the word length of a data type. For example, if a data type has four bits, its contiguous binary point must be understood to occur at one of the following five positions:

- .0000
- 0.000
- 00.00
- 000.0
- 0000.

See also data type, noncontiguous binary point, word length

convergent rounding

Rounding mode that rounds to the nearest allowable quantized value. Numbers that are exactly halfway between the two nearest allowable quantized values are rounded up only if the least significant bit (after rounding) would be set to 0.

See also ceiling (round toward), floor (round toward), nearest (round toward), rounding, truncation, zero (round toward)

data type

Set of characteristics that define a group of values. A fixed-point data type is defined by its word length, its fraction length, and whether it is signed or unsigned. A floating-point data type is defined by its word length and whether it is signed or unsigned.

See also fixed-point representation, floating-point representation, fraction length, word length

data type override

Parameter in the Fixed-Point Settings interface that allows you to set the output data type and scaling of fixed-point blocks on a system or subsystem level.

See also data type, scaling

exponent

Part of the numerical representation used to express a floating-point or fixed-point number.

1. Floating-point numbers are typically represented as

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

2. 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

$$\text{slope} = \text{fractional slope} \times 2^{\text{exponent}}$$

The exponent of a fixed-point number is equal to the negative of the fraction length:

$$\textit{exponent} = -1 \times \textit{fraction length}$$

See also bias, fixed-point representation, floating-point representation, fraction length, fractional slope, integer, mantissa, slope

fixed-point representation

Method for representing numerical values and data types that have a set range and precision.

1. Fixed-point numbers can be represented as

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

where the slope can be expressed as

$$\textit{slope} = \textit{fractional slope} \times 2^{\textit{exponent}}$$

The slope and the bias together represent the scaling of the fixed-point number.

2. Fixed-point data types can be defined by their word length, their fraction length, and whether they are signed or unsigned.

See also bias, data type, exponent, fraction length, fractional slope, integer, precision, range, scaling, slope, word length

floating-point representation

Method for representing numerical values and data types that can have changing range and precision.

1. Floating-point numbers can be represented as

$$\textit{real-world value} = \textit{mantissa} \times 2^{\textit{exponent}}$$

2. Floating-point data types are defined by their word length.

See also data type, exponent, mantissa, precision, range, word length

floor (round toward)

Rounding mode that rounds to the closest representable number in the direction of negative infinity.

See also ceiling (round toward), convergent rounding, nearest (round toward), rounding, truncation, zero (round toward)

fraction

Part of a fixed-point number represented by the bits to the right of the binary point. The fraction represents numbers that are less than one.

See also binary point, bit, fixed-point representation

fraction length

Number of bits to the right of the binary point in a fixed-point representation of a number.

See also binary point, bit, fixed-point representation, fraction

fractional slope

Part of the numerical representation used to express a fixed-point number. 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

$$\text{slope} = \text{fractional slope} \times 2^{\text{exponent}}$$

The term *slope adjustment* is sometimes used as a synonym for fractional slope.

See also bias, exponent, fixed-point representation, integer, slope

guard bits

Extra bits in either a hardware register or software simulation that are added to the high end of a binary word to ensure that no information is lost in case of overflow.

See also binary word, bit, overflow

integer

1. Part of a fixed-point number represented by the bits to the left of the binary point. The integer represents numbers that are greater than or equal to one.

2. Also called the "stored integer." The raw binary number, in which the binary point is assumed to be at the far right of the word. The integer is part of the numerical representation used to express a fixed-point number. Fixed-point numbers can be represented as

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

or

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

where the slope can be expressed as

$$\text{slope} = \text{fractional slope} \times 2^{\text{exponent}}$$

See also bias, fixed-point representation, fractional slope, integer, real-world value, slope

integer length

Number of bits to the left of the binary point in a fixed-point representation of a number.

