



# Deploy MATLAB and Simulink to NVIDIA GPUs using GPU Coder

---

Bill Chou

*Code Generation*

*MathWorks*

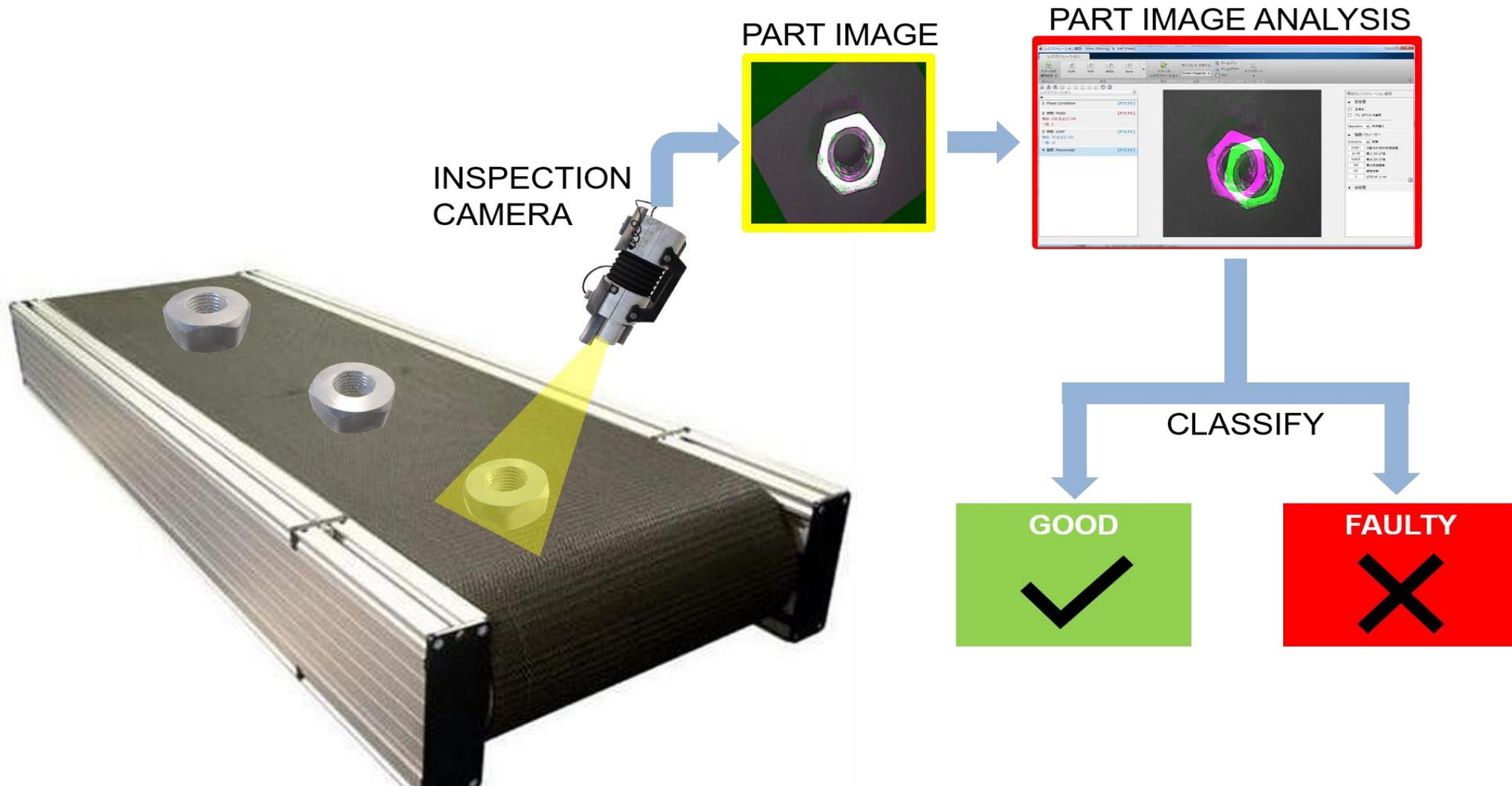
*bchou@mathworks.com*



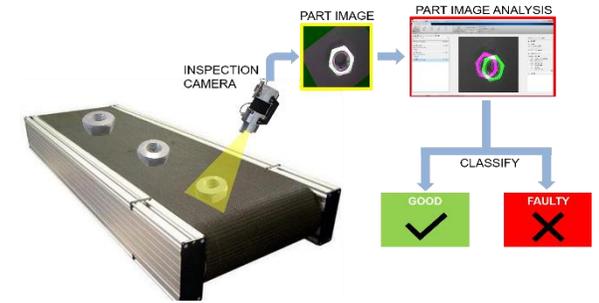
# Agenda

- Why translate MATLAB to CUDA?
- What is GPU Coder?
- Workflow for using GPU Coder
- Deep Learning workflow
- Wrap-up

# Automated Optical Inspection



# Automated Optical Inspection



```
tb_myNDNet_Jetson_Intro.mlx * x tb_myNDNet_Jetson.mlx x myNDNet_Preprocess.m x myNDNet_Postprocess.m x targetFunction.m x main_nutsDet.cpp x
```

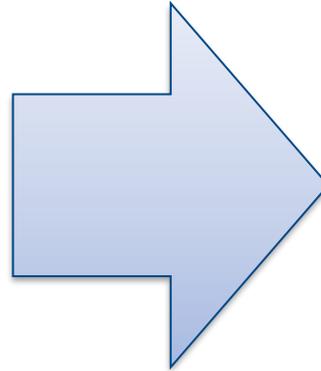
### Grab one frame from input video stream, run inference, identify defect in image

```
2 load('NDNet.mat');
3 img = imread('testImg.png');
4 |
5 wi = uint32(320);
6 he = uint32(240);
7 ch = uint32(3);
8 sz = [he wi ch];
9
10 imgOcv = mat2ocv(img);
11 Weights = convnet.Layers(67).Weights;
12
13 tic;
14 out = targetFunction(imgOcv,Weights,false);
15 toc
16
17 tic;
18 out = targetFunction_mex(imgOcv,Weights,false);
19 toc
20
21 out = ocv2mat(out,sz);
22
23 figure;
24 s1 = subplot(1,2,1);
25 s2 = subplot(1,2,2);
26 imshow(img, 'Parent',s1);
27 title(s1, 'Input');
28
29 imshow(out, 'Parent',s2);
```

# Why Use CUDA code generation ?

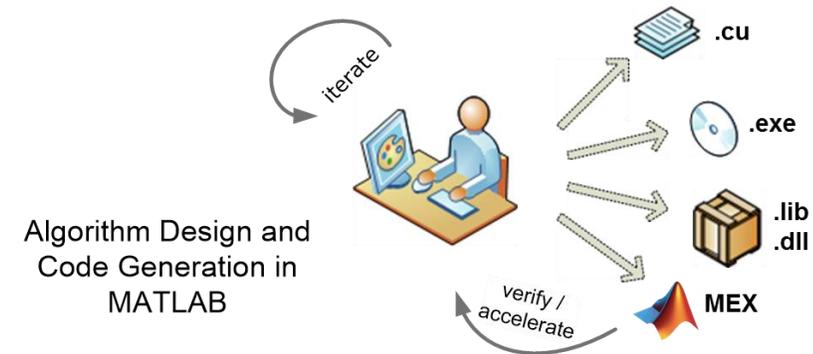
## Pains: Hand code

- **Cannot code in CUDA**
- Time consuming
- Manual Coding Errors
- Multiple implementations
- Expensive



