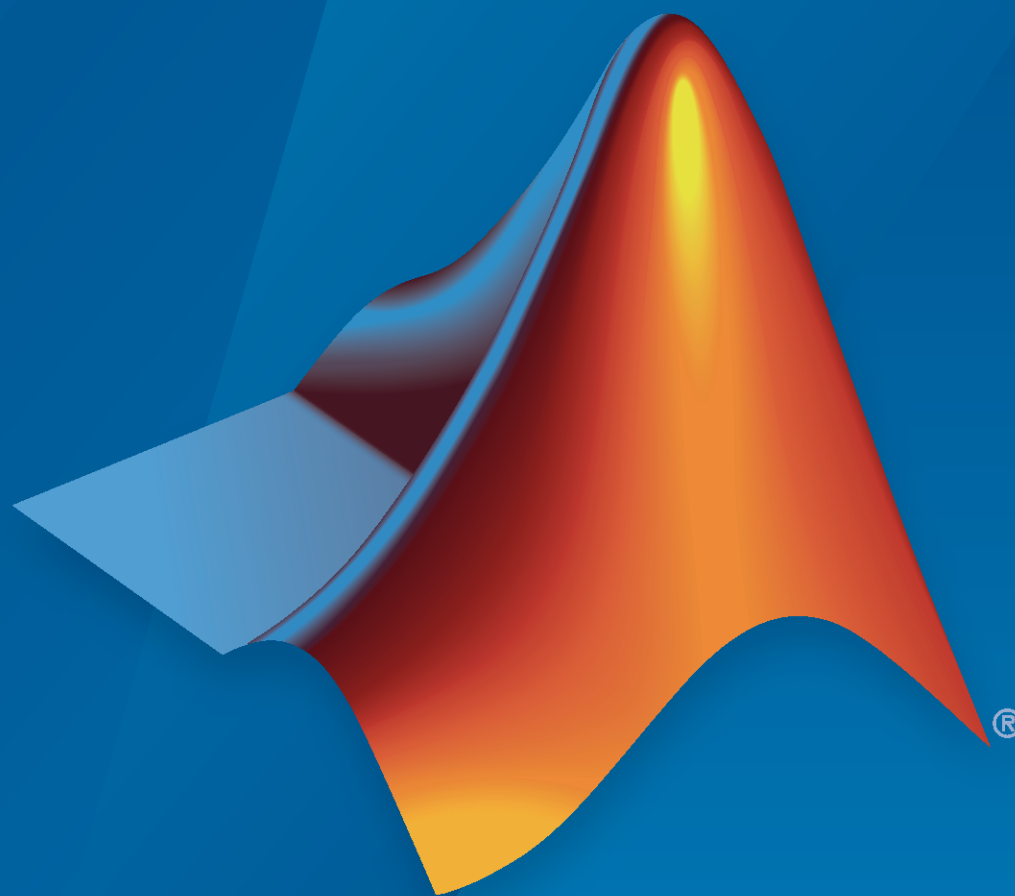


GPU Coder™

Getting Started Guide



MATLAB®

R2020b



How to Contact MathWorks



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



Phone: 508-647-7000



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

GPU Coder™ Getting Started Guide

© COPYRIGHT 2017–2020 by The MathWorks, Inc.

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

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

Trademarks

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

Patents

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

Revision History

September 2017	Online only	New for Version 1.0 (Release 2017b)
March 2018	Online only	Revised for Version 1.1 (Release 2018a)
September 2018	Online only	Revised for Version 1.2 (Release 2018b)
March 2019	Online only	Revised for Version 1.3 (Release 2019a)
September 2019	Online only	Revised for Version 1.4 (Release 2019b)
March 2020	Online only	Revised for Version 1.5 (Release 2020a)
September 2020	Online only	Revised for Version 2.0 (Release 2020b)

About GPU Coder

1

GPU Coder Product Description	1-2
Installing Prerequisite Products	1-3
MathWorks Products and Support Packages	1-3
Third-Party Hardware	1-3
Third-Party Software	1-3
Tips	1-5
Setting Up the Prerequisite Products	1-7
MEX Setup	1-7
Environment Variables	1-7
Verify Setup	1-10
The GPU Environment Check and Setup App	1-12
Hardware Setup	1-13
Board Settings	1-14
Workflow Checks	1-15
Environment Checks	1-16
GPU Code Generation Environment Check Report	1-17

Tutorials

2

Code Generation by Using the GPU Coder App	2-2
Learning Objectives	2-2
Tutorial Prerequisites	2-2
Example: The Mandelbrot Set	2-2
Tutorial Files	2-3
Run the Original MATLAB Code	2-4
Prepare MATLAB Code for Code Generation	2-5
Make the MATLAB Code Suitable for Code Generation	2-6
Generate Code by Using the GPU Coder App	2-7
Code Generation Using the Command Line Interface	2-17
Learning Objectives	2-17
Tutorial Prerequisites	2-17
Example: The Mandelbrot Set	2-17
Tutorial Files	2-18
Run the Original MATLAB Code	2-19
Make the MATLAB Code Suitable for Code Generation	2-20
Code Generation from the Command Line	2-21

Verify Correctness of the Generated Code	2-24
GPU Code Generation: The Mandelbrot Set	2-25

Verification

3

Verify Correctness of the Generated Code	3-2
Verify MEX Functions in the GPU Coder App	3-2
Verify MEX Functions at the Command Line	3-2
Code Verification Through Software-In-The-Loop	3-3
Numerical Differences Between CPU and GPU	3-4

About GPU Code Generation

4

GPU Programming Paradigm	4-2
GPU Code Generation Workflow	4-3

About GPU Coder

- “GPU Coder Product Description” on page 1-2
- “Installing Prerequisite Products” on page 1-3
- “Setting Up the Prerequisite Products” on page 1-7
- “The GPU Environment Check and Setup App” on page 1-12

GPU Coder Product Description

Generate CUDA code for NVIDIA GPUs

GPU Coder™ generates optimized CUDA® code from MATLAB® code and Simulink® models. The generated code includes CUDA kernels for parallelizable parts of your deep learning, embedded vision, and signal processing algorithms. For high performance, the generated code calls optimized NVIDIA® CUDA libraries, including TensorRT, cuDNN, cuFFT, cuSolver, and cuBLAS. The code can be integrated into your project as source code, static libraries, or dynamic libraries, and it can be compiled for desktops, servers, and GPUs embedded on NVIDIA Jetson, NVIDIA DRIVE, and other platforms. You can use the generated CUDA within MATLAB to accelerate deep learning networks and other computationally intensive portions of your algorithm. GPU Coder lets you incorporate handwritten CUDA code into your algorithms and into the generated code.

When used with Embedded Coder®, GPU Coder lets you verify the numerical behavior of the generated code via software-in-the-loop (SIL) and processor-in-the-loop (PIL) testing.

Installing Prerequisite Products

To use GPU Coder for CUDA C/C++ code generation, you must install the following products:

MathWorks Products and Support Packages

- MATLAB (required).
- MATLAB Coder™ (required).
- Parallel Computing Toolbox™ (required).
- Simulink (required for generating code from Simulink models).
- Computer Vision Toolbox™ (recommended).
- Deep Learning Toolbox™ (required for deep learning).
- Embedded Coder (recommended).
- Image Processing Toolbox™ (recommended).
- Simulink Coder (required for generating code from Simulink models).
- GPU Coder Interface for Deep Learning Libraries support package (required for deep learning).
- GPU Coder Support Package for NVIDIA GPUs (required for deployment to embedded targets such as NVIDIA Jetson and Drive).

For instructions on installing MathWorks® products, see the MATLAB installation documentation for your platform. If you have installed MATLAB and want to check which other MathWorks products are installed, enter `ver` in the MATLAB Command Window. To install the support packages, use Add-On Explorer in MATLAB.

If MATLAB is installed on a path that contains non 7-bit ASCII characters, such as Japanese characters, GPU Coder does not work because it cannot locate code generation library functions.

Third-Party Hardware

- NVIDIA GPU enabled for CUDA with compute capability 3.2 or higher and a compatible graphics driver. For more information, see [CUDA GPUs \(NVIDIA\)](#).

Deep learning applications in 8-bit integer precision requires a CUDA GPU with a compute capability of 6.1, 6.3 or higher. For half-precision (16-bit integer), a CUDA GPU with minimum compute capability of 7.0 is required.

- ARM® Mali graphics processor.

For the Mali device, GPU Coder supports C++ code generation for only deep learning networks.

Third-Party Software

Required

C/C++ Compiler:

Linux®	Windows®
GCC C/C++ compiler 6.3.x	Microsoft® Visual Studio® 2013

Linux®	Windows®
	Microsoft Visual Studio 2015
	Microsoft Visual Studio 2017
	Microsoft Visual Studio 2019

Optional

Standalone code (static library, dynamically linked library, or executable program) generation has additional software requirements. For CUDA MEX, the code generator uses the NVIDIA compiler and libraries installed with MATLAB.

Software Name	Information
CUDA toolkit	<p>GPU Coder has been tested with CUDA toolkit v9.x-v10.2.</p> <p>To download the CUDA toolkit, see CUDA Toolkit Archive (NVIDIA).</p> <p>It is recommended to select the default installation options that includes <code>nvcc</code> compiler, <code>cuFFT</code>, <code>cuBLAS</code>, <code>cuSOLVER</code>, Thrust libraries, and other tools.</p> <p>The NVIDIA <code>nvcc</code> compiler relies on tight integration with the host development environment, including the host compiler and runtime libraries. It is recommended that you follow the CUDA toolkit documentation for detailed information on compiler, libraries, and other platform specific requirements. See, CUDA Toolkit Documentation (NVIDIA)</p> <p>The <code>nvcc</code> compiler supports multiple versions of GCC and therefore you can generate CUDA code with other versions of GCC. However, there may be compatibility issues when executing the generated code from MATLAB as the C/C++ runtime libraries that are included with the MATLAB installation are compiled for GCC 6.3.</p>
NVIDIA CUDA deep neural network library (cuDNN) for NVIDIA GPUs	<p>cuDNN is a GPU-accelerated library of primitives for deep neural networks. For more information, see cuDNN (NVIDIA).</p> <p>GPU Coder has been tested with cuDNN v7.5.x.</p>
NVIDIA TensorRT	<p>TensorRT is a high performance inference optimizer and runtime library for NVIDIA GPUs. For more information, see TensorRT (NVIDIA).</p> <p>GPU Coder has been tested with TensorRT v7.0.x.</p>

Software Name	Information
ARM Compute Library	<p>The ARM Compute Library is a optimized computer vision and machine learning library for Mali GPUs. For more information, see Compute Library (ARM).</p> <p>GPU Coder has been tested with v19.05.</p> <p>This library must be installed on the ARM target hardware. Do not use a prebuilt library because it might be incompatible with the compiler on the ARM hardware. Instead, build the library from the source code. Build the library on either your host machine or directly on the target hardware. See instructions for building the library on GitHub®. You can also find information on building the library for CPUs in this post on MATLAB answers.</p> <p>When building the Compute Library, enable OpenCL support in the build options. See the ARM Compute Library documentation for instructions. OpenCL library (v1.2 or higher) on the ARM target hardware. See the ARM Compute Library documentation for version requirements. After the build is complete, rename the build folder containing the libraries as <code>lib</code>. Additionally, copy the OpenCL libraries present in the <code>build/openc1-1.2-stubs</code> folder into the <code>lib</code> folder. These steps are required so that the generated makefile can locate the libraries when building the generated code on the target hardware.</p>
Open Source Computer Vision Library (OpenCV)	<p>Required for deep learning examples. For more information, see OpenCV.</p> <p>For examples targeting NVIDIA GPUs on the host development computer, use OpenCV v3.1.0. For examples targeting ARM GPUs, use OpenCV v2.4.9 on the ARM target hardware.</p> <p>The examples require separate libs such as, <code>opencv_core.lib</code> and <code>opencv_video.lib</code>. You must download the OpenCV source and build the libraries. For more information, refer to the OpenCV documentation.</p>

Tips

- On Windows, a space or special character in the path to the tools, compilers, and libraries can create issues during the build process. You must install third-party software in locations that does

not contain spaces or change Windows settings to enable creation of short names for files, folders, and paths. For more information, see *Using Windows short names* solution in MATLAB Answers.

- The “Analyze Execution Profiles of the Generated Code” workflow depends on the `nvprof` tool from NVIDIA. From CUDA toolkit v10.1 onwards, NVIDIA restricts access to performance counters to only admin users. To enable GPU performance counters to be used by all users, see the instructions provided in https://developer.nvidia.com/nvidia-development-tools-solutions-ERR_NVGPUCTRPERM-permission-issue-performance-counters.
- GPU Coder does not support generating CUDA code by using CUDA toolkit version 8.
- Other versions of these deep learning libraries may have compatibility issues with all the features GPU Coder supports this release.
- On the target platforms, use the JetPack or the DriveInstall software appropriate for your board to install all the libraries. For more information, see “Install and Setup Prerequisites for NVIDIA Boards” (GPU Coder Support Package for NVIDIA GPUs).

See Also

Apps

GPU Coder | GPU Environment Check

Functions

`codegen` | `coder.checkGpuInstall`

Objects

`coder.gpuEnvConfig`

More About

- “Setting Up the Prerequisite Products” on page 1-7
- “The GPU Environment Check and Setup App” on page 1-12
- “Code Generation by Using the GPU Coder App” on page 2-2
- “Code Generation Using the Command Line Interface” on page 2-17
- “Code Generation for Deep Learning Networks by Using cuDNN”
- “Code Generation for Deep Learning Networks by Using TensorRT”
- “Code Generation for Deep Learning Networks Targeting ARM Mali GPUs”

Setting Up the Prerequisite Products

In this section...
“MEX Setup” on page 1-7
“Environment Variables” on page 1-7
“Verify Setup” on page 1-10

To use GPU Coder for CUDA C/C++ code generation, install the products specified in “Installing Prerequisite Products” on page 1-3.