See also binary point, bit, fixed-point representation, fraction length, integer

least significant bit (LSB)

Bit in a binary word that can represent the smallest value. The LSB is the rightmost bit in a big-endian-ordered binary word. The weight of the LSB is related to the fraction length according to

$$\text{weight of LSB} = 2^{-\text{fraction length}}$$

See also big-endian, binary word, bit, most significant bit

logical shift

Shift of the bits of a binary word, for which a zero is incorporated into the most significant bit for each bit shift to the right and into the least significant bit for each bit shift to the left.

See also arithmetic shift, binary point, binary word, bit, most significant bit

mantissa

Part of the numerical representation used to express a floating-point number. Floating-point numbers are typically represented as

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

See also exponent, floating-point representation

most significant bit (MSB)

Bit in a binary word that can represent the largest value. The MSB is the leftmost bit in a big-endian-ordered binary word.

See also binary word, bit, least significant bit

nearest (round toward)

Rounding mode that rounds to the closest representable number, with the exact midpoint rounded to the closest representable number in the direction of positive infinity. This is equivalent to the round mode in Fixed-Point Toolbox.

See also ceiling (round toward), convergent rounding, floor (round toward), rounding, truncation, zero (round toward)

noncontiguous binary point

Binary point that is understood to fall outside the word length of a data type. For example, the binary point for the following 4-bit word is understood to occur two bits to the right of the word length,

0000_ _.

thereby giving the bits of the word the following potential values:

$2^5 2^4 2^3 2^2$ _ _.

See also binary point, data type, word length

one's complement representation

Representation of signed fixed-point numbers. 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.

See also binary number, binary word, sign/magnitude representation, signed fixed-point, two's complement representation

overflow

Situation that occurs when the magnitude of a calculation result is too large for the range of the data type being used. In many cases you can choose to either saturate or wrap overflows.

See also saturation, wrapping

padding

Extending the least significant bit of a binary word with one or more zeros.

See also least significant bit

precision

1. Measure of the smallest numerical interval that a fixed-point data type and scaling can represent, determined by the value of the number's least significant bit. The precision is given by the slope, or the number of fractional bits. The term *resolution* is sometimes used as a synonym for this definition.

2. Measure of the difference between a real-world numerical value and the value of its quantized representation. This is sometimes called quantization error or quantization noise.

See also data type, fraction, least significant bit, quantization, quantization error, range, slope

Q format

Representation used by Texas Instruments to encode signed two's complement fixed-point data types. This fixed-point notation takes the form

$Qm.n$

where

- Q indicates that the number is in Q format.
- m is the number of bits used to designate the two's complement integer part of the number.

- n is the number of bits used to designate the two's complement fractional part of the number, or the number of bits to the right of the binary point.

In Q format notation, the most significant bit is assumed to be the sign bit.

See also binary point, bit, data type, fixed-point representation, fraction, integer, two's complement

quantization

Representation of a value by a data type that has too few bits to represent it exactly.

See also bit, data type, quantization error

quantization error

Error introduced when a value is represented by a data type that has too few bits to represent it exactly, or when a value is converted from one data type to a shorter data type. Quantization error is also called quantization noise.

See also bit, data type, quantization

radix point

Symbol in the shape of a period that separates the integer and fractional parts of a number in any base system. Bits to the left of the radix point are integer and/or sign bits, and bits to the right of the radix point are fraction bits.

See also binary point, bit, fraction, integer, sign bit

range

Span of numbers that a certain data type can represent.

See also data type, precision

real-world value

Stored integer value with fixed-point scaling applied. Fixed-point numbers can be represented as

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

or

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

where the slope can be expressed as

$$\text{slope} = \text{fractional slope} \times 2^{\text{exponent}}$$

See also integer

resolution

See **precision**

rounding

Limiting the number of bits required to express a number. One or more least significant bits are dropped, resulting in a loss of precision. Rounding is necessary when a value cannot be expressed exactly by the number of bits designated to represent it.