## Solution: GPU Coder

- Automatically convert to CUDA
- Get to optimized CUDA faster
- Eliminate manual coding errors
- Maintain Single “Truth”
- Stay within MATLAB/Simulink at a higher level

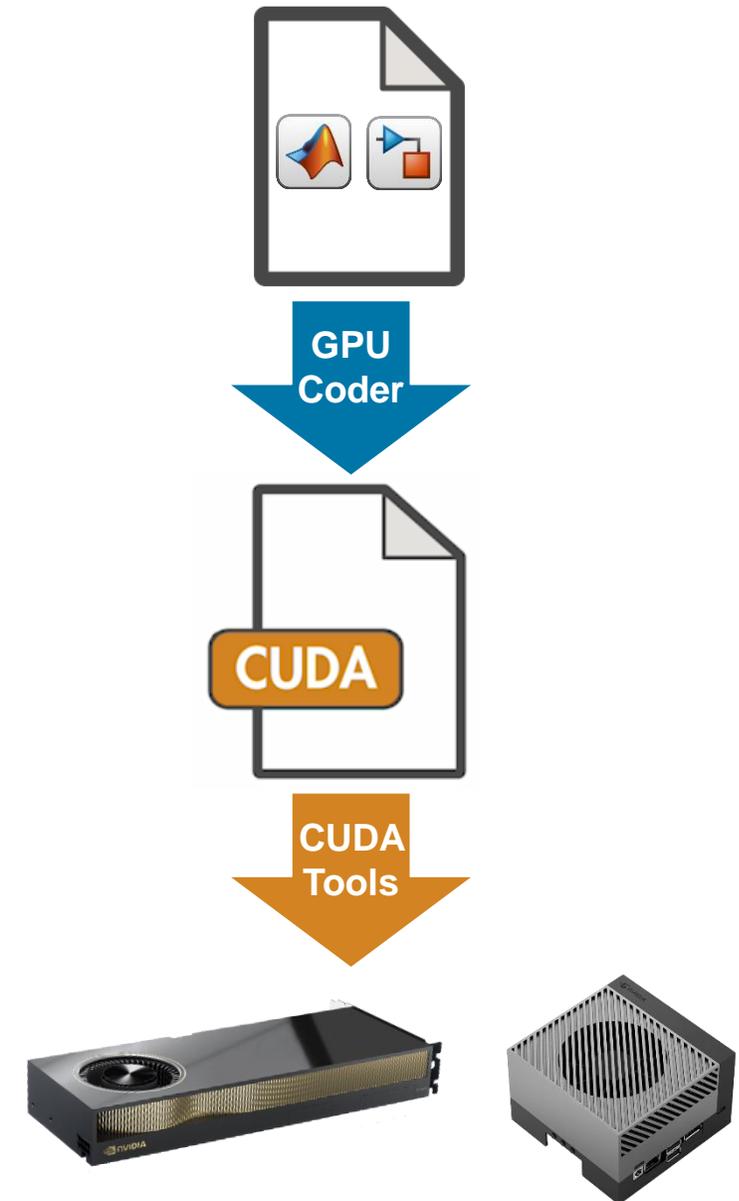


# Agenda

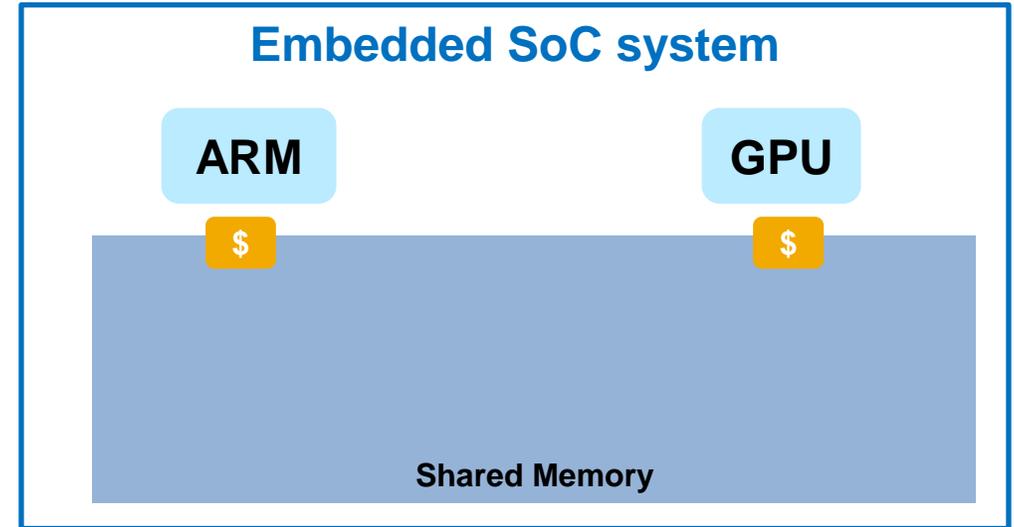
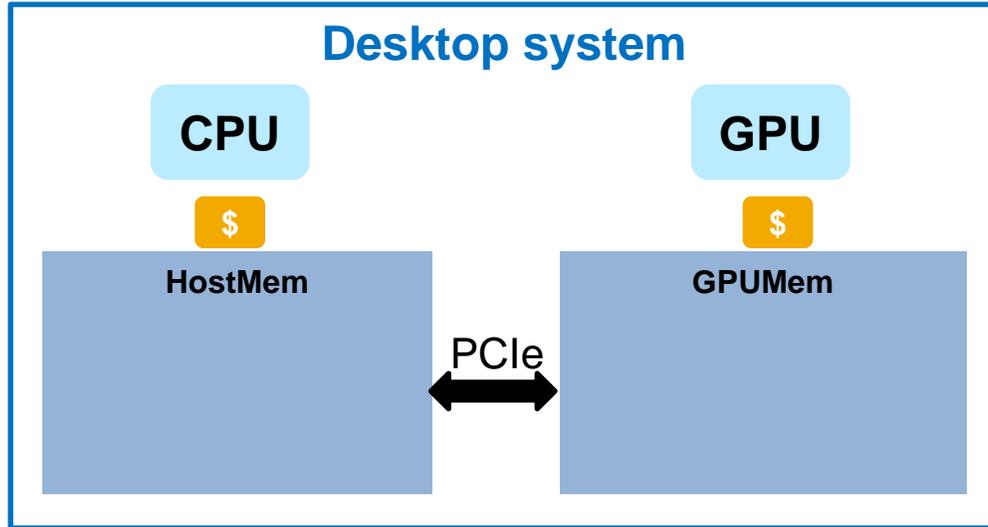
- Why translate MATLAB to CUDA?
- **What is GPU Coder?**
- Workflow for using GPU Coder
- Deep Learning workflow
- Wrap-up