MEX Setup

When generating CUDA MEX with GPU Coder, the code generator uses the NVIDIA compiler and libraries included with MATLAB. Depending on the operating system on your development computer, you only need to set up the MEX code generator.

Windows Systems

If you have multiple versions of Microsoft Visual Studio compilers for the C/C++ language installed on your Windows system, MATLAB selects one as the default compiler. If the selected compiler is not compatible with the version supported by GPU Coder, change the selection. For supported Microsoft Visual Studio versions, see “Installing Prerequisite Products” on page 1-3.

To change the default compiler, use the `mex -setup C++` command. When you call `mex -setup C++`, MATLAB displays a message with links to set up a different compiler. Select a link and change the default compiler for building MEX files. The compiler that you choose remains the default until you call `mex -setup C++` to select a different default. For more information, see “Change Default Compiler”. The `mex -setup C++` command changes only the C++ language compiler. You must also change the default compiler for C by using `mex -setup C`.

Linux Platform

MATLAB and the CUDA toolkit support only the GCC/G++ compiler for the C/C++ language on Linux platforms. For supported GCC/G++ versions, see “Installing Prerequisite Products” on page 1-3.

Environment Variables

Standalone code (static library, dynamically linked library, or executable program) generation has additional set up requirements. GPU Coder uses environment variables to locate the necessary tools, compilers, and libraries required for code generation.

Note On Windows, a space or special character in the path to the tools, compilers, and libraries can create issues during the build process. You must install third-party software in locations that does not contain spaces or change Windows settings to enable creation of short names for files, folders, and paths. For more information, see *Using Windows short names* solution in MATLAB Answers.

Platform	Variable Name	Description
Windows	CUDA_PATH	Path to the CUDA toolkit installation. For example: C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v10.0\
	NVIDIA_CUDNN	Path to the root folder of cuDNN installation. The root folder contains the bin, include, and lib subfolders. For example: C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\cuDNN\
	NVIDIA_TENSORRT	Path to the root folder of TensorRT installation. The root folder contains the bin, data, include, and lib subfolders. For example: C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\TensorRT\
	OPENCV_DIR	Path to the build folder of OpenCV on the host. This variable is required for building and running deep learning examples. For example: C:\Program Files\opencv\build
	PATH	Path to the CUDA executables. Generally, the CUDA toolkit installer sets this value automatically. For example: C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v10.0\bin Path to the cudnn.dll dynamic library. The name of this library may be different on your installation. For example: C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\cuDNN\bin Path to the nvinfer* dynamic libraries of TensorRT. The name of this library may be different on your installation. For example: C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\TensorRT\lib

Platform	Variable Name	Description
		<p>Path to the Dynamic-link libraries (DLL) of OpenCV. This variable is required for running deep learning examples.</p> <p>For example:</p> <p>C:\Program Files\opencv\build\x64\vc15\bin</p>
Linux	PATH	<p>Path to the CUDA toolkit executable.</p> <p>For example:</p> <p>/usr/local/cuda-10.0/bin</p>
		<p>Path to the OpenCV libraries. This variable is required for building and running deep learning examples.</p> <p>For example:</p> <p>/usr/local/lib/</p>
		<p>Path to the OpenCV header files. This variable is required for building deep learning examples.</p> <p>For example:</p> <p>/usr/local/include/opencv</p>
	LD_LIBRARY_PATH	<p>Path to the CUDA library folder.</p> <p>For example:</p> <p>/usr/local/cuda-10.0/lib64</p>
		<p>Path to the cuDNN library folder.</p> <p>For example:</p> <p>/usr/local/cuDNN/lib64/</p>
		<p>Path to the TensorRT library folder.</p> <p>For example:</p> <p>/usr/local/TensorRT/lib/</p>
<p>Path to the ARM Compute Library folder on the target hardware.</p> <p>For example:</p> <p>/usr/local/arm_compute/lib/</p> <p>Set LD_LIBRARY_PATH on the ARM target hardware.</p>		

Platform	Variable Name	Description
	NVIDIA_CUDNN	Path to the root folder of cuDNN library installation. For example: /usr/local/cuDNN/
	NVIDIA_TENSORRT	Path to the root folder of TensorRT library installation. For example: /usr/local/TensorRT/
	ARM_COMPUTELIB	Path to the root folder of the ARM Compute Library installation on the ARM target hardware. Set this value on the ARM target hardware. For example: /usr/local/arm_compute

Verify Setup

To verify that your development computer has all the tools and configuration needed for GPU code generation, use the `coder.checkGpuInstall` function. This function performs checks to verify if your environment has the all third-party tools and libraries required for GPU code generation. You must pass a `coder.gpuEnvConfig` object to the function. This function verifies the GPU code generation environment based on the properties specified in the given configuration object.

You can also use the equivalent GUI-based application that performs the same checks and can be launched using the command, **Check GPU Install**.

In the MATLAB Command Window, enter:

```
gpuEnvObj = coder.gpuEnvConfig;
gpuEnvObj.BasicCodegen = 1;
gpuEnvObj.BasicCodeexec = 1;
gpuEnvObj.DeepLibTarget = 'tensorrt';
gpuEnvObj.DeepCodeexec = 1;
gpuEnvObj.DeepCodegen = 1;
results = coder.checkGpuInstall(gpuEnvObj)
```

The output shown here is representative. Your results might differ.

```
Compatible GPU           : PASSED
CUDA Environment        : PASSED
  Runtime               : PASSED
  cuFFT                 : PASSED
  cuSOLVER              : PASSED
  cuBLAS                : PASSED
cuDNN Environment       : PASSED
TensorRT Environment    : PASSED
Basic Code Generation    : PASSED
Basic Code Execution     : PASSED
Deep Learning (TensorRT) Code Generation: PASSED
Deep Learning (TensorRT) Code Execution: PASSED
```

```
results =  
  
    struct with fields:  
  
        gpu: 1  
        cuda: 1  
        cudnn: 1  
        tensorrt: 1  
        basiccodegen: 1  
        basiccodeexec: 1  
        deepcodegen: 1  
        deepcodeexec: 1  
        tensorrtdatatype: 1  
        profiling: 0
```

See Also

Apps

GPU Coder | GPU Environment Check

Functions

codegen | coder.checkGpuInstall

Objects

coder.gpuEnvConfig

More About

- “Installing Prerequisite Products” on page 1-3
- “The GPU Environment Check and Setup App” on page 1-12
- “Code Generation by Using the GPU Coder App” on page 2-2
- “Code Generation Using the Command Line Interface” on page 2-17
- “Code Generation for Deep Learning Networks by Using cuDNN”
- “Code Generation for Deep Learning Networks by Using TensorRT”

The GPU Environment Check and Setup App

The **GPU Environment Check** app is an interactive tool to verify and set up the GPU code generation environment on your development computer and embedded hardware platforms such as the NVIDIA DRIVE and Jetson.


To start the app, in the MATLAB Command Window, enter:

```
gpusetup
```

Using the **GPU Environment Check** app, you can:

- Verify the host development computer environment for the NVIDIA compilers and libraries necessary for GPU code generation.
- Perform basic code generation and test the execution of the generated code on the GPU device in the host computer. The tests validate code execution by comparing the results with MATLAB simulation.
- Perform deep learning code generation and execution tests on the development computer. You can target the NVIDIA cuDNN or TensorRT libraries. Requires the GPU Coder Interface for Deep Learning Libraries support package.
- Connect to embedded NVIDIA boards such as DRIVE and Jetson to perform code generation and execution tests. Requires the GPU Coder Support Package for NVIDIA GPUs.
- Specify the location of the libraries and generate a MATLAB script that sets up the environment variables required by GPU Coder.

Note The `gpusetup` app generates a report file in the current folder. If you do not have write permissions in the current folder, before running the app, change the folder by using the MATLAB `cd` command.



GPU Code Generation Environment Check / Setup

Select Hardware Host (for MEX) ▼

Select GPU GPU0-Tesla K20c ▼

Workflow Checks

Basic Code Generation

Generate Code

Generate and Execute

SIL profiling

Deep Learning Code Generation

Generate Code





Generate and Execute


Target TensorRT ▼


Data Type FP32 ▼

Check

Environment Checks

CUDA Installation Path	<input style="width: 90%;" type="text" value="C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v10.1"/>	
cuDNN	<input style="width: 90%;" type="text" value="C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v10.1\cu"/>	
TensorRT	<input style="width: 90%;" type="text" value="C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v10.1\Te"/>	
NVTX Library Path	<input style="width: 90%;" type="text" value="C:\Program Files\NVIDIA Corporation\NvToolsExt\"/>	

 Run Checks

 Export Settings

Hardware Setup

The **Check/Setup** panel of the app provides drop-down lists that enable you to select a GPU device on the host development computer or hardware platforms such as the NVIDIA DRIVE and Jetson.

Option	Values	Description
Select Hardware	Host (for MEX)	<p>Perform code generation, code execution, and environment checks on the host development computer. The app generates CUDA MEX to perform tests.</p> <p>If your development computer has multiple GPU devices, use the Select GPU option to select an appropriate GPU device.</p>
	Drive	<p>Perform code generation and code execution checks on an NVIDIA DRIVE target platform.</p> <p>After installing the GPU Coder Support Package for NVIDIA GPUs, use the Board Settings panel to specify connection parameters for the target.</p>
	Jetson	<p>Perform code generation and code execution checks on an NVIDIA Jetson target platform.</p> <p>After installing the GPU Coder Support Package for NVIDIA GPUs, use the Board Settings panel to specify connection parameters for the target.</p>
Select GPU	<i>GPU<idx>-<device name></i>	<p>Select the GPU device to run tests. When there are multiple devices, the first device is the default.</p> <p>This option is visible only when the Select Hardware option is set to Host (for MEX). For DRIVE or Jetson hardware, use the GPU Device ID option available in the Board Settings panel to select a specific GPU device.</p>

Board Settings

Specify the connection parameters for hardware platforms such as the NVIDIA DRIVE and Jetson. The app uses the `jetson` or `drive` functions of the GPU Coder Support Package for NVIDIA GPUs to create a live hardware connection object. The support package software uses an SSH connection over TCP/IP to execute commands while building and running the generated CUDA code on the DRIVE or Jetson platforms. The target platform must be on the same network as the host computer.

Alternatively, you can use an Ethernet crossover cable to connect the board directly to the host computer. For more information on requirements, setup, and configuration steps for your NVIDIA boards, see “Install and Setup Prerequisites for NVIDIA Boards” (GPU Coder Support Package for NVIDIA GPUs).

Option	Description
Device Address	IP address or host name of the hardware. For example, 169.254.0.2 or gpuCoder-tegratx2-1 To use the host name, you must connect an Ethernet cable to the Ethernet port of the board. Then, use Linux commands to configure the hardware IP address and associate the host name with the IP address.
Username	Valid Linux user name for the operating system on the board.
Password	Valid password for the Linux user name specified.
Execution Timeout	Specify the time in seconds that the app waits for before validating the execution tests on the target. The default value is 10 seconds.
GPU Device ID	In a multi GPU environment such as NVIDIA Drive platforms, specify the CUDA GPU device to target.

Workflow Checks

There are two types of workflow checks that you can perform by using the app:

Basic code generation and execution tests on the development computer. These tests validate code execution by comparing the results with MATLAB simulation.

Option	Description
Generate Code	Test basic code generation and building. This test requires a valid CUDA code generation environment on the specified hardware.
Generate Code and Execute	Test basic code generation, building, and execution on the device in Specified Hardware . This test requires a valid CUDA code generation environment and GPU device on the specified hardware.
SIL Profiling	Perform basic SIL profiling tests on the host computer.

Deep learning code generation and execution tests on the development computer. You can target the cuDNN or TensorRT libraries.

Option	Description
Generate Code	Test deep learning code generation and building. This test requires a valid CUDA code generation environment on the specified hardware.
Generate Code and Execute	Test deep learning code generation, building, and execution on the device in Specified Hardware . This test requires a valid CUDA code generation environment and GPU device on the specified hardware.
Target	Specify the deep learning library to generate code for. Valid options are cuDNN or TensorRT.
Data Type Check	Specify the precision of the tensor data type input to the network or the tensor output of a layer. To perform inference in 32-bit floats, use 'FP32'. For half-precision, use 'FP16'. For 8-bit integer, use 'INT8'. Default value is 'FP32'. INT8 precision requires a CUDA GPU with minimum compute capability of 6.1. FP16 precision requires a CUDA GPU with minimum compute capability of 7.0.

Environment Checks

Specify the location of the libraries for checking the CUDA development environment on your host computer. Generate a MATLAB script `gpuEnvSettings.m` that sets up the environment variables required by GPU Coder. For more information, see “Setting Up the Prerequisite Products” on page 1-7.

Option	Description
CUDA Installation Path	Path to the CUDA toolkit installation. For example: <code>/usr/local/cuda-10.1/bin</code>
cuDNN	Path to the root folder of cuDNN library installation. For example: <code>/usr/local/cuDNN/</code>
TensorRT	Path to the root folder of TensorRT library installation. For example: <code>/usr/local/TensorRT/</code>

Option	Description
NVTX Library Path	Path to the nvtx libraries required for profiling. To enable this item, select SIL Profiling . On a standard CUDA toolkit installation, this path is usually the CUDA library folder. For example: <code>/usr/local/cuda-10.1/lib64</code>

GPU Code Generation Environment Check Report

When you select **Run Checks**, the `gpucoderSetup` app performs environment, code generation, and other checks based on the settings that you have selected. It then generates `gpucoderSetupReport` report that indicates if a test has passed and provides additional information for tests that have failed. A HTML version of the report of the same name is created in the current folder.