See also bit, ceiling (round toward), convergent rounding, floor (round toward), least significant bit, nearest (round toward), precision, truncation, zero (round toward)

saturation

Method of handling numeric overflow that represents positive overflows as the largest positive number in the range of the data type being used, and negative overflows as the largest negative number in the range.

See also overflow, wrapping

scaling

1. Format used for a fixed-point number of a given word length and signedness. The slope and bias together form the scaling of a fixed-point number.
2. Changing the slope and/or bias of a fixed-point number without changing the stored integer.

See also bias, fixed-point representation, integer, slope

shift

Movement of the bits of a binary word either toward the most significant bit ("to the left") or toward the least significant bit ("to the right"). Shifts to the right can be either logical, where the spaces emptied at the front of the word with each shift are filled in with zeros, or arithmetic, where the word is sign extended as it is shifted to the right.

See also arithmetic shift, logical shift, sign extension

sign bit

Bit (or bits) in a signed binary number that indicates whether the number is positive or negative.

See also binary number, bit

sign extension

Addition of bits that have the value of the most significant bit to the high end of a two's complement number. Sign extension does not change the value of the binary number.

See also binary number, guard bits, most significant bit, two's complement representation, word

sign/magnitude representation

Representation of signed fixed-point or floating-point numbers. In sign/magnitude representation, 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.

See also binary word, bit, fixed-point representation, floating-point representation, one's complement representation, sign bit, signed fixed-point, two's complement representation

signed fixed-point

Fixed-point number or data type that can represent both positive and negative numbers.

See also data type, fixed-point representation, unsigned fixed-point

slope

Part of the numerical representation used to express a fixed-point number. Along with the bias, the slope forms the scaling of a fixed-point number. 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

$$\text{slope} = \text{fractional slope} \times 2^{\text{exponent}}$$

See also bias, fixed-point representation, fractional slope, integer, scaling, [Slope Bias]

slope adjustment

See **fractional slope**

[Slope Bias]

Representation used to define the scaling of a fixed-point number.

See also bias, scaling, slope

stored integer

See **integer**

trivial scaling

Scaling that results in the real-world value of a number being simply equal to its stored integer value:

$$\text{real-world value} = \text{integer}$$

In [Slope Bias] representation, fixed-point numbers can be represented as

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

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:

$$\text{real-world value} = \text{integer} \times 2^{-\text{fraction length}} = \text{integer} \times 2^0$$

Scaling is always trivial for pure integers, such as `int8`, and also for the true floating-point types `single` and `double`.

See also bias, binary point, binary point-only scaling, fixed-point representation, fraction length, integer, least-significant bit, scaling, slope, [Slope Bias]

truncation

Rounding mode that drops one or more least significant bits from a number.

See also ceiling (round toward), convergent rounding, floor (round toward), nearest (round toward), rounding, zero (round toward)

two's complement representation

Common representation of signed fixed-point numbers. Negation using signed two's complement representation consists of a translation into one's complement followed by the binary addition of a one.

See also binary word, one's complement representation, sign/magnitude representation, signed fixed-point

unsigned fixed-point

Fixed-point number or data type that can only represent numbers greater than or equal to zero.

See also data type, fixed-point representation, signed fixed-point

word

Fixed-length sequence of binary digits (1's and 0's). In digital hardware, numbers are stored in words. The way hardware components or software functions interpret this sequence of 1's and 0's is described by a data type.

See also binary word, data type

word length

Number of bits in a binary word or data type.

See also binary word, bit, data type

wrapping

Method of handling overflow. Wrapping uses modulo arithmetic to cast a number that falls outside of the representable range the data type being used back into the representable range.

See also data type, overflow, range, saturation

zero (round toward)

Rounding mode that rounds to the closest representable number in the direction of zero. This is equivalent to the `fix` mode in Fixed-Point Toolbox.