# What is GPU Coder?

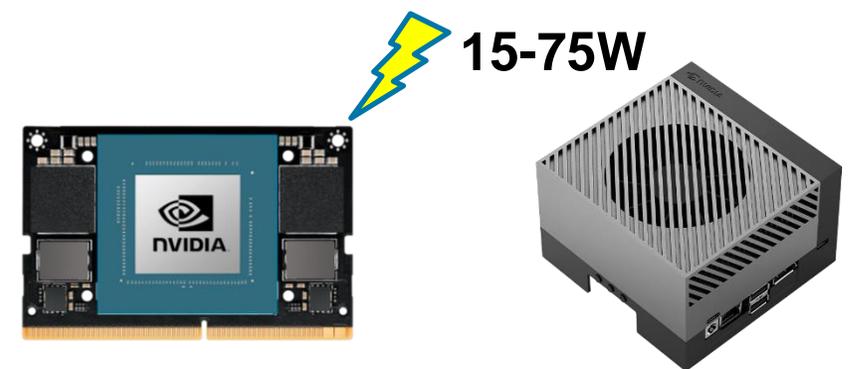
- Tool that automatically converts MATLAB to CUDA
- Generated CUDA is portable across all modern NVIDIA GPUs
- Automates deployment to cloud, desktop, and embedded GPUs
  - Including desktop RTX and embedded Jetson AGX Orin
- Can call optimization libraries like TensorRT
- Profile the generated CUDA to find opportunities to optimize performance



# Types of GPUs

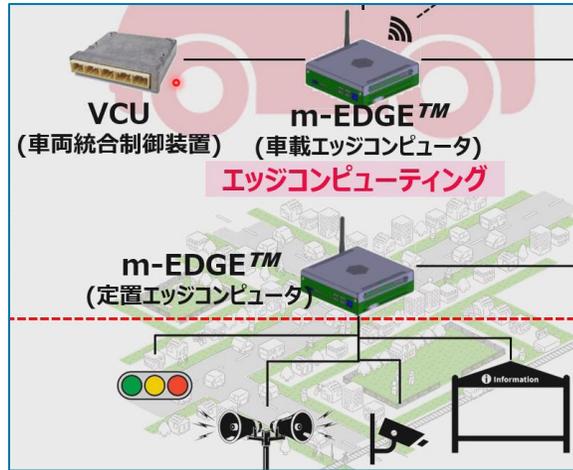


Desktop GPUs  
(and Cloud GPUs)



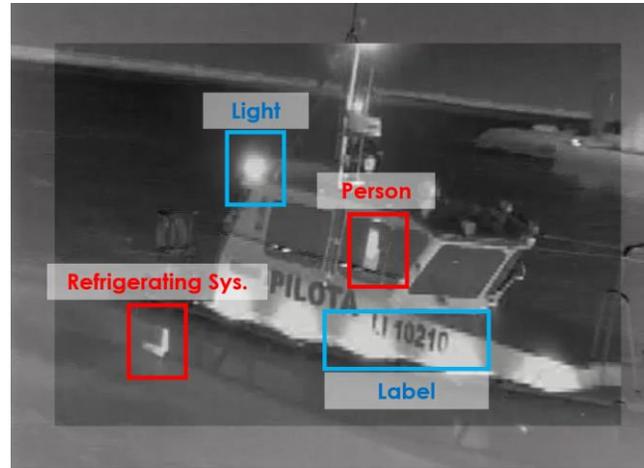
Embedded GPUs

# Wide Range of Applications



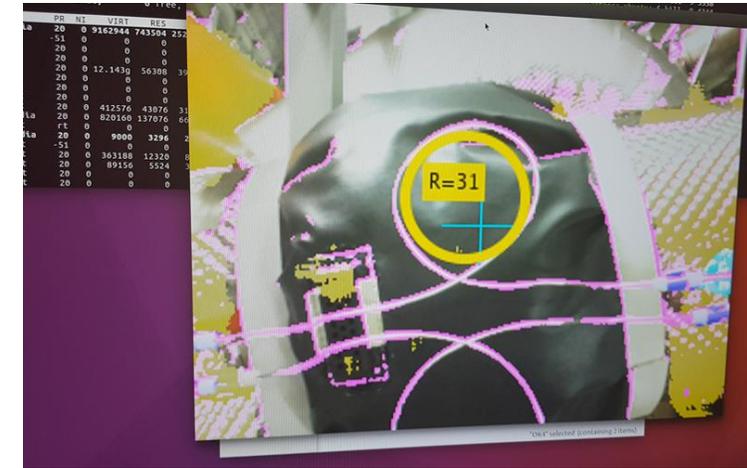
## DENSO Ten

Vehicle detection and traffic flow analysis



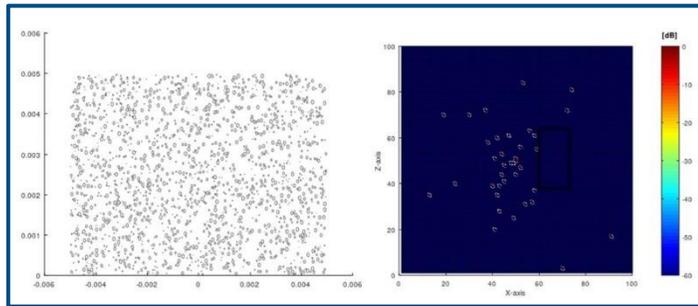
## Drass

Maritime object tracking on NVIDIA GPU



## Airbus

Aircraft inspection on Jetson



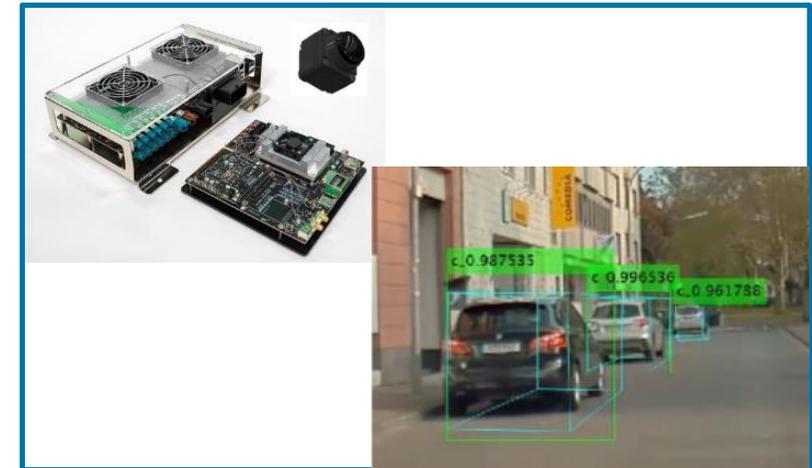
## M&R Technology

Medical Ray-Tracing Algorithms



## JINGCE Electronics

Display measurement instrument on Jetson



## Continental

Computer Vision Stack for AD

# Drass Develops Deep Learning System for Real-Time Object Detection in Maritime Environments

## Challenge

Help ship operators monitor sea environments and detect objects, obstacles, and other ships

## Solution

Create an object-detection deep learning model that can be deployed on ships and run in real time

## Results

- Data labeling automated
- Development time reduced
- Flexible and reproducible framework established



First day of object detection tests with optronic system prototype.