GPU Code Generation Environment Check Report

Results for Host GPU : TITAN Xp

ENVIRONMENT CHECKS

Check	Result	Message
Compatible GPU	✓	
CUDA Environment a. Runtime ✓ b. cuFFT ✓ c. cuSOLVER ✓ d. cuBLAS ✓	✓	
cuDNN Environment	✓	
TensorRT Environment	✓	
Profiling Environment	✓	

BASIC CODE GENERATION CHECKS

Check	Result	Message
Basic Code Generation	✓	
Basic Code Execution	✓	

DEEP LEARNING CODE GENERATION CHECKS - TensorRT

See Also

Apps

GPU Coder | GPU Environment Check

Functions

codegen | coder.checkGpuInstall

Objects

coder.gpuEnvConfig

More About

- “Installing Prerequisite Products” on page 1-3
- “Setting Up the Prerequisite Products” on page 1-7
- “Code Generation by Using the GPU Coder App” on page 2-2
- “Code Generation Using the Command Line Interface” on page 2-17
- “Code Generation for Deep Learning Networks by Using cuDNN”
- “Code Generation for Deep Learning Networks by Using TensorRT”

Tutorials

- “Code Generation by Using the GPU Coder App” on page 2-2
- “Code Generation Using the Command Line Interface” on page 2-17
- “GPU Code Generation: The Mandelbrot Set” on page 2-25

Code Generation by Using the GPU Coder App

The easiest way to create CUDA kernels is to place the `coder.gpu.kernel fun` pragma into your primary MATLAB function. The primary function is also known as the top-level or entry-point function. When GPU Coder encounters the `kernel fun` pragma, it attempts to parallelize all the computation within this function and then maps it to the GPU. For more information about GPU kernels, see “GPU Programming Paradigm” on page 4-2.

Learning Objectives

In this tutorial, you learn how to:

- Prepare your MATLAB code for CUDA code generation by using the `kernel fun` pragma.
- Create and set up a GPU Coder project.
- Define function input properties.
- Check for code generation readiness and run-time issues.
- Specify code generation properties.
- Generate CUDA C code by using the GPU Coder app.

Tutorial Prerequisites

This tutorial requires the following products:

- MATLAB
- MATLAB Coder
- GPU Coder
- C compiler
- NVIDIA GPU enabled for CUDA
- CUDA toolkit and driver
- Environment variables for the compilers and libraries. For more information, see “Environment Variables” on page 1-7.

Example: The Mandelbrot Set

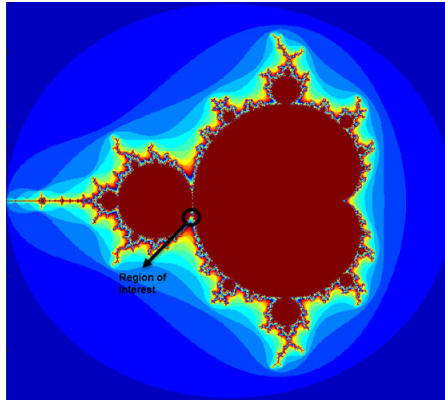
Description

You do not have to be familiar with the algorithm in the example to complete the tutorial.

The Mandelbrot set is the region in the complex plane consisting of the values z_0 for which the trajectories defined by this equation remain bounded at $k \rightarrow \infty$.

$$z_{k+1} = z_k^2 + z_0, \quad k = 0, 1, \dots$$

The overall geometry of the Mandelbrot set is shown in the figure. This view does not have the resolution to show the richly detailed structure of the fringe just outside the boundary of the set. At increasing magnifications, the Mandelbrot set exhibits an elaborate boundary that reveals progressively finer recursive detail.



Algorithm

For this tutorial, pick a set of limits that specify a highly zoomed part of the Mandelbrot set in the valley between the main cardioid and the p/q bulb to its left. A 1000x1000 grid of real parts (x) and imaginary parts (y) is created between these two limits. The Mandelbrot algorithm is then iterated at each grid location. An iteration number of 500 is enough to render the image in full resolution.

```
maxIterations = 500;
gridSize = 1000;
xlim = [-0.748766713922161, -0.748766707771757];
ylim = [0.123640844894862, 0.123640851045266];
```

This tutorial uses an implementation of the Mandelbrot set by using standard MATLAB commands running on the CPU. This implementation is based on the code provided in the Experiments with MATLAB e-book by Cleve Moler. This calculation is vectorized such that every location is updated simultaneously.

Tutorial Files

Create a MATLAB function called `mandelbrot_count.m` with the following lines of code. This code is a baseline vectorized MATLAB implementation of the Mandelbrot set. For every point ($xGrid, yGrid$) in the grid, it calculates the iteration index count at which the trajectory defined by the equation reaches a distance of 2 from the origin. It then returns the natural logarithm of count, which is used generate the color coded plot of the Mandelbrot set. Later in this tutorial, you modify this file to make it suitable for code generation.

```
function count = mandelbrot_count(maxIterations,xGrid,yGrid)
% mandelbrot computation

z0 = xGrid + 1i*yGrid;
count = ones(size(z0));

z = z0;
for n = 0:maxIterations
    z = z.*z + z0;
    inside = abs(z)<=2;
    count = count + inside;
end
count = log(count);
```

Create a MATLAB script called `mandelbrot_test.m` with the following lines of code. The script generates a 1000 x 1000 grid of real parts (x) and imaginary parts (y) between the limits specified by `xlim` and `ylim`. It also calls the `mandelbrot_count` function and plots the resulting Mandelbrot set.

```
maxIterations = 500;
gridSize = 1000;
xlim = [-0.748766713922161, -0.748766707771757];
ylim = [0.123640844894862, 0.123640851045266];

x = linspace(xlim(1),xlim(2),gridSize);
y = linspace(ylim(1),ylim(2),gridSize);
[xGrid,yGrid] = meshgrid(x,y);

%% Mandelbrot computation in MATLAB
count = mandelbrot_count(maxIterations,xGrid,yGrid);

% Show
figure(1)
imagesc(x,y,count);
colormap([jet();flipud(jet());0 0 0]);
axis off
title('Mandelbrot set with MATLAB');
```

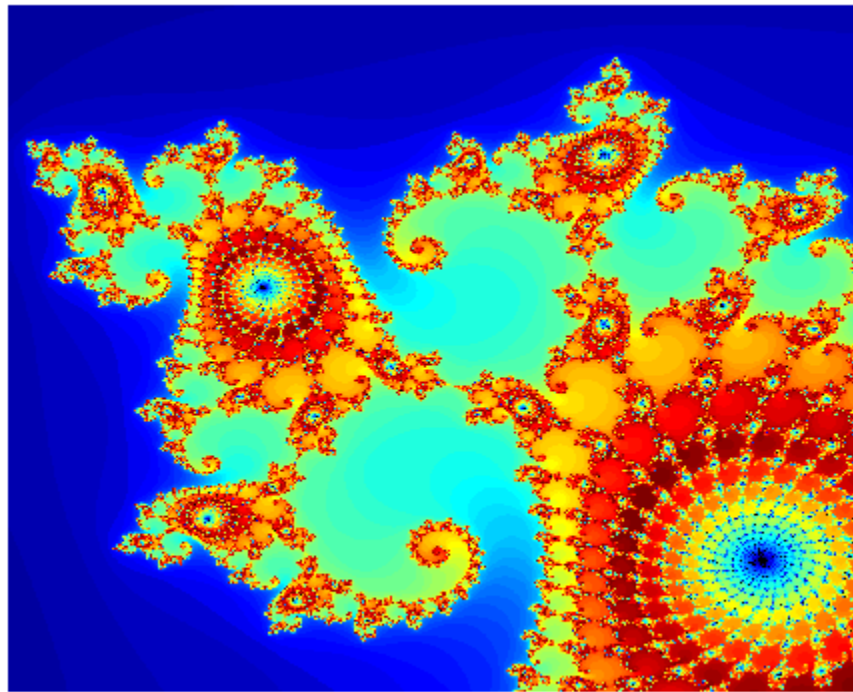
Run the Original MATLAB Code

Run the Mandelbrot Example

Before making the MATLAB version of the Mandelbrot set algorithm suitable for code generation, you can test the functionality of the original code.

- 1 Change the current MATLAB working folder to the location that contains `mandelbrot_count.m` and `mandelbrot_test.m`. GPU Coder places generated code in this folder. Change your current working folder if you do not have full access to this folder.
- 2 Run the `mandelbrot_test` script.

The test script runs and shows the geometry of the Mandelbrot within the boundary set by the variables `xlim` and `ylim`.

Mandelbrot Set with MATLAB

Prepare MATLAB Code for Code Generation

Before you generate code with GPU Coder, check for coding issues in the original MATLAB code.

Check for Issues at Design Time

There are two tools that help you detect code generation issues at design time:

- Code Analyzer tool
- Code generation readiness tool

The Code Analyzer is a tool incorporated into the MATLAB Editor that continuously checks your code as you enter it. The Code Analyzer reports issues and recommends modifications to maximize performance and maintainability of your code. To identify the warnings and errors specific to code generation from your MATLAB code, add the `%#codegen` directive to your MATLAB file. For more information, see Code Analyzer preferences.

Note The Code Analyzer does not detect all code generation issues. After eliminating the errors or warnings that the Code Analyzer detects, compile your code with GPU Coder to determine if the code has other compliance issues.

The code generation readiness tool screens the MATLAB code for features and functions that are not supported for code generation. This tool provides a report that lists issues and recommendations for

making the MATLAB code suitable for code generation. You can access the code generation readiness tool in these ways:

- In the current folder browser — right-click the MATLAB file that contains the entry-point function.
- At the command line — by using the `coder.screener()` function.
- In the GPU Coder app — after specifying the entry-point files, the app runs the Code Analyzer and the code generation readiness tool.

Check for Issues at Code Generation Time

You can use GPU Coder to check for issues at code generation time. When GPU Coder detects errors or warnings, it generates an error report that describes the issues and provides links to the problematic MATLAB code. For more information, see “Code Generation Reports”.

Make the MATLAB Code Suitable for Code Generation

To begin the process of making your MATLAB code suitable for code generation, use the file `mandelbrot_count.m`.

- 1 Set your MATLAB current folder to the work folder that contains your files for this tutorial.
- 2 In the MATLAB Editor, open `mandelbrot_count.m`. The Code Analyzer message indicator at the top right corner of the MATLAB Editor is green. The analyzer did not detect errors, warnings, or opportunities for improvement in the code.
- 3 After the function declaration, add the `%#codegen` directive to turn on the error checking that is specific to code generation.

```
function count = mandelbrot_count(maxIterations,xGrid,yGrid) %#codegen
```

The Code Analyzer message indicator remains green, indicating that it has not detected any code generation issues.

- 4 To map the `mandelbrot_count` function to a CUDA kernel, modify the original MATLAB code by placing the `coder.gpu.kernelfun` pragma in the body of the function.

```
function count = mandelbrot_count(maxIterations,xGrid,yGrid) %#codegen
% Add kernelfun pragma to trigger kernel creation
coder.gpu.kernelfun;

% mandelbrot computation
z0 = xGrid + 1i*yGrid;
count = ones(size(z0));

z = z0;
for n = 0:maxIterations
    z = z.*z + z0;
    inside = abs(z)<=2;
    count = count + inside;
end
count = log(count);
```

If you use the `coder.gpu.kernelfun` pragma, GPU Coder attempts to map the computations in the function `mandelbrot_count` to the GPU.

- 5 Save the file. You are now ready to compile your code by using the GPU Coder app.

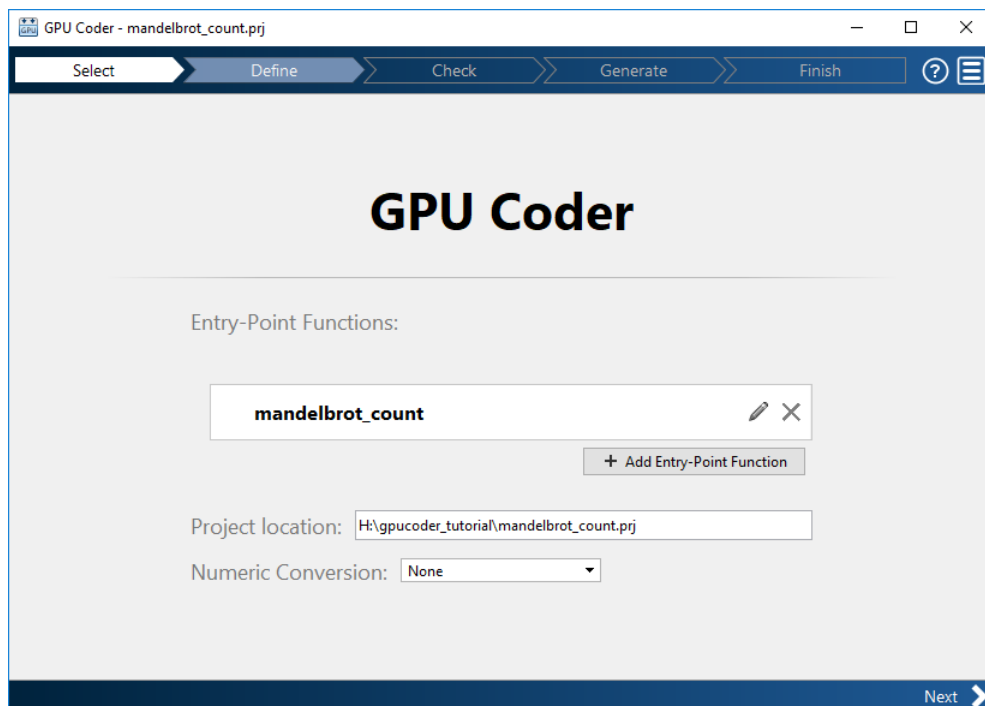
Generate Code by Using the GPU Coder App

Open the GPU Coder App

On the MATLAB toolstrip **Apps** tab, under **Code Generation**, click the GPU Coder app icon. You can also open the app by typing `gpcoder` in the MATLAB Command Window. The app opens the **Select** source files page.