See also ceiling (round toward), convergent rounding, floor (round toward), nearest (round toward), rounding, truncation

A

- abs function 11-2
- add function 11-5
- all function 11-7
- and function 11-8
- ANSI C
 - compared with `fi` objects 2-20
- any function 11-9
- area function 11-10
- arithmetic operations
 - fixed-point 2-8

B

- bar function 11-11
- barh function 11-12
- Bias property 9-12
- bin function 11-13
- bin property 9-2
- bin2num function 11-14
- binary conversions 2-23
- bitand function 11-16
- bitcmp function 11-17
- bitget function 11-18
- bitor function 11-19
- bitset function 11-20
- bitshift function 11-21
- bitxor function 11-22
- buffer function 11-23

C

- CastBeforeSum property 9-5
- casts
 - fixed-point 2-16
- clabel function 11-24
- comet function 11-25
- comet3 function 11-26
- compass function 11-27
- complex function 11-28

- complex multiplication
 - fixed-point 2-11
- coneplot function 11-29
- conj function 11-30
- contour function 11-31
- contour3 function 11-32
- contourc function 11-33
- contourf function 11-34
- convergent function 11-35
- copyobj function 11-36
- ctranspose function 11-37

D

- Data property 9-2
- DataType property 9-12
- DataTypeMode property 9-12
- dec function 11-38
- demons 1-7
- denormalmax function 11-39
- denormalmin function 11-40
- diag function 11-41
- disp function 11-42
- display preferences
 - setting 5-5
- display settings 1-5
- div function 11-43
- double function 11-46
- double property 9-2

E

- end function 11-47
- eps function 11-48
- eq function 11-49
- errorbar function 11-50
- etreeplot function 11-51
- exponentbias function 11-52
- exponentlength function 11-53
- exponentmax function 11-54

- exponentmin function 11-55
- ezcontour function 11-56
- ezcontourf function 11-57
- ezmesh function 11-58
- ezplot function 11-59
- ezplot3 function 11-60
- ezpolar function 11-61
- ezsurf function 11-62
- ezsurf function 11-63

F

- feather function 11-64
- fi function 11-65
- fi objects
 - constructing 3-2
 - properties
 - bin 9-2
 - Data 9-2
 - double 9-2
 - hex 9-3
 - int 9-3
 - NumericType 9-3
 - oct 9-4
- fimath function 11-72
- fimath objects 2-13
 - constructing 4-2
 - properties
 - CastBeforeSum 9-5
 - MaxProductWordLength 9-5
 - MaxSumWordLength 9-5
 - OverflowMode 9-5
 - ProductFractionLength 9-6
 - ProductMode 9-6
 - ProductWordLength 9-7
 - RoundMode 9-7
 - SumFractionLength 9-8
 - SumMode 9-8
 - SumWordLength 9-9
- fimath property 9-2
- FimathDisplay property 9-10
- find function 11-75
- fiobjects
 - properties
 - fimath 9-2
- fipref function 11-76
- fipref objects
 - constructing 5-2
 - properties
 - FimathDisplay 9-10
 - LoggingMode 9-10
 - NumberDisplay 9-11
 - NumericTypeDisplay 9-10
- fixed-point data
 - reading from workspace 8-2
 - writing to workspace 8-2
- fixed-point data types
 - addition 2-10
 - arithmetic operations 2-8
 - casts 2-16
 - complex multiplication 2-11
 - modular arithmetic 2-8
 - multiplication 2-11
 - overflow handling 2-5
 - precision 2-5
 - range 2-5
 - rounding 2-6
 - saturation 2-5
 - scaling 2-4
 - subtraction 2-10
 - two's complement 2-9
 - wrapping 2-5
- fixed-point run-time API 8-6
- fixed-point signal logging 8-6
- FixedExponent property 9-13
- format
 - rat 9-11
- Format property 9-16
- fplot function 11-78
- fractionlength function 11-79