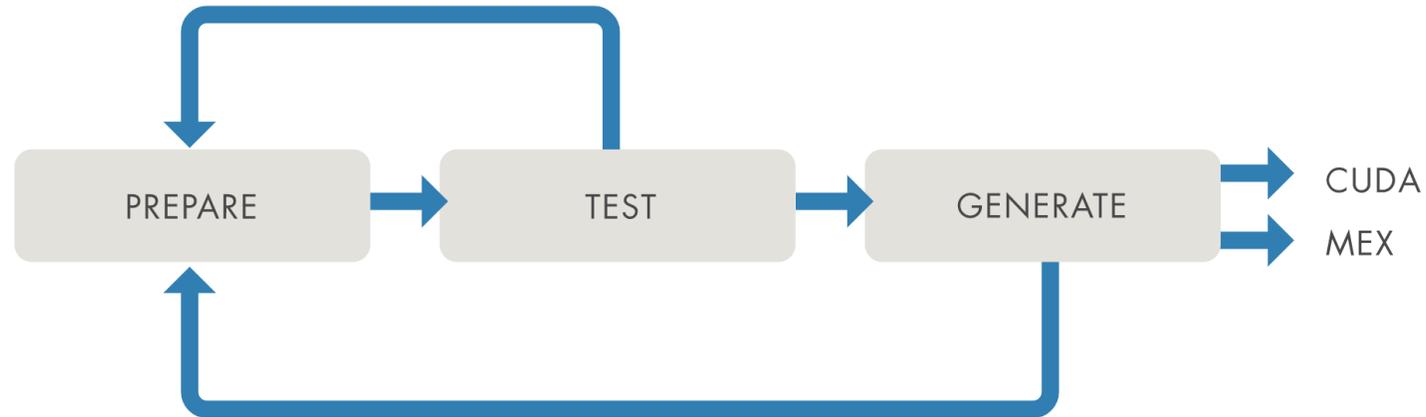
*“From data annotation to choosing, training, testing, and fine-tuning our deep learning model, MATLAB had all the tools we needed—and GPU Coder enabled us to rapidly deploy to our NVIDIA GPUs even though we had limited GPU experience.”*

*- Valerio Imbriolo, Drass Group*

# Agenda

- Why translate MATLAB to CUDA?
- What is GPU Coder?
- **Workflow for using GPU Coder**
- Deep Learning workflow
- Wrap-up

# Working with GPU Coder: Three-Step Workflow



## **Prepare** your MATLAB algorithm for code generation

- Make implementation choices
- Use supported language features

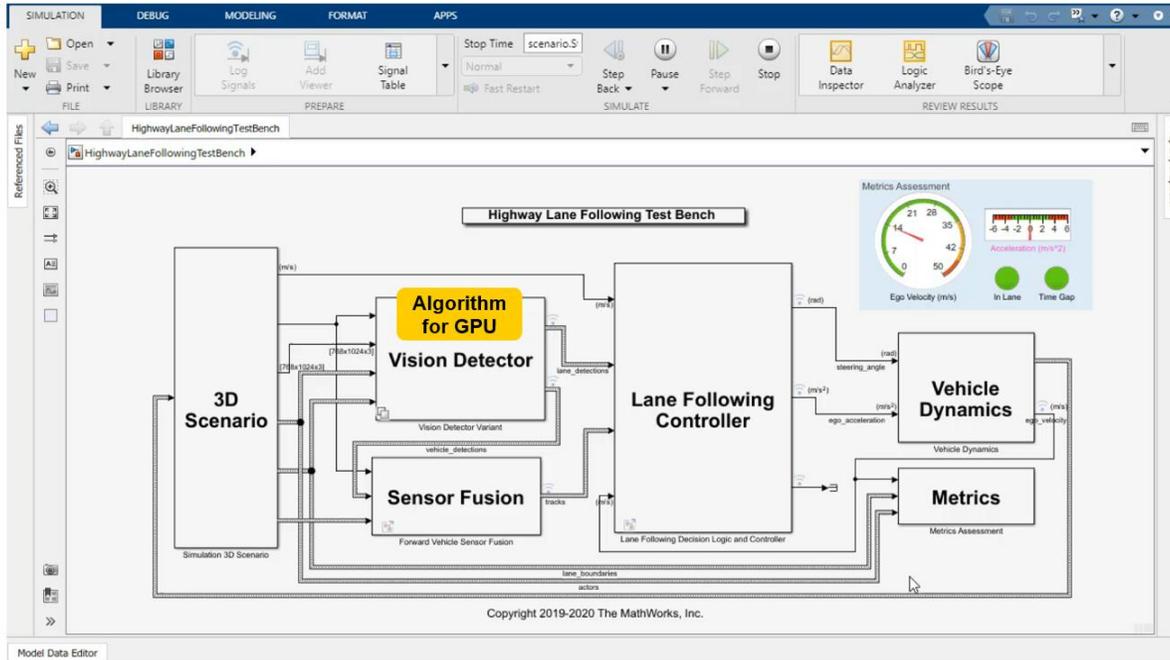
## **Test** if your MATLAB code is ready for code generation

- Validate that MATLAB program generates code
- Accelerate execution of user-written algorithm

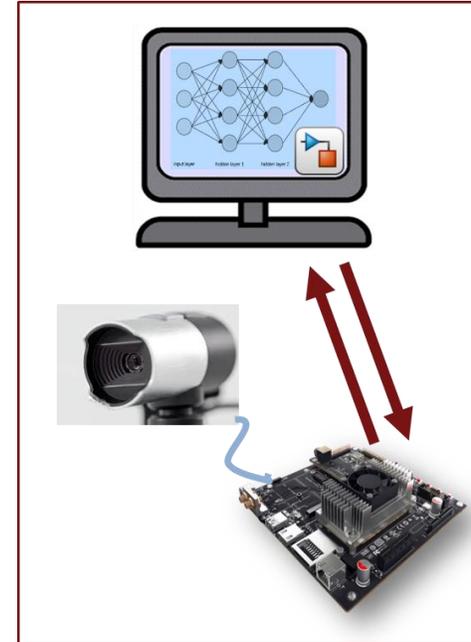
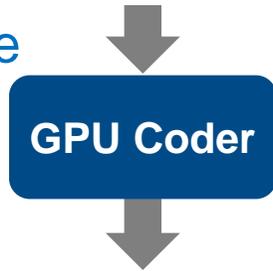
## **Generate** source code or MEX for final use

- Iterate your MATLAB code and use profiling tools to help optimize further
- Implement as source, executable, or library