Select Source Files

- 1 On the **Select** source files page, enter or select the name of the primary function, `mandelbrot_count`. The primary function is also known as the top-level or entry-point function. The app creates a project with the default name `mandelbrot_count.prj` in the current folder.



- 2 Click **Next** and go to the **Define Input Types** step. The app analyzes the function for coding issues and code generation readiness. If the app identifies issues, it opens the **Review Code Generation Readiness** page where you can review and fix issues. In this example, because the app does not detect issues, it opens the **Define Input Types** page.

Define Input Types

The code generator must determine the data types of all the variables in the MATLAB files at compile time. Therefore, you must specify the data types of all the input variables. You can specify the input data types in one of these two ways:

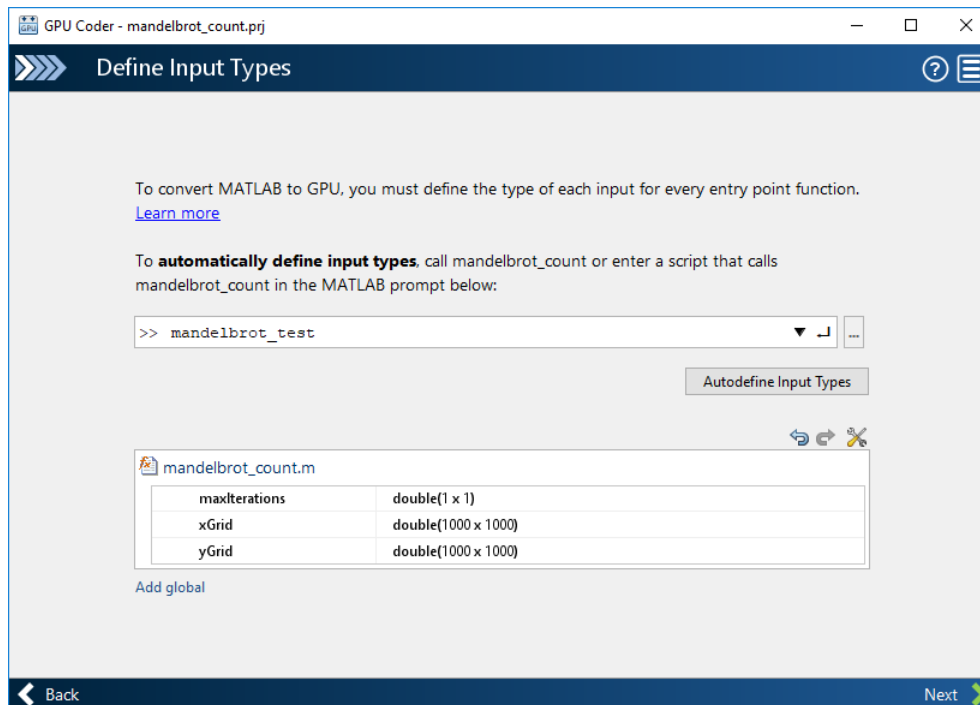
- Provide a test file that calls the project entry-point functions. The GPU Coder app can infer the input argument types by running the test file.
- Enter the input types directly.

For more information about input specifications, see “Input Specification”.

In this example, to define the properties of the inputs `maxIterations`, `xGrid`, and `yGrid`, specify the test file `mandelbrot_test.m`:

- 1 Enter or select the test file `mandelbrot_test.m`.
- 2 Click **Autodefine Input Types**.

The test file `mandelbrot_test.m` calls the entry-point function, `mandelbrot_count.m` with the expected input types. The app infers that the input `maxIterations` is `double(1x1)` and the inputs `xGrid` and `yGrid` are `double(1000x1000)`.



- 3 Click **Next** go to the **Check for Run-Time Issues** step.

Check for Run-Time Issues

The **Check for Run-Time Issues** step generates a MEX file from your entry-point functions, runs the MEX function, and reports issues. This step is optional. However, it is a best practice to perform this step. Using this step, you can detect and fix defects that are harder to diagnose in the generated GPU code.

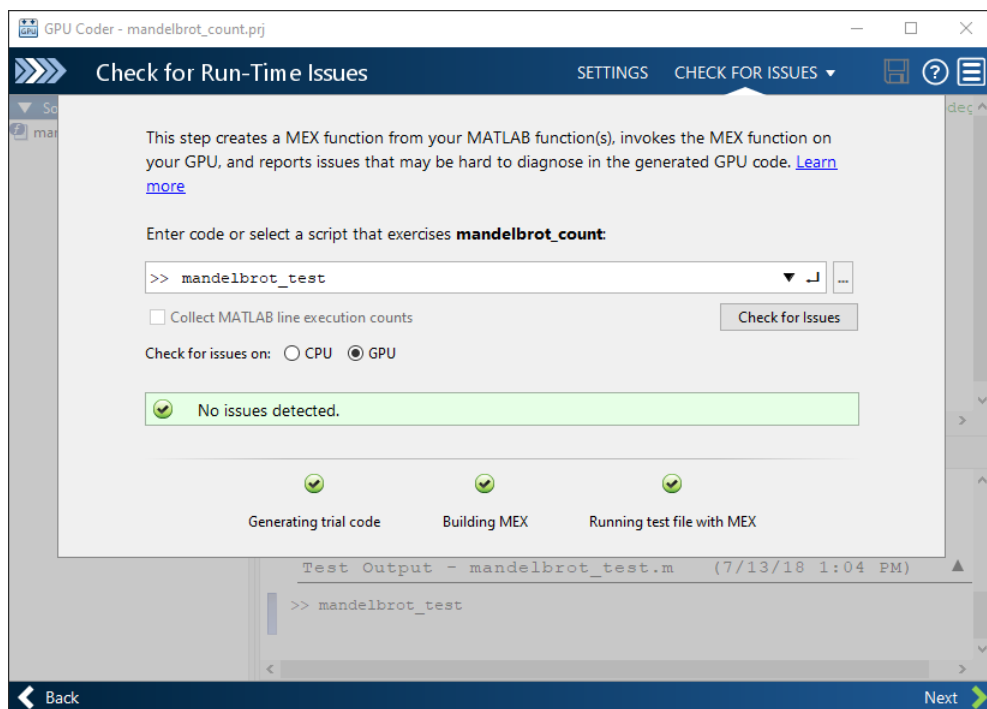
GPU Coder provides the option to perform GPU-specific checks at this point. When you select this option, GPU Coder generates CUDA C code and a MEX file from your entry-point functions, runs the MEX function, and reports issues. Some of the GPU-specific run-time checks include:

- Checks for register spills.
- Stack size conformance checks.

Note There may be certain MATLAB constructs in your code that cause the **Check for Run-Time Issues** to fail CPU-specific checks but pass the GPU-specific checks.

- 1 To open the **Check for Run-Time Issues** dialog box, click the **Check for Issues** arrow.
- 2 In the **Check for Run-Time Issues** dialog box, specify a test file or enter code that calls the entry-point function with example inputs. For this example, use the test file `mandelbrot_test.m` that you used to define the input types.
- 3 To enable GPU-specific checks, select the **GPU** option button. Click **Check for Issues**.

The app generates a MEX function. It runs the test script `mandelbrot_test` replacing calls to `mandelbrot_count` with calls to the generated MEX. If the app detects issues during the MEX function generation or execution, it provides warning and error messages. You can click these messages to navigate to the problematic code and fix the issue. In this example, the app does not detect issues. The MEX function has the same functionality as the original `mandelbrot_count` function.

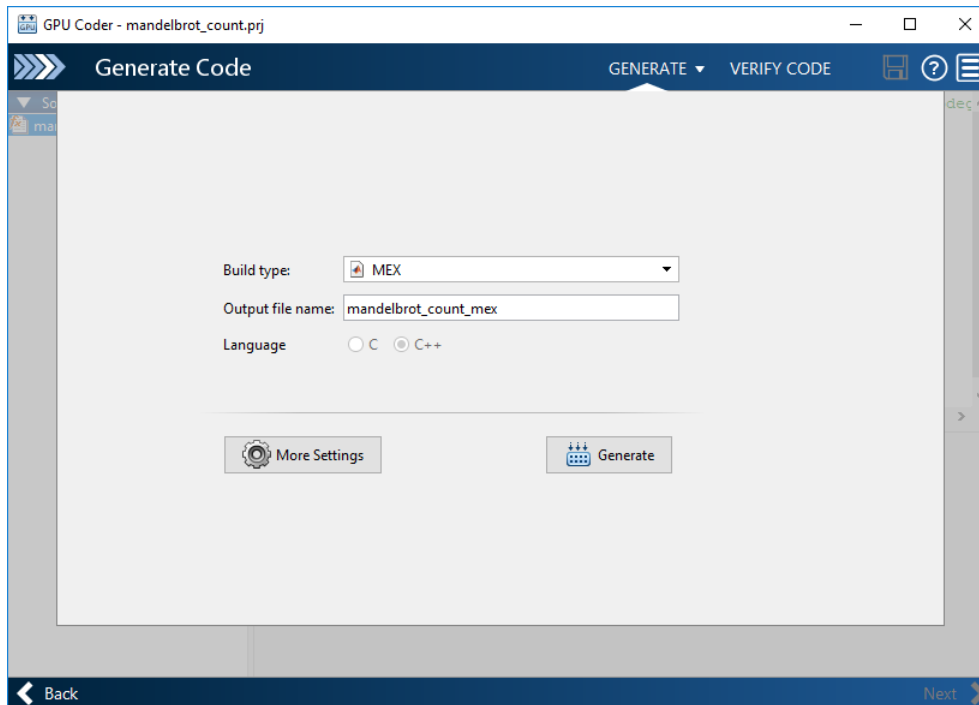


Note There may be certain MATLAB constructs in your code that cause the **Check for Run-Time Issues** to fail CPU-specific checks but pass the GPU-specific checks.

- 4 Click **Next** go to the **Generate Code** step.

Generate CUDA Code

- 1 To open the **Generate** dialog box, click the **Generate** arrow.

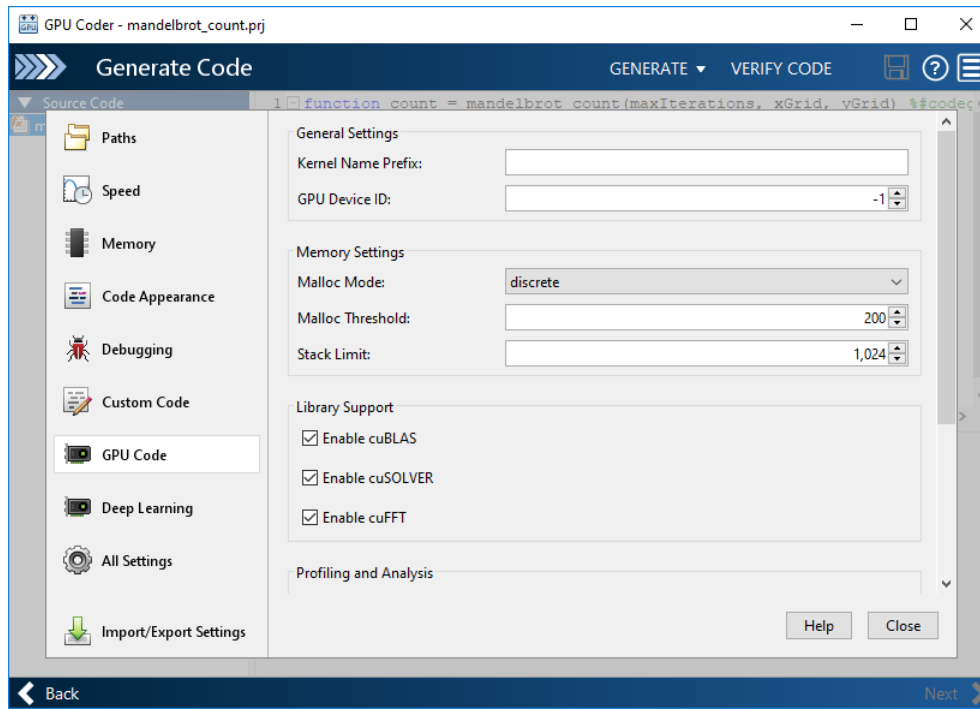


- 2 In the **Generate** dialog box, you can select the type of build that you want GPU Coder to perform. The available options are listed in this table.

Build Type	Description
Source code	CUDA C Source code to integrate with an external project.
MEX	Compiled code to run inside MATLAB.
Static Library	Binary library for static linking with an external project.
Dynamic Library	Binary library for dynamic linking with an external project.
Executable	Standalone program (requires a separate main file written in C).

For this tutorial, set **Build type** to MEX (.mex). By generating a MEX output, you can check the correctness of the generated CUDA code from within MATLAB. The MEX build type does not require additional settings like **Toolchain** and **Hardware Board**. It also does not provide the option to generate only the source code. GPU Coder can automatically select an available CUDA toolchain as long as the “Environment Variables” on page 1-7 are set properly.

To view advanced options, select **More Settings**. To the **Compiler Flags** option, add `--fmad=false`. This flag, when passed to the `nvcc`, instructs the compiler to disable Floating-point Multiply-add (FMAD) optimization. This option is set to prevent numerical mismatch in the generated code because of architectural differences between the CPU and the GPU. For more information, see “Numerical Differences Between CPU and GPU” on page 3-4.



This table describes the settings specific to GPU Coder.