FractionLength property 9-13
function
 line 11-118
functions
 abs 11-2
 add 11-5
 all 11-7
 and 11-8
 any 11-9
 area 11-10
 bar 11-11
 barh 11-12
 bin 11-13
 bin2num 11-14
 bitand 11-16
 bitcmp 11-17
 bitget 11-18
 bitor 11-19
 bitset 11-20
 bitshift 11-21
 bitxor 11-22
 buffer 11-23
 clabel 11-24
 comet 11-25
 comet3 11-26
 compass 11-27
 complex 11-28
 coneplot 11-29
 conj 11-30
 contour 11-31
 contour3 11-32
 contourc 11-33
 contourf 11-34
 convergent 11-35
 copyobj 11-36
 ctranspose 11-37
 dec 11-38
 denormalmax 11-39
 denormalmin 11-40
 diag 11-41
 disp 11-42
 div 11-43
 double 11-46
 end 11-47
 eps 11-48
 eq 11-49
 errorbar 11-50
 etreeplot 11-51
 exponentbias 11-52
 exponentlength 11-53
 exponentmax 11-54
 exponentmin 11-55
 ezcontour 11-56
 ezcontourf 11-57
 ezmesh 11-58
 ezplot 11-59
 ezplot3 11-60
 ezpolar 11-61
 ezsurf 11-62
 ezsurfc 11-63
 feather 11-64
 fi 11-65
 fimath 11-72
 find 11-75
 fipref 11-76
 fplot 11-78
 fractionlength 11-79
 ge 11-81
 get 11-82
 gplot 11-83
 gt 11-84
 hankel 11-85
 hex 11-86
 hex2num 11-87
 hist 11-88
 histc 11-89
 horzcat 11-90
 imag 11-91
 inspect 11-93
 int 11-94

int16 11-97
int32 11-98
int8 11-96
intmax 11-99
intmin 11-100
ipermute 11-101
iscolumn 11-102
isempty 11-103
isequal 11-104
isfi 11-105
isnumeric 11-107
isobject 11-109
isreal 11-111
isrow 11-112
isscalar 11-113
assigned 11-114
isvector 11-115
le 11-116
length 11-117
logical 11-119
loglog 11-120
lowerbound 11-121
lsb 11-122
lt 11-123
max 11-124
mesh 11-126
meshc 11-127
meshz 11-128
min 11-129
minus 11-130
mpy 11-131
mtimes 11-133
ndims 11-134
ne 11-135
noperations 11-137
not 11-136
noverflows 11-138
num2bin 11-139
num2hex 11-140
num2int 11-142
numberofelements 11-143
numerictype 11-144
nunderflows 11-149
oct 11-150
or 11-151
patch 11-152
pcolor 11-153
permute 11-154
plot 11-155
plot3 11-156
plotmatrix 11-157
plotyy 11-158
plus 11-159
polar 11-160
pow2 11-161
quantize 11-163
quantizer 11-166
quiver 11-170
quiver3 11-171
randquant 11-172
range 11-174
real 11-176
realmax 11-177
realmin 11-178
repmat 11-179
reset 11-182
reshape 11-183
rgbplot 11-184
ribbon 11-185
rose 11-186
round 11-187
savefigpref 11-189
scatter 11-190
scatter3 11-191
sdec 11-192
semilogx 11-193
semilogy 11-194
set 11-195
sign 11-197
single 11-198

size 11-199
slice 11-200
spy 11-201
squeeze 11-202
stairs 11-203
stem 11-204
stem3 11-205
streamribbon 11-206
streamslice 11-207
streamtube 11-208
stripscaling 11-209
sub 11-210
subsasgn 11-212
subsref 11-213
sum 11-214
surf 11-215
surfc 11-216
surf1 11-217
surfnorm 11-218
text 11-219
times 11-220
toeplitz 11-221
tostring 11-222
transpose 11-224
treemap 11-225
tril 11-226
trimesh 11-227
triplot 11-228
trisurf 11-229
triu 11-230
uint16 11-232
uint32 11-233
uint8 11-231
uminus 11-234
uplus 11-235
upperbound 11-236
vertcat 11-237
voronoi 11-238
voronoin 11-239
waterfall 11-240

wordlength 11-241
xlim 11-242
ylim 11-243
zlim 11-244

G

ge function 11-81
get function 11-82
gplot function 11-83
gt function 11-84

H

hankel function 11-85
help
 getting 1-3
hex function 11-86
hex property 9-3
hex2num function 11-87
hist function 11-88
histc function 11-89
horzcat function 11-90