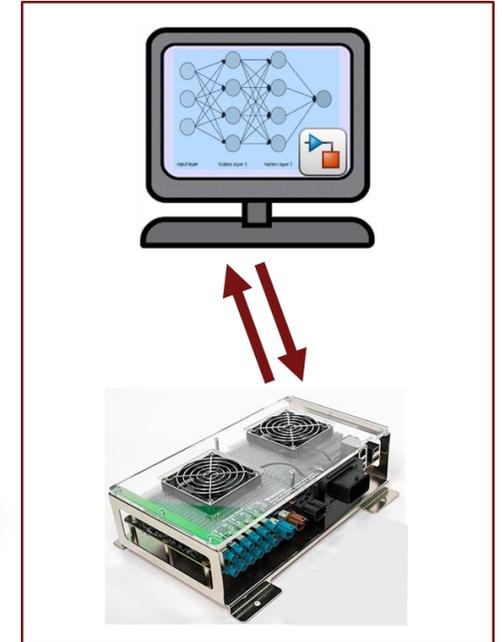
# GPU Coder for Simulink Workflows



Deploy Standalone Application

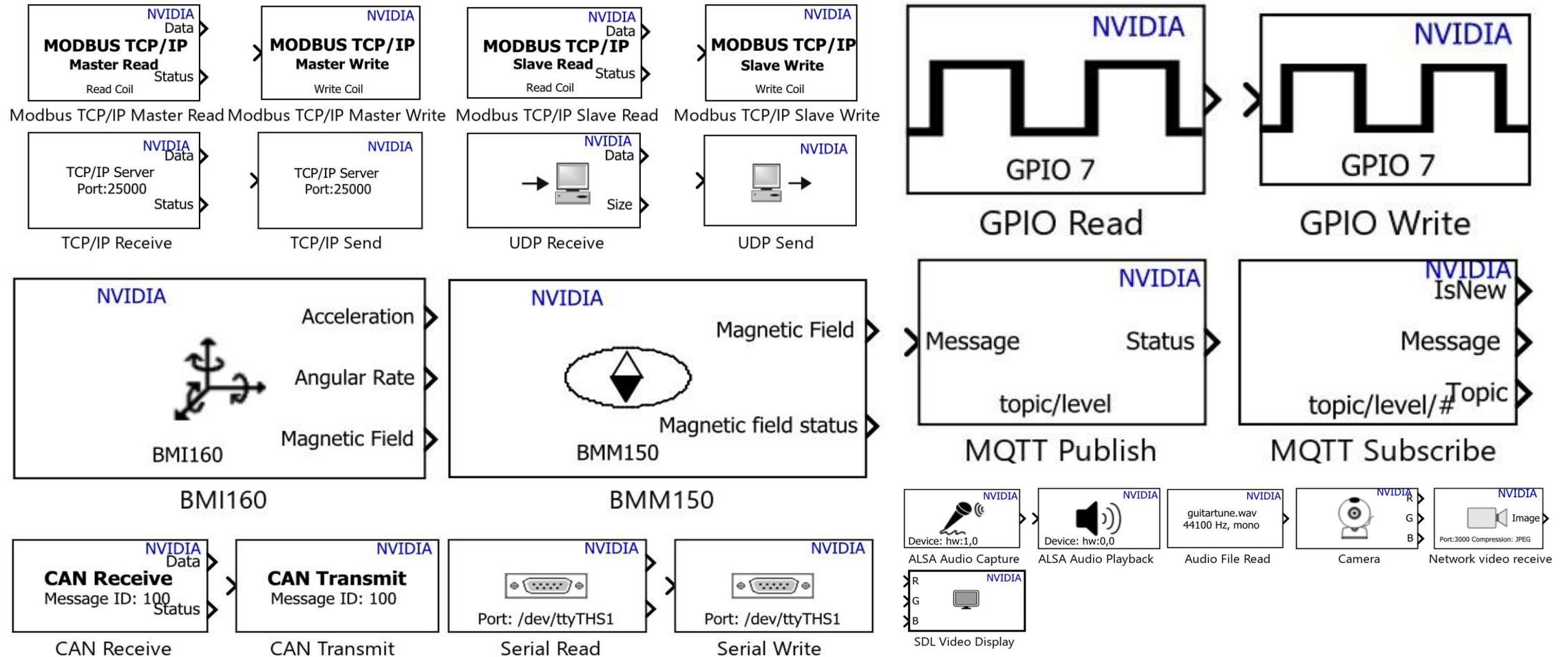


Access Peripheral from Simulink



Software-in-Loop, Processor-in-Loop, External Mode

# NVIDIA Peripheral Support – block library



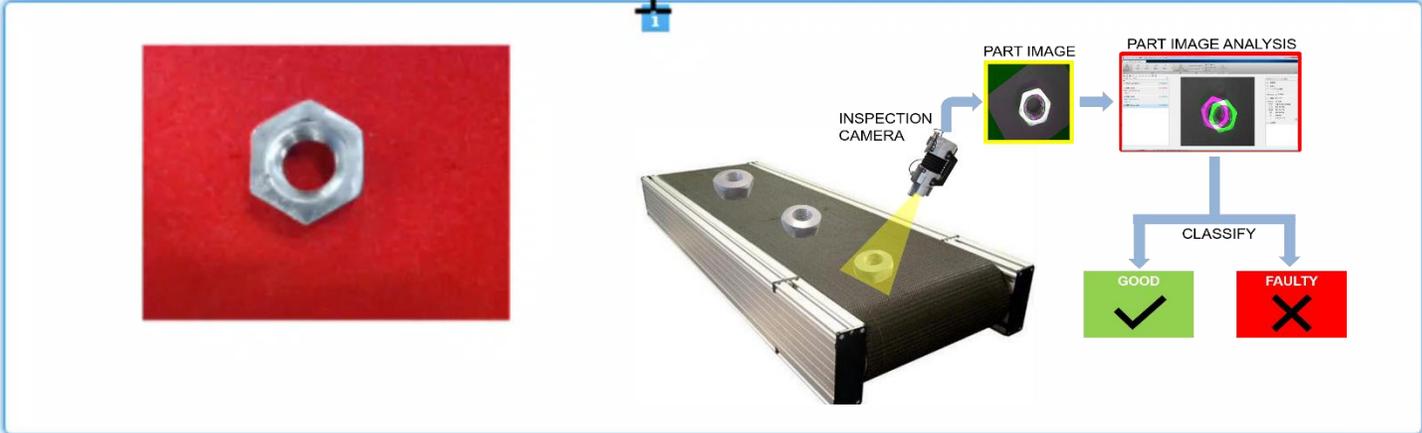
# Using GPU Coder with Automated Optical Inspection Example

tb\_myNDNet\_Jetson.mlx x myNDNet\_Preprocess.m x myNDNet\_Postprocess.m x targetFunction.m x +

## CNN and post-processing

Read new image captured by `webCam`. The images used to train network were `pre-processed` image to fit in input image size of network. Original image captured by `webCam` is like this;

```
50 img = imread('testImg.png');  
51 figure, imshow(img)
```



```
52 % So, we need to extract where the nuts are as an ROI before passing image  
53 % to the network.
```

## Pre-processing for captured image from webcam

myNDNet\_Preprocess is to extract ROI by using traditional image processing. Since generated C code from this m code will be