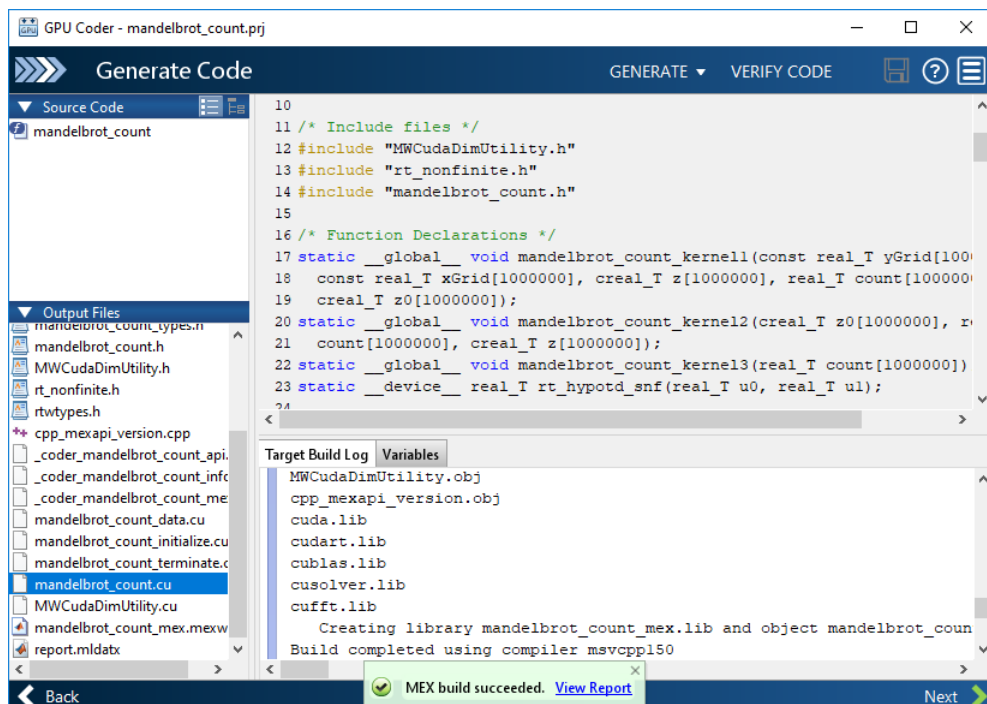
GPU Coder Configuration Properties

UI Setting	Value Type	Description
Kernel Name Prefix	String	Specify custom name prefix for kernel names in the generated code. For example, entering 'CUDA_' creates kernels with names CUDA_kernel1, CUDA_kernel2, and so on. If no name is provided, GPU Coder prepends the kernel name with the name of the entry-point function. Kernel names can contain upper-case letters, lowercase letters, digits 0-9, and underscore character _. GPU Coder removes unsupported characters from the kernel names and appends <code>alpha</code> to prefixes that do not begin with an alphabetic letter.
Malloc Mode	Enumerated 'Discrete' 'Unified'	Selects the type of GPU memory allocation: Discrete or Unified.
Malloc Threshold	Integer	Size above which the private variables are allocated on the heap instead of the stack.
Stack Limit	Integer	Available stack limit per GPU thread.
Enable cuSOLVER	Boolean 'True' 'False'	Allows GPU Coder to utilize cuSOLVER library calls where appropriate.
Benchmarking	Boolean 'True' 'False'	Generates CUDA code with benchmarking options such as <code>cudaEvent</code> API to accurately time kernel, <code>memcpy</code> , and other events.
Safe Build	Boolean 'True' 'False'	Generates code with error-checking for CUDA API and kernel calls.
Minimum Compute Capability	Enumerated '3.2' '3.5' '3.7' '5.0' '5.2' '5.3' '6.0' '6.1'	Select the minimum compute capability for code generation. The compute capability identifies the features supported by the GPU hardware and is used by applications at run time to determine which hardware features, instructions are available on the present GPU. If you specify custom compute capability, GPU Coder ignores this setting.
Custom Compute Capability	String	Specify the name of the NVIDIA virtual GPU architecture for which the CUDA input files must be compiled. For example, to specify a virtual architecture type - <code>arch=compute_50</code> . You can specify a real architecture using <code>arch=sm_50</code> . For more information, see the <i>Options for Steering GPU Code Generation</i> topic in the CUDA toolkit documentation.

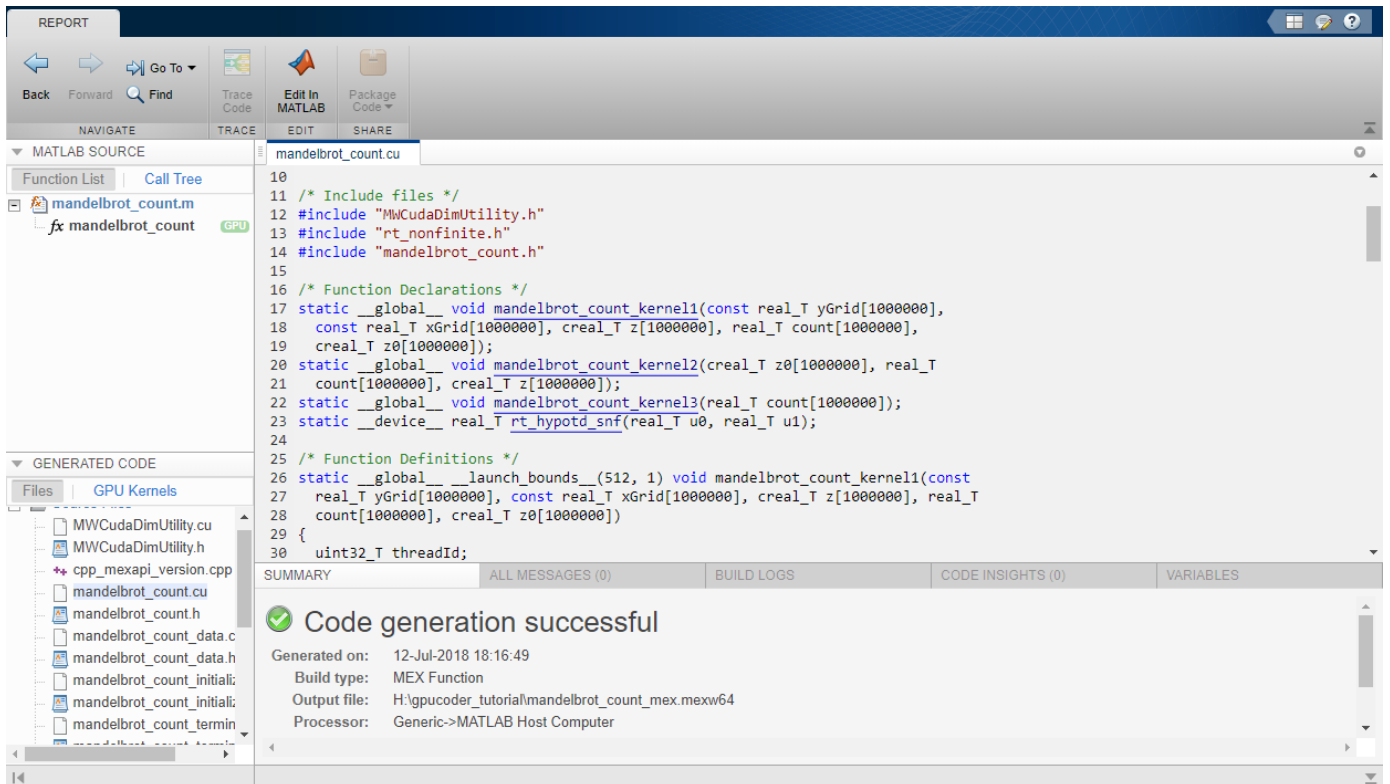
UI Setting	Value Type	Description
Compiler Flags	String	Pass additional flags to the GPU compiler. For example, <code>--fmad=false</code> instructs the <code>nvcc</code> compiler to disable contraction of floating-point multiply and add to a single Floating-Point Multiply-Add (FMAD) instruction. For similar NVIDIA compiler options, see the topic on <i>NVCC Command Options</i> in the CUDA toolkit documentation.
SelectCudaDevice	Integer	In a multi GPU environment such as NVIDIA Drive platforms, specify the CUDA device to target.

3 Click **Generate**.

GPU Coder generates the MEX executable `mandelbrot_count_mex` in your working folder. The `<pwd>\codegen\mex\mandelbrot_count` folder contains all other the generated files including the CUDA source (`*.cu`) and header files. The GPU Coder app indicates that the code generation succeeded. It displays the source MATLAB files and generated output files on the left side of the page. On the **Variables** tab, it displays information about the MATLAB source variables. On the **Target Build Log** tab, it displays the build log, including compiler warnings and errors. By default, in the code window, the app displays the CUDA source file `mandelbrot_count.cu`. To view a different file, in the **Source Code** or **Output Files** pane, click the file name.



4 To view the code generation report, click **View Report**. The report provides links to your MATLAB code and the generated CUDA (`*.cu`) files. It also provides compile-time information for the variables and expressions in your MATLAB code. This information helps you to find sources of error and warnings. It also helps you to debug code generation issues in your code. For more information, see “Code Generation Reports”.



The **GPU Kernels** section on the **Generated Code** tab provides a list of kernels created during GPU code generation. The items in this list link to the relevant source code. For example, when you click **mandelbrot_count_kernel1**, the code section for this kernel is shown in the code browser window.

After you review the report, you can close the **Code Generation Report** window. To view the report later, open `report.mldatx` in `<pwd>\codegen\mex\mandelbrot_cout\html` folder.

- 5 The `<pwd>\codegen\mex\mandelbrot_count` contains the `gpu_codegen_info.mat` MAT-file that contains the statistics for the generated GPU code. This MAT-file contains the `cuda_Kernel` variable that has information about the thread and block sizes, shared and constant memory usage, and input and output arguments of each kernel. The `cudaMalloc` and `cudaMemcpy` variables contain information about the size of all the GPU variables and the number of `memcpy` calls between the host and the device.

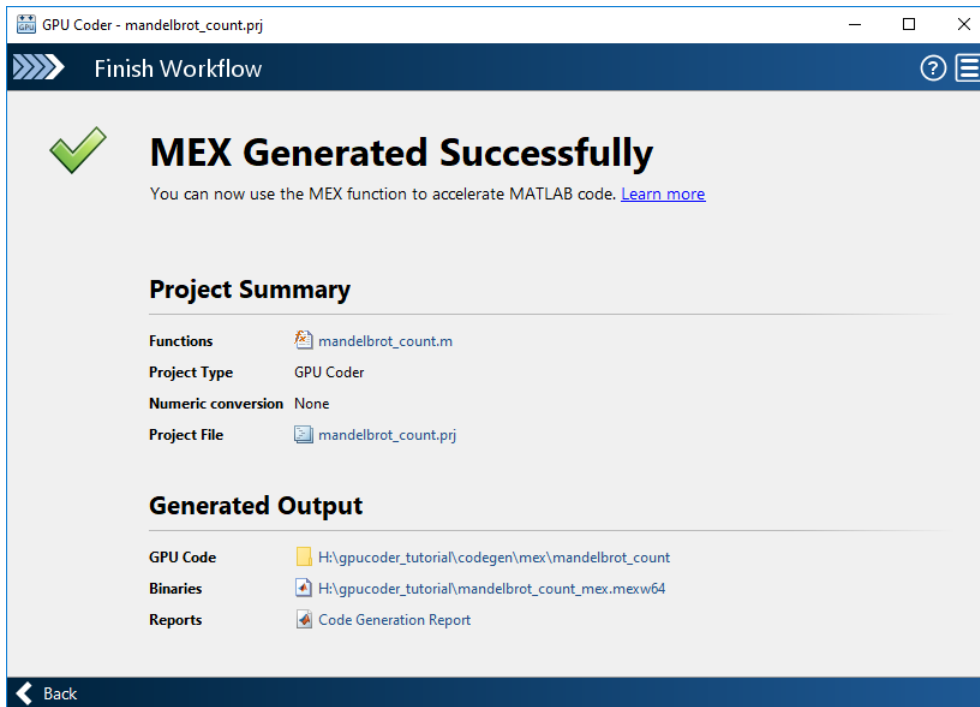
The screenshot shows the 'Variables - cudaMemcpy' window in the GPU Coder app. It displays three variable configuration panels:

- cudaKernel**: A 3x1 struct with 11 fields. The first three rows show configurations for 'mandelbrot_cou...' with threads [512;1;1], blocks [1954;1;1], and various input/output arguments like 'gpu_yGrid,g...' and 'gpu_z,gpu_co...'. The stream is 0, sharedMemory is [], minBlocksPerSM is 1, and the filename is 'H:\gpucod...'. The constantMemory is an empty array, and sourceFileLocationData is a 9x1 struct.
- cudaMalloc**: A 5x1 struct with 2 fields. The first three rows show configurations for 'gpu_yGrid', 'gpu_xGrid', and 'gpu_z' with sizes 8000000, 8000000, and 16000000 respectively. The fourth row shows 'gpu_count' with size 8000000, and the fifth row shows 'gpu_z0' with size 16000000.
- cudaMemcpy**: A 3x1 struct with 6 fields. The first three rows show configurations for 'gpu_yGrid', 'gpu_xGrid', and 'count'. The destinationvariable, sourcevariable, size, direction, conditionalvariable, and stream fields are populated. For 'gpu_yGrid' and 'gpu_xGrid', the direction is 'host->device'. For 'count', the direction is 'device->host'. The conditionalvariable is 'NO_ENCLOSING_CONDITION' and the stream is 0.

6 In the GPU Coder app, click **Next** to open the **Finish Workflow** page.

Review the Finish Workflow Page

The **Finish Workflow** page indicates that the code generation succeeded. It provides a project summary and links to the MATLAB source files, the code generation report, and the generated output binaries. You can save the configuration parameters of the current GPU Coder project as a MATLAB script. See “Convert MATLAB Coder Project to MATLAB Script”.



Verify Correctness of the Generated Code

To verify the correctness of the generated MEX file, see “Verify Correctness of the Generated Code” on page 3-2.