I

imag function 11-91
inspect function 11-93
int function 11-94
int property 9-3
int16 function 11-97
int32 function 11-98
int8 function 11-96
interoperability
 fi objects with Filter Design Toolbox 8-12
 fi objects with Signal Processing
 Blockset 8-7
 fi objects with Simulink 8-2
intmax function 11-99
intmin function 11-100
ipermute function 11-101

iscolumn function 11-102
isempty function 11-103
isequal function 11-104
isfi function 11-105
isnumeric function 11-107
isobject function 11-109
isreal function 11-111
isrow function 11-112
isscalar function 11-113
assigned function 11-114
isvector function 11-115

L

le function 11-116
length function 11-117
line function 11-118
logging modes
 setting 5-7
LoggingMode property 9-10
logical function 11-119
loglog function 11-120
lowerbound function 11-121
lsb function 11-122
lt function 11-123

M

max function 11-124
Max property 9-17
MaxProductWordLength property 9-5
MaxSumWordLength property 9-5
mesh function 11-126
meshc function 11-127
meshz function 11-128
min function 11-129
Min property 9-17
minus function 11-130
Mode property 9-16
Model Explorer

 setting embedded.numericity
 properties 6-7
modular arithmetic 2-8
mpy function 11-131
mtimes function 11-133
multiplication
 fixed-point 2-11

N

ndims function 11-134
ne function 11-135
NOperations property 9-18
nopnerations function 11-137
not function 11-136
noverflows function 11-138
NOverflows property 9-18
num2bin function 11-139
num2hex function 11-140
num2int function 11-142
NumberDisplay property 9-11
numberofelements function 11-143
numericity function 11-144
numericity objects
 constructing 6-2
 properties
 Bias 9-12
 DataType 9-12
 DataTypeMode 9-12
 FixedExponent 9-13
 FractionLength 9-13
 Scaling 9-14
 setting in the Model Explorer 6-7
 Signed 9-14
 Slope 9-14
 SlopeAdjustmentFactor 9-14
 WordLength 9-15
 setting properties in the Model
 Explorer 6-7
NumericType property 9-3

NumericTypeDisplay property 9-10
 nunderflows function 11-149
 NUnderflows property 9-18

O

oct function 11-150
 oct property 9-4
 one's complement 2-10
 or function 11-151
 overflow handling 2-5
 compared with ANSI C 2-25
 OverflowMode property 9-5 9-18

P

padding 2-17
 patch function 11-152
 pcolor function 11-153
 permute function 11-154
 plot function 11-155
 plot3 function 11-156
 plotmatrix function 11-157
 plotyy function 11-158
 plus function 11-159
 polar function 11-160
 pow2 function 11-161
 precision
 fixed-point data types 2-5
 ProductFractionLength property 9-6
 ProductMode property 9-6
 ProductWordLength property 9-7
 properties
 Bias, numerictype objects 9-12
 bin, fi objects 9-2
 CastBeforeSum, fimath objects 9-5
 Data, fi objects 9-2
 DataType, numerictype objects 9-12
 DataTypeMode, numerictype objects 9-12
 double, fi objects 9-2