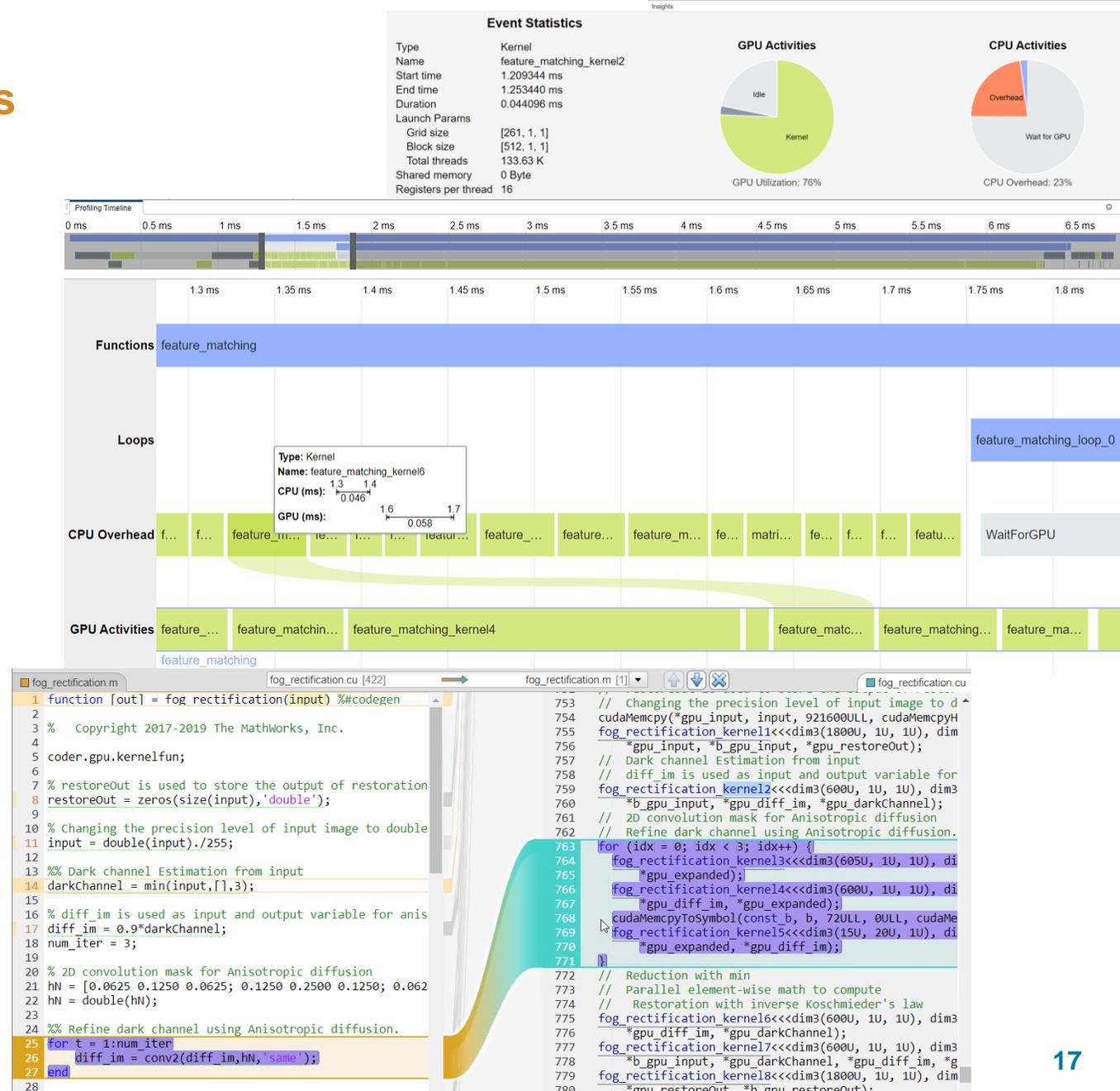
## Example in MATLAB

- Generate MEX file and accelerate in MATLAB
- Generate DLL/LIB/EXE and deploy
  - Desktop
  - Cloud
  - Embedded
- Generate source code and integrate into existing project
- Connected with Jetson board using MATLAB Coder Support Package for NVIDIA Jetson and NVIDIA DRIVE Platforms
- File transfer to/from Jetson
- Peripheral support
  - GPIO, audio input/output, USB/CSI camera, webcam, LiDAR sensors, ...
- Processor-in-loop (PIL) testing

# How to Optimize Further?

## Profiling and Bidirectional Traceability Tools

- Use the GPU Coder Performance Analyzer to profile the generated CUDA code
  - Identify bottlenecks & opportunities to optimize performance
- Use bidirectional traceability to map to/from CUDA code back to MATLAB code
  - Helps you to understand how CUDA kernels are created from your MATLAB algorithms



# Using the GPU Coder Performance Analyzer

PAToolDemoUsingSURF.mlx SurfDetect1.m PAToolDemoUsingSURF.m +

## Performance Analyzer

*Surf Features are an important detection mechanism in computer vision problems to help facilitate registration and object detection algorithms. Let's take a look. In this peppers example, each feature is marked in the image and contains a coordinate represented by the position on the image, a scale represented by the circle diameters, and orientation represented by the vector in the circles, and by the sign of Laplacian. This implementation designed for deployment takes a little over 1 second to run on a modern laptop. We will generate CUDA code for this algorithm to demonstrate how the Performance Analyzer Tool can help you modify your MATLAB code to achieve better performance.*

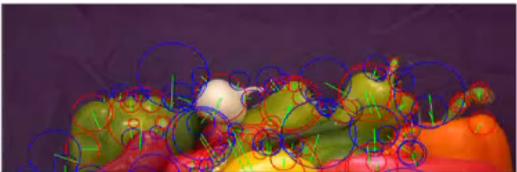
```
1  imageFile = 'peppers.png';
2  img = imread(imageFile);
3  tic;
4  intPoints = SurfDetect(img);
5  toc
```

Elapsed time is 1.240863 seconds.

```
6
7  % display a Surf Point structure
8  intPoints(3)
```

```
ans = struct with fields:
        x: 83.8932
        y: 203.8722
        scale: 2.0504
        orientation: 2.8450
        laplacian: 1
```

```
9  DrawIpoints(imageFile,intPoints);
```



# CUDA Kernel Optimizations Heuristics

- **Memory Allocation**
  - Supports various GPU memory spaces (local, global, shared, constant)
- **Data Transfer Minimization**
  - Analyzes data dependency between the CPU and GPU partitions to determine minimum set of locations where data must be copied between CPU and GPU using `cudaMemcpy`
  - Data shared between the CPU and GPU is allocated on GPU memory using `cudaMalloc` or `cudaMallocManaged` (for unified memory, GPU Coder determines minimum number of `cudaDeviceSync` calls)