See Also

Apps GPU Coder

Functions

`codegen` | `coder.gpu.kernelfun` | `coder.gpuConfig` | `gpcoder`

Objects

`coder.CodeConfig` | `coder.EmbeddedCodeConfig` | `coder.gpuConfig`

More About

- “GPU Programming Paradigm” on page 4-2
- “Installing Prerequisite Products” on page 1-3
- “Setting Up the Prerequisite Products” on page 1-7
- “GPU Code Generation Workflow” on page 4-3
- “The GPU Environment Check and Setup App” on page 1-12

Code Generation Using the Command Line Interface

The easiest way to create CUDA kernels is to place the `coder.gpu.kernel fun` pragma into your primary MATLAB function. The primary function is also known as the top-level or entry-point function. When the GPU Coder encounters `kernel fun` pragma, it attempts to parallelize all the computation within this function and then maps it to the GPU.

Learning Objectives

In this tutorial, you learn how to:

- Prepare your MATLAB code for CUDA code generation by using the `kernel fun` pragma.
- Create and set up a GPU Coder project.
- Define function input properties.
- Check for code generation readiness and run-time issues.
- Specify code generation properties.
- Generate CUDA code by using the `codegen` command.

Tutorial Prerequisites

This tutorial requires the following products:

- MATLAB
- MATLAB Coder
- GPU Coder
- C compiler
- NVIDIA GPU enabled for CUDA
- CUDA toolkit and driver
- Environment variables for the compilers and libraries. For more information, see “Environment Variables” on page 1-7

Example: The Mandelbrot Set

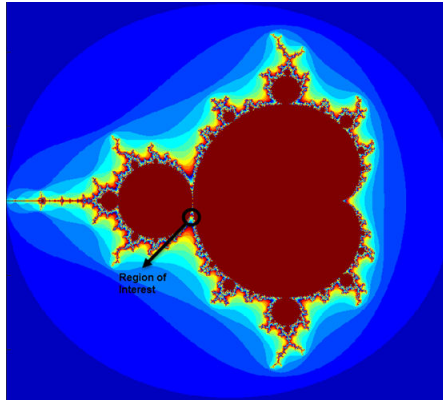
Description

You do not have to be familiar with the algorithm in the example to complete the tutorial.

The Mandelbrot set is the region in the complex plane consisting of the values z_0 for which the trajectories defined by

$$z_{k+1} = z_k^2 + z_0, \quad k = 0, 1, \dots$$

remain bounded at $k \rightarrow \infty$. The overall geometry of the Mandelbrot set is shown in the figure. This view does not have the resolution to show the richly detailed structure of the fringe just outside the boundary of the set. At increasing magnifications, the Mandelbrot set exhibits an elaborate boundary that reveals progressively finer recursive detail.



Algorithm

For this tutorial, pick a set of limits that specify a highly zoomed part of the Mandelbrot set in the valley between the main cardioid and the p/q bulb to its left. A 1000x1000 grid of real parts (x) and imaginary parts (y) is created between these two limits. The Mandelbrot algorithm is then iterated at each grid location. An iteration number of 500 is enough to render the image in full resolution.

```
maxIterations = 500;
gridSize = 1000;
xlim = [-0.748766713922161, -0.748766707771757];
ylim = [ 0.123640844894862,  0.123640851045266];
```

An implementation of the Mandelbrot set by using standard MATLAB commands running on the CPU is shown. This implementation is based on the code provided in the “*Experiments with MATLAB*” e-book by Cleve Moler. This calculation is vectorized such that every location is updated simultaneously.

Tutorial Files

Create a MATLAB script called `mandelbrot_count.m` with the following lines of code. This code is a baseline vectorized MATLAB implementation of the Mandelbrot set. Later in this tutorial, you modify this file to make it suitable for code generation.

```
function count = mandelbrot_count(maxIterations, xGrid, yGrid)
% mandelbrot computation

z0 = xGrid + 1i*yGrid;
count = ones(size(z0));

z = z0;
for n = 0:maxIterations
    z = z.*z + z0;
    inside = abs(z)<=2;
    count = count + inside;
end
count = log(count);
```

Create a MATLAB script called `mandelbrot_test.m` with the following lines of code. The script generates 1000 x 1000 grid of real parts (x) and imaginary parts (y) between the limits specified by `xlim` and `ylim`. It also calls the `mandelbrot_count` function and plots the resulting Mandelbrot set.

```
maxIterations = 500;
gridSize = 1000;
```

```
xlim = [-0.748766713922161, -0.748766707771757];
ylim = [ 0.123640844894862,  0.123640851045266];

x = linspace( xlim(1), xlim(2), gridSize );
y = linspace( ylim(1), ylim(2), gridSize );
[xGrid,yGrid] = meshgrid( x, y );


%% Mandelbrot computation in MATLAB
count = mandelbrot_count(maxIterations, xGrid, yGrid);

% Show
figure(1)
imagesc( x, y, count );
colormap([jet();flipud( jet() );0 0 0]);
axis off
title('Mandelbrot set with MATLAB');
```

Run the Original MATLAB Code

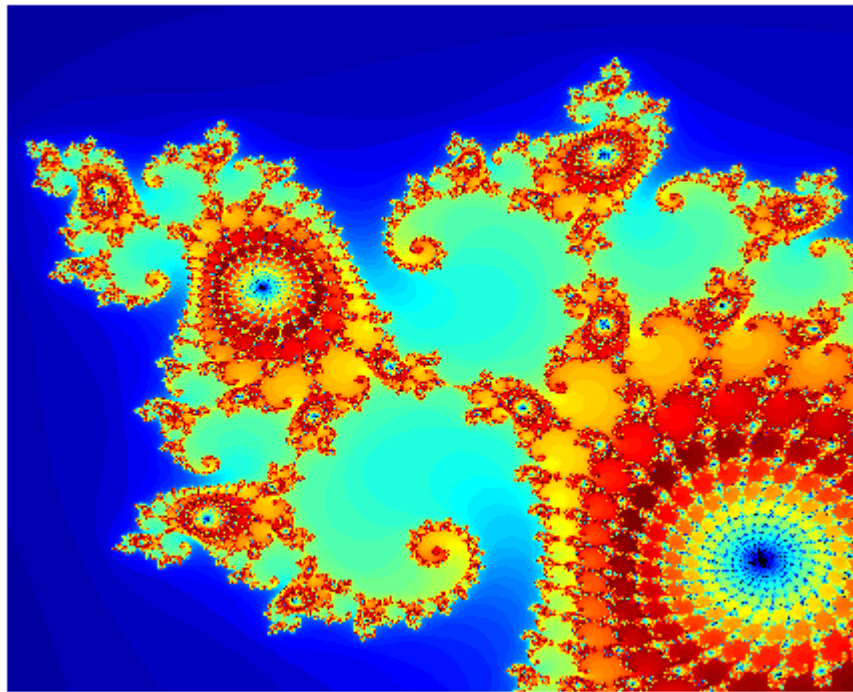
Run the Mandelbrot Example

Before making the MATLAB version of the Mandelbrot set algorithm suitable for code generation, you can test the functionality of the original code.

- 1 Change the current working folder of MATLAB to the location that contains the two files you created in the previous step. GPU Coder places generated code in this folder, change your current working folder if you do not have full access to this folder.
- 2 Open the `mandelbrot_test` script in the MATLAB Editor.
- 3 Run the test script by clicking the run button  or by entering `mandelbrot_test` in the MATLAB Command Window.

The test script runs and shows the geometry of the Mandelbrot within the boundary set by the variables `xlim` and `ylim`.

Mandelbrot Set with MATLAB



Make the MATLAB Code Suitable for Code Generation

To begin the process of making your MATLAB code suitable for code generation, use the file `mandelbrot_count.m`.

- 1 Set your MATLAB current folder to the work folder that contains your files for this tutorial.
- 2 In the MATLAB Editor, open `mandelbrot_count.m`. The file opens in the MATLAB Editor. The Code Analyzer message indicator in the top right corner of the MATLAB Editor is green. The analyzer did not detect errors, warnings, or opportunities for improvement in the code.
- 3 Turn on MATLAB for code generation error checking. After the function declaration, add the `%#codegen` directive.

```
function count = mandelbrot_count(maxIterations, xGrid, yGrid) %#codegen
```

The Code Analyzer message indicator remains green, indicating that it has not detected code generation issues.

- 4 To map the `mandelbrot_count` function to a CUDA kernel, modify the original MATLAB code by placing the `coder.gpu.kernelfun` pragma outside the `for`-loop body.

```
function count = mandelbrot_count(maxIterations, xGrid, yGrid) %#codegen
% mandelbrot computation

z0 = xGrid + 1i*yGrid;
count = ones(size(z0));

% Add Kernelfun pragma to trigger kernel creation
coder.gpu.kernelfun;
```

```

z = z0;
for n = 0:maxIterations
    z = z.*z + z0;
    inside = abs(z)<=2;
    count = count + inside;
end
count = log(count);

```

When using the `coder.gpu.kernelfun` pragma, GPU Coder attempts to map the computations in the function `mandelbrot_count` to the GPU.

- 5 Save the file. You are now ready to compile your code by using the command-line interface.

Code Generation from the Command Line

You can use the `codegen` command to translate MATLAB functions to a CUDA compatible C/C++ static or dynamic library, executable, or MEX function, instead of using the GPU Coder app.

Define Input Types

At compile time, GPU Coder must know the data types of all the inputs to the entry-point function. Therefore, if your entry-point function has inputs, you must specify its data type at the time that you compile the file with the `codegen` function.

You can generate inputs and then use the `-args` option in the `codegen` function to let GPU Coder determine the class, size, and complexity of the input parameters. To generate inputs for `mandelbrot_count` function, use these commands:

```

maxIterations = 500;
gridSize = 1000;
xlim = [-0.748766713922161, -0.748766707771757];
ylim = [ 0.123640844894862,  0.123640851045266];

x = linspace( xlim(1), xlim(2), gridSize );
y = linspace( ylim(1), ylim(2), gridSize );
[xGrid,yGrid] = meshgrid( x, y );

```

Alternatively, you can specify the size, type and complexity of the inputs to the entry-point functions without generating input data by using the `coder.typeof` function.

```

ARGS = cell(1,1);
ARGS{1} = cell(3,1);
ARGS{1}{1} = coder.typeof(0);
ARGS{1}{2} = coder.typeof(0,[1000 1000]);
ARGS{1}{3} = coder.typeof(0,[1000 1000]);

```

Build Configuration

To configure build settings such as output file name, location, type, you have to create `coder` configuration objects. To create the objects, use the `coder.gpuConfig` function. For example, to create a `coder.MexCodeConfig` code generation object for use with `codegen` when generating a MEX function, use:

```
cfg = coder.gpuConfig('mex');
```

Other available options are:

- `cfg = coder.gpuConfig('lib');`, to create a code generation configuration object for use with `codegen` when generating a CUDA C/C++ static library.

- `cfg = coder.gpuConfig('dll');` to create a code generation configuration object for use with codegen when generating a CUDA C/C++ dynamic library.
- `cfg = coder.gpuConfig('exe');` to create a code generation configuration object for use with codegen when generating a CUDA C/C++ executable.

For more information, see `coder.gpuConfig`.

Each configuration object comes with a set of parameters, initialized to default values. You can use dot notation to modify the value of one configuration object parameter at a time. Use this syntax:

```
configuration_object.property = value
```

You can enable the same settings as in the “Code Generation by Using the GPU Coder App” on page 2-2 by using the following command-line equivalents:

```
cfg = coder.gpuConfig('mex');
cfg.GpuConfig.CompilerFlags = '--fmad=false';
cfg.GenerateReport = true;
```

The `cfg` configuration object has configuration parameters that are common to MATLAB Coder and GPU Coder and parameters that are GPU Coder-specific. You can see all the GPU-specific properties available in the `cfg` configuration object by typing `cfg.GpuConfig` in the MATLAB Command Window.

```
>> cfg.GpuConfig

ans =

    config with properties:

        Enabled: 1
        MallocMode: 'discrete'
        KernelNamePrefix: ''
        EnableCUBLAS: 1
        EnableCUSOLVER: 1
        EnableCUFFT: 1
        Benchmarking: 0
        SafeBuild: 0
        ComputeCapability: '3.5'
        CustomComputeCapability: ''
        CompilerFlags: ''
        StackLimitPerThread: 1024
        MallocThreshold: 200
        SelectCudaDevice: -1
```

The `--fmad=false` flag when passed to the `nvcc`, instructs the compiler to disable Floating-Point Multiply-Add (FMAD) optimization. This option is set to prevent numerical mismatch in the generated code because of architectural differences in the CPU and the GPU. For more information, see “Numerical Differences Between CPU and GPU” on page 3-4.

For more information on configuration parameters that are common to MATLAB Coder and GPU Coder, see `coder.CodeConfig` class.