 fimath, fi objects 9-2
 FimathDisplay, fipref objects 9-10
 FixedExponent, numerictype objects 9-13
 Format, quantizers 9-16
 FractionLength, numerictype
 objects 9-13
 hex, fi objects 9-3
 int, fi objects 9-3
 LoggingMode, fipref objects 9-10
 Max, quantizers 9-17
 MaxProductWordLength, fimathobjects 9-5
 MaxSumWordLength, fimath objects 9-5
 Min, quantizers 9-17
 Mode, quantizers 9-16
 NOperations, quantizers 9-18
 NOverflows, quantizers 9-18
 NumberDisplay, fipref objects 9-11
 NumericType, fi objects 9-3
 NumericTypeDisplay, fipref objects 9-10
 NUnderflows, quantizers 9-18
 oct, fi objects 9-4
 OverflowMode, fimath objects 9-5
 OverflowMode, quantizers 9-18
 ProductFractionLength, fimath
 objects 9-6
 ProductMode, fimath objects 9-6
 ProductWordLength, fimath objects 9-7
 RoundMode, fimath objects 9-7
 RoundMode, quantizers 9-19
 Scaling, numerictype objects 9-14
 Signed, numerictype objects 9-14
 Slope, numerictype objects 9-14
 SlopeAdjustmentFactor, numerictype
 objects 9-14
 SumFractionLength, fimath objects 9-8
 SumMode, fimath objects 9-8
 SumWordLength, fimath objects 9-9
 WordLength, numerictype objects 9-15
 property values
 quantizer objects 7-4

Q

- quantize function 11-163
- quantizer function 11-166
- quantizer objects
 - constructing 7-2
 - property values 7-4
- quantizers
 - properties
 - Format 9-16
 - Max 9-17
 - Min 9-17
 - Mode 9-16
 - NOperations 9-18
 - NOverflows 9-18
 - NUnderflows 9-18
 - OverflowMode 9-18
 - RoundMode 9-19
- quiver function 11-170
- quiver3 function 11-171

R

- randquant function 11-172
- range
 - fixed-point data types 2-5
- range function 11-174
- rat format 9-11
- reading fixed-point data from workspace 8-2
- real function 11-176
- realmax function 11-177
- realmin function 11-178
- repmat function 11-179
- reset function 11-182
- reshape function 11-183
- rgbplot function 11-184
- ribbon function 11-185
- rose function 11-186
- round function 11-187
- rounding
 - fixed-point data types 2-6

- RoundMode property 9-7 9-19
- run-time API
 - fixed-point data 8-6

S

- saturation 2-5
- savefigpref function 11-189
- scaling 2-4
- Scaling property 9-14
- scatter function 11-190
- scatter3 function 11-191
- sdec function 11-192
- semilogx function 11-193
- semilogy function 11-194
- set function 11-195
- sign function 11-197
- signal logging
 - fixed-point 8-6
- Signed property 9-14
- single function 11-198
- size function 11-199
- slice function 11-200
- Slope property 9-14
- SlopeAdjustmentFactor property 9-14
- spy function 11-201
- squeeze function 11-202
- stairs function 11-203
- stem function 11-204
- stem3 function 11-205
- streamribbon function 11-206
- streamslice function 11-207
- streamtube function 11-208
- stripscaling function 11-209
- sub function 11-210
- subsasgn function 11-212
- subsref function 11-213
- sum function 11-214
- SumFractionLength property 9-8
- SumMode property 9-8

SumWordLength property 9-9
surf function 11-215
surfc function 11-216
surf1 function 11-217
surfnorm function 11-218

T

text function 11-219
times function 11-220
toeplitz function 11-221
tostring function 11-222
transpose function 11-224
treemap function 11-225
tril function 11-226
trimesh function 11-227
triplot function 11-228
trisurf function 11-229
triu function 11-230
two's complement 2-9

U

uint16 function 11-232
uint32 function 11-233
uint8 function 11-231
uminus function 11-234
unary conversions 2-22
uplus function 11-235

upperbound function 11-236

V

vertcat function 11-237
voronoi function 11-238
voronoin function 11-239

W

waterfall function 11-240
wordlength function 11-241
WordLength property 9-15
wrapping
 fixed-point data types 2-5
writing fixed-point data to workspace 8-2

X

xlim function 11-242

Y

ylim function 11-243

Z

zlim function 11-244