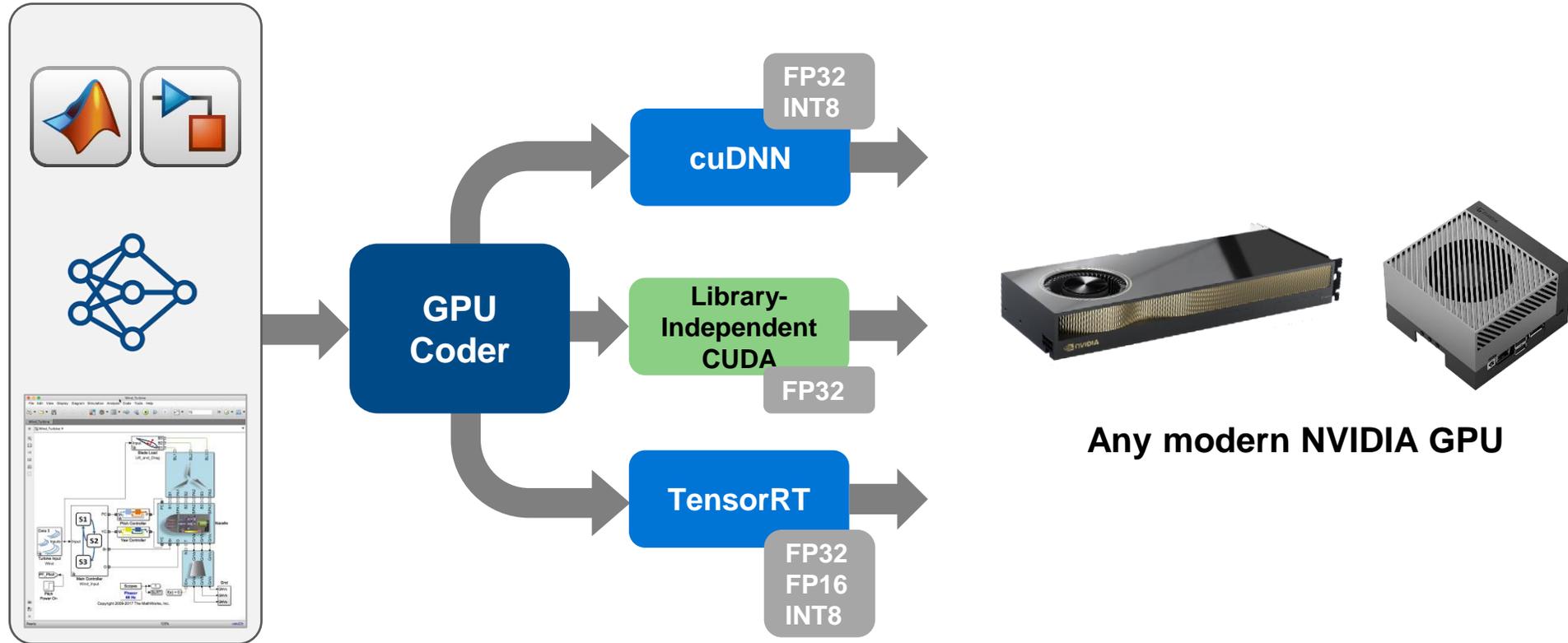
# Optimizations for Generated CUDA Code

- Accelerated library support
  - cuFFT, cuBLAS, cuSolver, Thrust, cuDNN, and TensorRT
- Design patterns (apply manually):
  - [Stencil processing](#) can be used for operations such as convolution, median filtering, and finite element methods
  - [Matrix-matrix processing](#) can be used for operations such as sum of absolute differences (SAD) and sum of squared differences (SSD)
  - [Reductions](#) can be used for operations such as vector summation

# Agenda

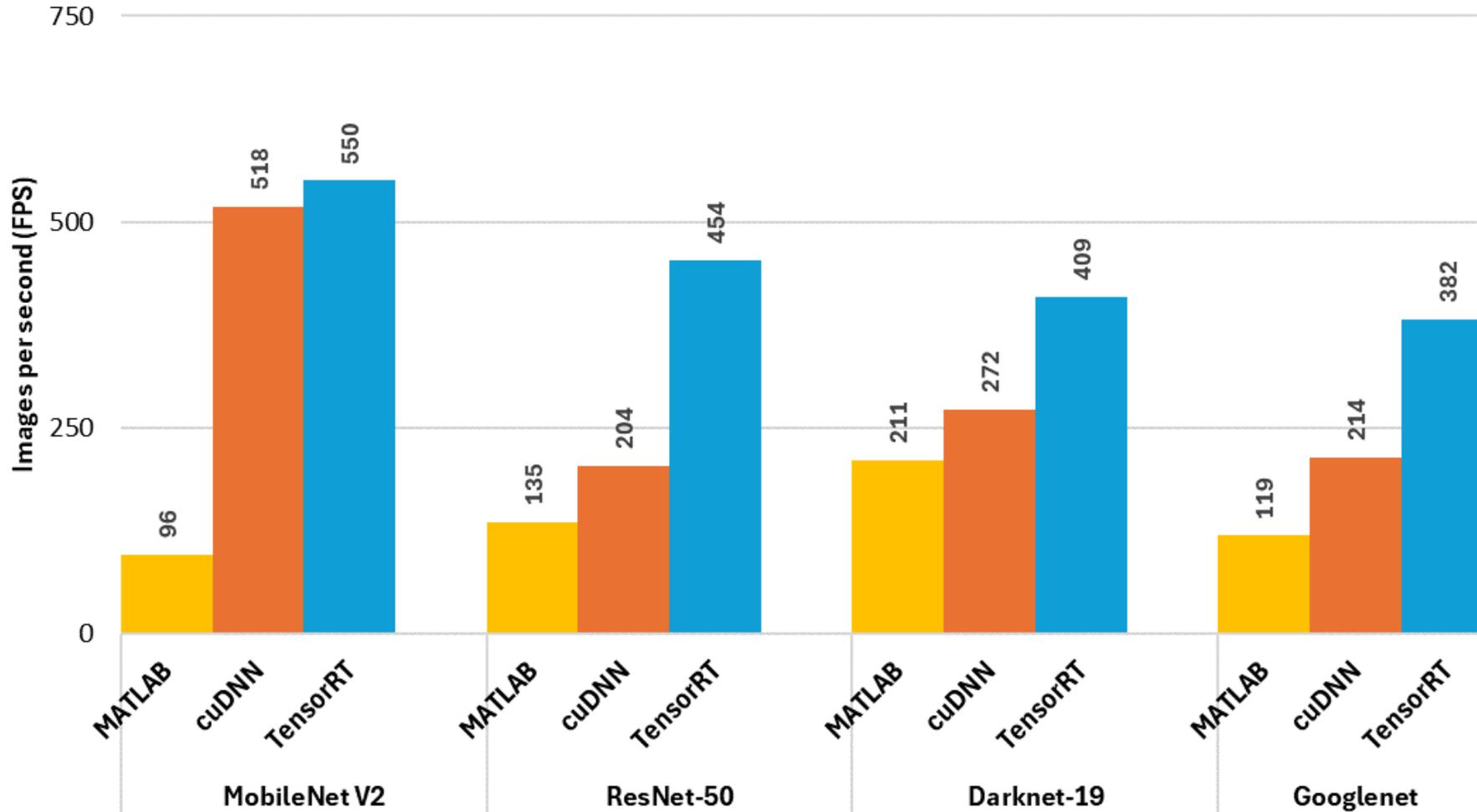
- Why translate MATLAB to CUDA?
- What is GPU Coder?
- Workflow for using GPU Coder
- **Deep Learning workflow**
- Wrap-up