Build Script

You can create a build script `mandelbrot_codegen.m` that automates the series of commands mentioned previously.

```
% GPU code generation for getting started example (mandelbrot_count.m)
%% Create configuration object of class 'coder.MexCodeConfig'.
cfg = coder.gpuConfig('mex');
cfg.GenerateReport = true;
cfg.GpuConfig.CompilerFlags = '--fmad=false';

%% Define argument types for entry-point 'mandelbrot_count'.
ARGS = cell(1,1);
ARGS{1} = cell(3,1);
ARGS{1}{1} = coder.typeof(0);
ARGS{1}{2} = coder.typeof(0,[1000 1000]);
ARGS{1}{3} = coder.typeof(0,[1000 1000]);

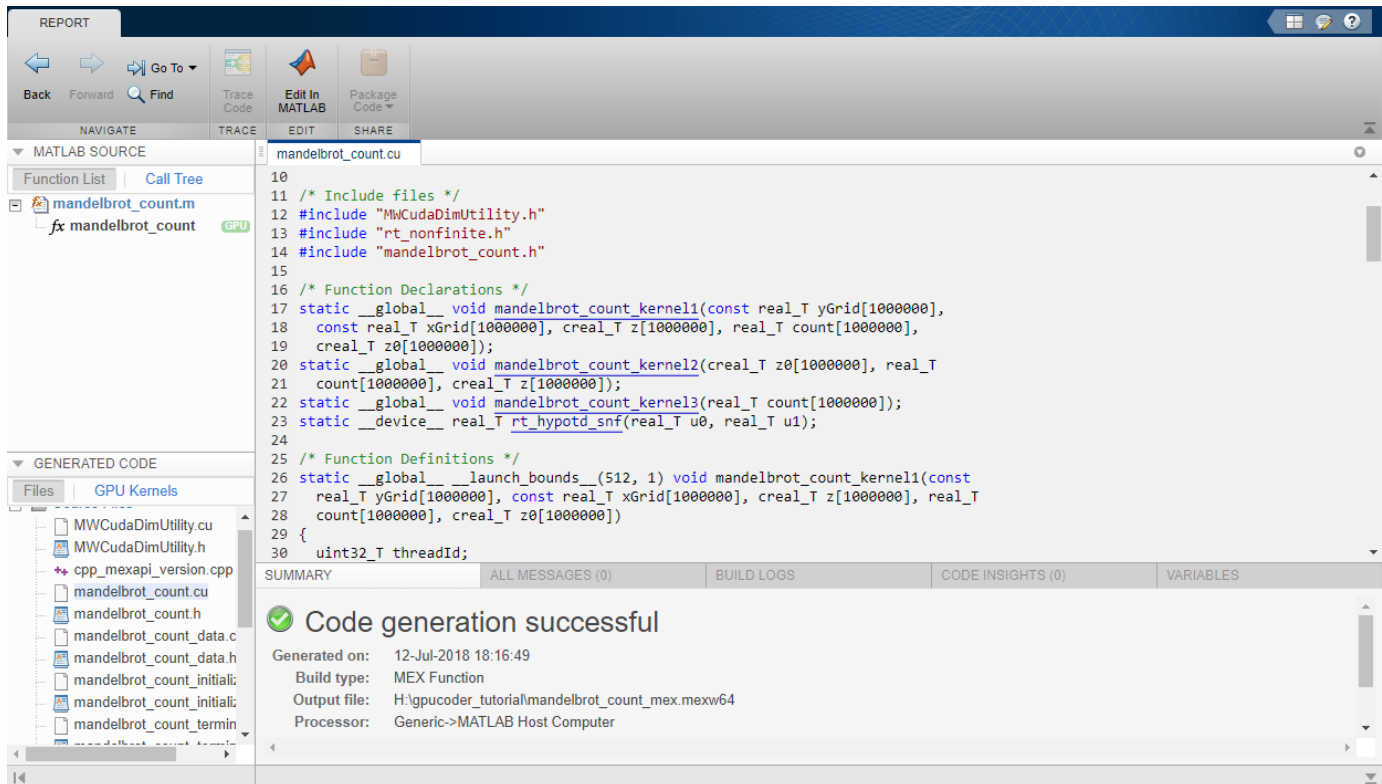
%% Invoke GPU Coder.
codegen -config cfg mandelbrot_count -args ARGS{1}
```

The `codegen` command opens the file `mandelbrot_count.m` and translates the MATLAB code into CUDA code.

- The `-report` option instructs `codegen` to generate a code generation report that you can use to debug your MATLAB code.
- The `-args` option instructs `codegen` to compile the file `mandelbrot_count.m` by using the class, size, and complexity of the input parameters *maxIterations*, *xGrid*, and *yGrid*.
- The `-config` option instructs `codegen` to use the specified configuration object for code generation.

When code generation is successful, you can view the resulting code generation report by clicking **View Report** in the MATLAB Command Window.

```
>> mandelbrot_codegen
Code generation successful: View report
```



Verify Correctness of the Generated Code

To verify correctness of the generated MEX file, see “Verify Correctness of the Generated Code” on page 3-2.

See Also

Apps

GPU Coder

Functions

codegen | coder.gpu.kernelfun | coder.gpuConfig | gpucoder

Objects

coder.CodeConfig | coder.EmbeddedCodeConfig | coder.gpuConfig

More About

- “GPU Programming Paradigm” on page 4-2
- “Installing Prerequisite Products” on page 1-3
- “Setting Up the Prerequisite Products” on page 1-7
- “GPU Code Generation Workflow” on page 4-3
- “The GPU Environment Check and Setup App” on page 1-12

GPU Code Generation: The Mandelbrot Set

This example shows how to generate CUDA® code from a simple MATLAB® function by using GPU Coder™. A Mandelbrot set implementation by using standard MATLAB commands acts as the entry-point function. This example uses the `codegen` command to generate a MEX function that runs on the GPU. You can run the MEX function to check for run-time errors.

Third-Party Prerequisites

Required

This example generates CUDA MEX and has the following third-party requirements.

- CUDA enabled NVIDIA® GPU and compatible driver.

Optional

For non-MEX builds such as static, dynamic libraries or executables, this example has the following additional requirements.

- NVIDIA toolkit.
- Environment variables for the compilers and libraries. For more information, see “Third-Party Hardware” on page 1-3 and “Setting Up the Prerequisite Products” on page 1-7.

Verify GPU Environment

To verify that the compilers and libraries necessary for running this example are set up correctly, use the `coder.checkGpuInstall` function.

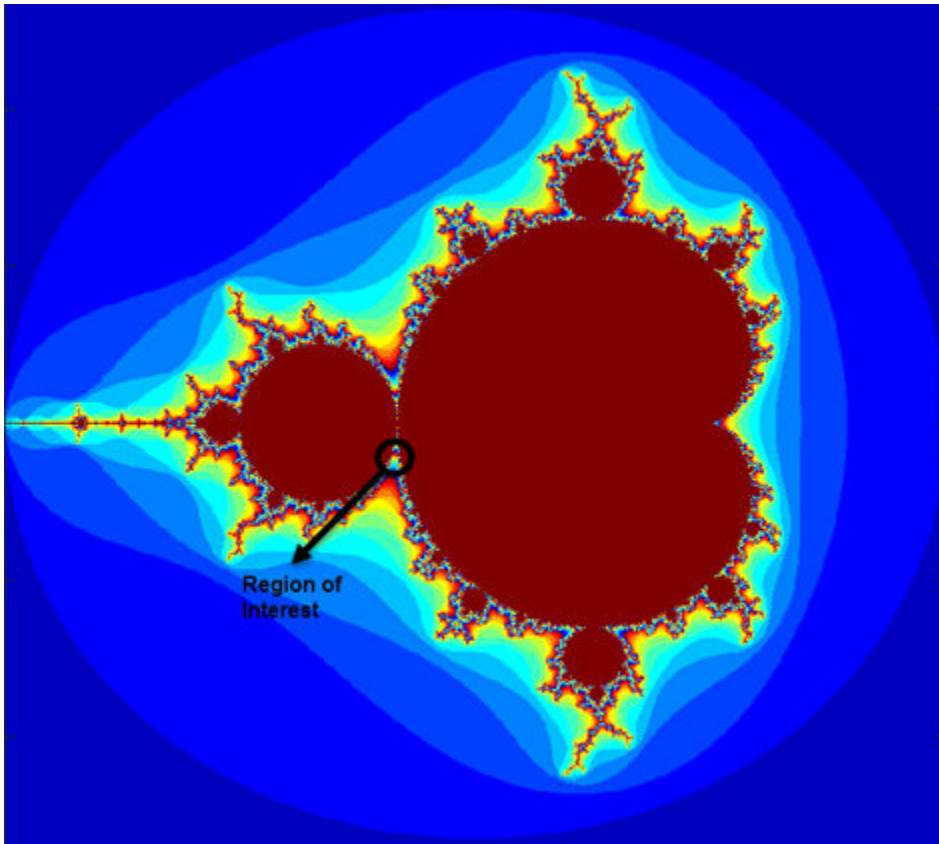
```
envCfg = coder.gpuEnvConfig('host');  
envCfg.BasicCodegen = 1;  
envCfg.Quiet = 1;  
coder.checkGpuInstall(envCfg);
```

Mandelbrot Set

The Mandelbrot set is the region in the complex plane consisting of the values z_0 for which the trajectories defined by

$$z_{k+1} = z_k^2 + z_0, k = 0, 1, \dots$$

remain bounded as $k \rightarrow \infty$. The overall geometry of the Mandelbrot set is shown in the figure. This view does not have the resolution to show the richly detailed structure of the fringe just outside the boundary of the set.



Define Input Regions

Pick a set of limits that specify a highly zoomed part of the Mandelbrot set in the valley between the main cardioid and the p/q bulb to its left. A 1000×1000 grid of $Re\{x\}$ and $Im\{y\}$ is created between these two limits. The Mandelbrot algorithm is then iterated at each grid location. An iteration number of 500 is enough to render the image in full resolution.

```
maxIterations = 500;
gridSize = 1000;
xlim = [-0.748766713922161, -0.748766707771757];
ylim = [ 0.123640844894862,  0.123640851045266];

x = linspace( xlim(1), xlim(2), gridSize );
y = linspace( ylim(1), ylim(2), gridSize );
[xGrid,yGrid] = meshgrid( x, y );
```

The Mandelbrot Entry-Point Function

The `mandelbrot_count.m` entry-point function contains a vectorized implementation of the Mandelbrot set based on the code provided in the e-book *Experiments with MATLAB* by Cleve Moler. The `%#codegen` directive turns on MATLAB for code generation error checking. When GPU Coder encounters the `coder.gpu.kernelFun` pragma, it attempts to parallelize all the computation within this function, and then maps it to the GPU.

```
type mandelbrot_count
```

```
function count = mandelbrot_count(maxIterations, xGrid, yGrid) %#codegen
% Copyright 2016-2019 The MathWorks, Inc.

z0 = xGrid + 1i*yGrid;
count = ones(size(z0));

% Map computation to GPU.
coder.gpu.kernelfun;

z = z0;
for n = 0:maxIterations
    z = z.*z + z0;
    inside = abs(z)<=2;
    count = count + inside;
end
count = log(count);
```

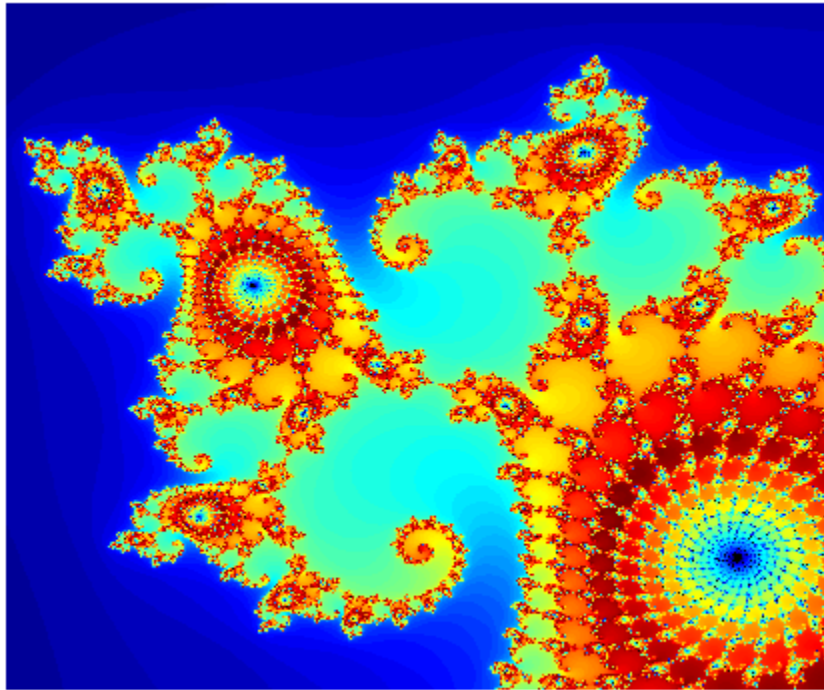
Test the Functionality of mandelbrot_count

Run the `mandelbrot_count` function with the `xGrid`, `yGrid` values that were previously generated, and then plot the results.

```
count = mandelbrot_count(maxIterations, xGrid, yGrid);

figure(2), imagesc( x, y, count );
colormap( [jet();flipud( jet() );0 0 0] );
title('Mandelbrot Set on MATLAB');
axis off
```

Mandelbrot Set on MATLAB



Generate CUDA MEX for the Function

To generate CUDA MEX for the `mandelbrot_count` function, create a GPU code configuration object and run the `codegen` command. Because of architectural differences between the CPU and GPU, numeric verification does not always match. This scenario is true when using the single data type in your MATLAB code and performing accumulation operations on these single data type values. Like this Mandelbrot example even the double data types cause numeric errors. One reason for this mismatch is that the GPU floating-point units use fused Floating-point Multiply-Add (FMAD) instructions and the CPU does not use these instructions. The `fmad=false` option that is passed to the `nvcc` compiler turns off this FMAD optimization.

```
cfg = coder.gpuConfig('mex');
cfg.GpuConfig.CompilerFlags = '--fmad=false';
codegen -config cfg -args {maxIterations,xGrid,yGrid} mandelbrot_count
```

Code generation successful: To view the report, open('codegen/mex/mandelbrot_count/html/report.m')

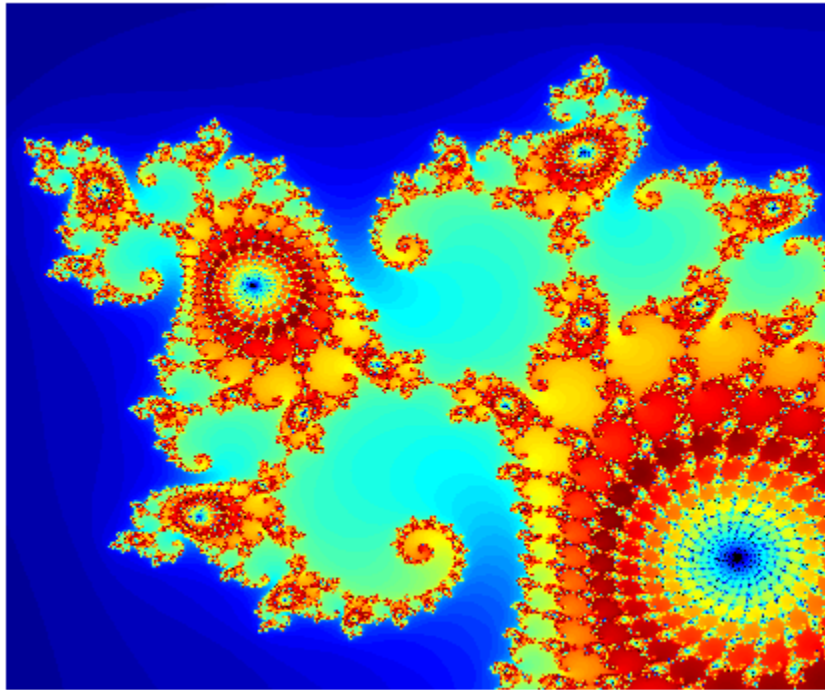
Run the MEX Function

After generating a MEX function, verify that it has the same functionality as the original MATLAB entry-point function. Run the generated `mandelbrot_count_mex` and plot the results.

```
countGPU = mandelbrot_count_mex(maxIterations, xGrid, yGrid);

figure(2), imagesc( x, y, countGPU );
colormap( [jet();flipud( jet() );0 0 0] );
title('Mandelbrot Set on GPU');
axis off
```

Mandelbrot Set on GPU



Conclusion

CUDA code was generated for a simple MATLAB function implementing the Mandelbrot set. Implementation was accomplished by using the `coder.gpu.kernel fun` pragma and invoking the `codegen` command to generate MEX function. Additional compiler flags, `FMAD=false`, was passed to the `nvcc` compiler to disable the FMAD optimization that the NVIDIA compilers performed.

Verification

Verify Correctness of the Generated Code

After you generate code, GPU Coder provides you multiple options to inspect the source code and test the correctness of the generated code.

- The code generation report provides an interactive interface for inspecting the generated CUDA C++ source files, generated data types, and other code insights. For more information, see “Code Generation Reports”.
- Verify generated MEX functions in the GPU Coder app.
- Verify generated MEX functions at the command line.
- With the GPU Coder Support Package for NVIDIA GPUs, you can use the processor-in-the-loop (PIL) execution to check the numerical behavior of the CUDA code that you generate from MATLAB functions. For more information, see “Processor-In-The-Loop Execution from Command Line” (GPU Coder Support Package for NVIDIA GPUs) and “Processor-In-The-Loop Execution with the GPU Coder App” (GPU Coder Support Package for NVIDIA GPUs).
- If you have Embedded Coder, you can verify the numerical behavior of the generated CUDA C++ code by using software-in-the-loop (SIL) execution.

Verify MEX Functions in the GPU Coder App

In the GPU Coder app, after you generate a MEX function, you can verify that it has the same functionality as the original MATLAB entry-point function. Provide a test file that generates test vectors for the function under test and then calls the original MATLAB entry-point function. The test file can be a MATLAB function or script but must be in the same folder as the original entry-point function.

- On the **Generate Code** page of the GPU Coder app, click **Verify Code**.
- Type or select the test file. For example, `myfunction_test`.
- To run the test file without replacing calls to the original MATLAB function with calls to the MEX function, select **MATLAB code** for the **Run using** option. Click **Run Generated Code**.
- To run the test file, replacing calls to the original MATLAB function with calls to the MEX function, select **Generated code** for the **Run using** option. Click **Run Generated Code**.
- Compare the results of running the original MATLAB function with the results of running the generated CUDA MEX function.

Verify MEX Functions at the Command Line

You can verify the generated CUDA MEX file at the command line by using the `coder.runTest` function. The `coder.runTest` function runs the test file by replacing calls to the original MATLAB function with calls to the MEX function. For example, to test `myfunction` function with `myfunction_test` test file, enter the following code in the MATLAB Command Window.

```
coder.runTest('myfunction_test','myfunction')
```

Compare the results with the results of running the original MATLAB function. If errors occur during the run, call stack information is available for debugging. Alternatively, you can use `codegen` with the `-test` option.

```
codegen myfunction -test'myfunction_test'
```


The test file can be a MATLAB function, script, or class-based unit test.

Code Verification Through Software-In-The-Loop



GPU Coder supports software-in-the-loop (SIL) execution, which enables you to verify source code and compiled object code. During SIL execution through a MATLAB SIL interface, the software compiles and using the test vectors that you provide, runs library code on your development computer. You can reuse test vectors developed for your MATLAB functions to verify the numerical behavior of library code.

Note

- On a Microsoft Windows system, the Windows Firewall can potentially block a SIL execution. Change the Windows Firewall settings to allow access.
 - For SIL on a Windows platform, it is recommended to set the code generation language to C++ instead of C.
 - When using SIL execution, make sure that the **Benchmarking** option in GPU Coder settings is `false`. Executing SIL with benchmarking results in compilation errors.
-

SIL Execution with the GPU Coder App

The software-in-the-loop (SIL) execution is supported only for static and dynamic library output types. If you generated a MEX function, you must change the project settings to use `Static Library` or `Dynamic Library` for **Build type** and run **Generate code** again.

- 1 To open the GPU Coder app, on the MATLAB toolstrip **Apps** tab, under **Code Generation**, click the app icon.
- 2 To open your project, click , and then click `Open existing project`. Select the project, `myproject.prj`. On the **Generate Code** page, click the **Generate** arrow .
- 3 In the **Generate** dialog box:
 - Set **Build type** to `Static Library` or `Dynamic Library`.
 - Clear the **Generate code only** check box.
 - You can leave the other settings to their default values.
- 4 To generate the CUDA C++ code, click **Generate**. Click **Verify Code**.
- 5 In the command field, specify the test file (for example, `myfunction_test.m`) that calls the original MATLAB functions (for example, `myfunction`).
- 6 To start the SIL execution, click **Run Generated Code**. The GPU Coder app:
 - Generates a standalone library in `codegen\lib\myfunction`.
 - Generates SIL interface code in `codegen\lib\myfunction\sil`.
 - Runs the test file, replacing calls to the MATLAB function with calls to the generated code in the library.
 - Displays messages from the SIL execution in the **Test Output** tab.

- 7 Verify that the results from the SIL execution match the results from the original MATLAB functions.
- 8 To terminate the SIL execution process, click **Stop SIL Verification**. Alternatively, on the **Test Output** tab, click the link that follows To terminate execution.

SIL Execution from Command Line

To set up and start a SIL execution from the command line, you create a GPU Coder configuration object for library code generation, enable `config_obj.VerificationMode = 'SIL'`, use `codegen` function to generate the library code and the SIL interface, and use `coder.runTest` function to run the test file for your original MATLAB function. The following is a build script that automates the series of commands to perform SIL execution.

```

%% Create configuration object for static library.
config = coder.gpuConfig('lib');
config.GenerateReport = true;
config.VerificationMode = 'SIL';

%% Define argument types for entry-point 'mandelbrot_count'.
ARGS = cell(1,1);
ARGS{1} = cell(3,1);
ARGS{1}{1} = coder.typeof(0);
ARGS{1}{2} = coder.typeof(0,[1000 1000]);
ARGS{1}{3} = coder.typeof(0,[1000 1000]);

%% Invoke GPU Coder.
codegen -config config myfunction -args ARGS{1}

%% Run the test file with the sil interface
coder.runTest('myfunction_test', ['myfunction_sil.' mexext]);

%% Terminate SIL execution
clear myfunction_sil;

```

Numerical Differences Between CPU and GPU

Because of architectural differences between the CPU and GPU, numerical verification does not always match. This scenario is specially true when using `single` data type in your MATLAB code and performing accumulation operations on these `single` data type values. However, there are cases like the Mandelbrot example where even `double` data types cause numerical errors. One reason for this mismatch is that the GPU floating-point units use fused Floating-Point Multiply-Add (FMAD) instructions while the CPU does not use these instructions. It is also important to note that the CUDA compiler performs these instruction-level optimizations by default impacting the accuracy of the computed results. For example, the CUDA compiler fuses floating-point multiply and add instructions into a single instruction. This Floating-point Multiply-Add (FMAD) operation executes twice as fast compared to two single instructions but results in the loss of numerical accuracy. You can achieve tighter control over these optimizations by using intrinsic functions and compiler flags. See GPU Coder Configuration Properties to use compiler flags and “Legacy Code Integration” for using intrinsics.

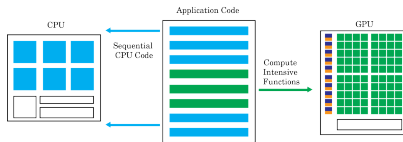
See Also

About GPU Code Generation

- “GPU Programming Paradigm” on page 4-2
- “GPU Code Generation Workflow” on page 4-3

GPU Programming Paradigm

GPU-accelerated computing follows a heterogeneous programming model. Highly parallelizable portions of the software application are mapped into kernels that execute on the physically separate GPU device, while the remainder of the sequential code still runs on the CPU. Each kernel is allocated several workers or threads, which are organized in blocks and grids. Every thread within the kernel executes concurrently with respect to each other.



The objective of GPU Coder is to take a sequential MATLAB program and generate partitioned, optimized CUDA code from it. This process involves:

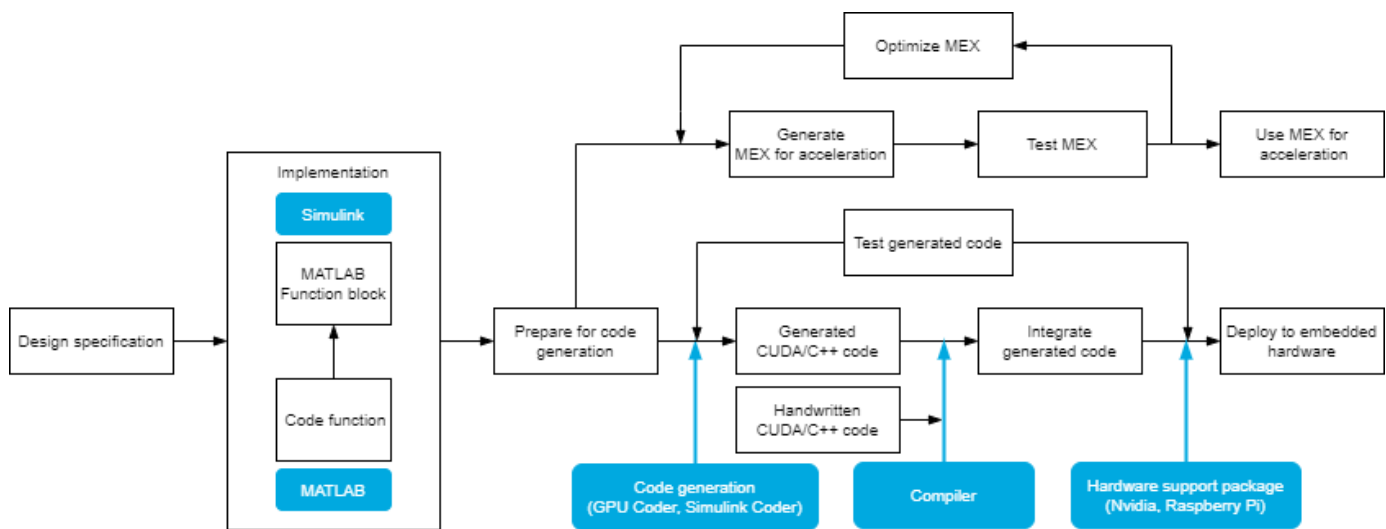
- CPU/GPU partitioning — Identifying segments of code that run on the CPU and segments that run on the GPU. For the different ways GPU Coder identifies CUDA kernels, see “Kernel Creation”. Memory transfer costs between CPU and GPU are a significant consideration in the kernel creation algorithm.
- After kernel partitioning is complete, GPU Coder analyzes the data dependency between the CPU and GPU partitions. Data that is shared between the CPU and GPU are allocated on GPU memory (by using `cudaMalloc` or `cudaMallocManaged` APIs). The analysis also determines the minimum set of locations where data has to be copied between CPU and GPU by using `cudaMemcpy`. If using Unified Memory in CUDA, then the same analysis pass also determines the minimum locations in the code where `cudaDeviceSync` calls must be inserted to get the right functional behavior.
- Next, within each kernel, GPU Coder can choose to map data to shared memory or constant memory. If used wisely, these memories are part of the GPU memory hierarchy structure and can potentially result in greater memory bandwidth. For information on how GPU Coder chooses to map to shared memory, see “Stencil Processing”. For information on how GPU Coder chooses to map to constant memory, see `coder.gpu.constantMemory`.
- Once partitioning and memory allocation and transfer statements are in place, GPU Coder generates CUDA code that follows the partitioning and memory allocation decisions. The generated source code can be compiled into a MEX target to be called from within MATLAB or into a shared library to be integrated with an external project. For GPU Coder options, see GPU Coder Configuration Properties.

See Also

GPU Code Generation Workflow

GPU Coder code generation technology produces CUDA C++ code and executable programs for algorithms. You can write algorithms programmatically by using MATLAB or graphically in the Simulink environment. You can generate code for MATLAB functions and Simulink MATLAB Function blocks that are useful for real-time and embedded applications. Because code generation is tightly integrated with the MATLAB and Simulink execution and simulation engines, the generated source code and executable programs match the functional behavior of MATLAB code executions and Simulink simulations to a high degree of accuracy.

The code generator supports two workflows for designing, implementing, and verifying generated CUDA code. The following figure shows the design and deployment environment options.



See Also