# Deep Learning Code Generation



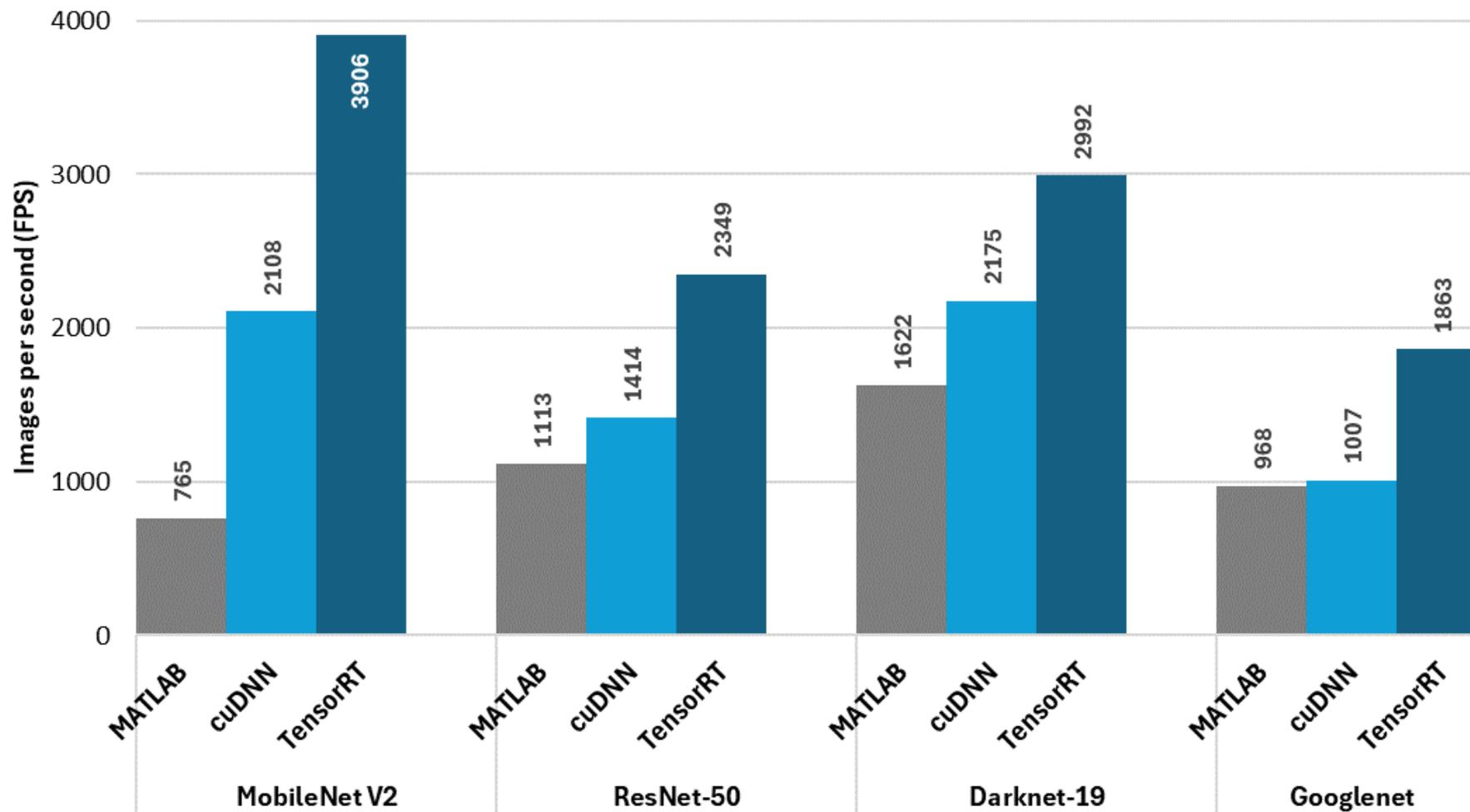
# Performance of Generated Code from GPU Coder

## Batch Size 1



# Performance of Generated Code from GPU Coder

## Batch Size 8



# Agenda

- Why translate MATLAB to CUDA?
- What is GPU Coder?
- Workflow for using GPU Coder
- Deep Learning workflow
- **Wrap-up**

# Key Takeaways

- GPU Coder generates CUDA code from MATLAB & Simulink
  - Deploy signal processing, radar processing, image processing, and computer vision algorithms
  - Deploy complete AI applications (with Deep Learning Toolbox)
- Integrate handwritten CUDA code for use within MATLAB/Simulink and generated code
- Accelerate MATLAB & Simulink simulations using MEX-files

# Shipping Examples



Products Solutions Academia Support Community Events



## Help Center

Search Help Center

Help Center

### CONTENTS

- « Documentation Home
- « Examples
- « Code Generation

#### Category

Fixed-Point Designer

#### GPU Coder

- Get Started with GPU Coder 2
- Kernel Creation 13
- Performance 5
- Deep Learning with GPU Coder 28
- Deployment 21
- HDL Coder
- HDL Verifier
- IEC Certification Kit
- MATLAB Coder
- Simulink Code Inspector
- Simulink Coder
- Simulink PLC Coder

#### Type

- All 66
- MATLAB 51
- Simulink 15

#### New Example

- Added in R2024a

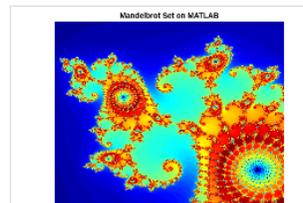
Documentation Examples Functions Apps Videos Answers

Trial Software Product Updates

## GPU Coder – Examples

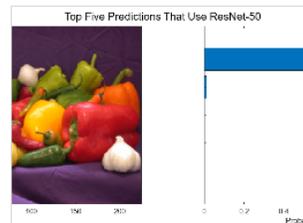
<https://www.mathworks.com/help/gpucoder/examples.html>

### Get Started with GPU Coder



#### GPU Code Generation: The Mandelbrot Set

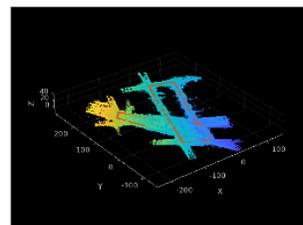
Generate CUDA® code from a simple MATLAB® function by using GPU Coder™. A Mandelbrot set implementation by using standard



#### Code Generation for Deep Learning Networks

Get started with CUDA code generation for image classification networks such as ResNet.

### Kernel Creation



#### Build a Map from Lidar Data



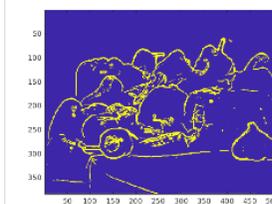
#### Feature Extraction Using



#### Feature Matching



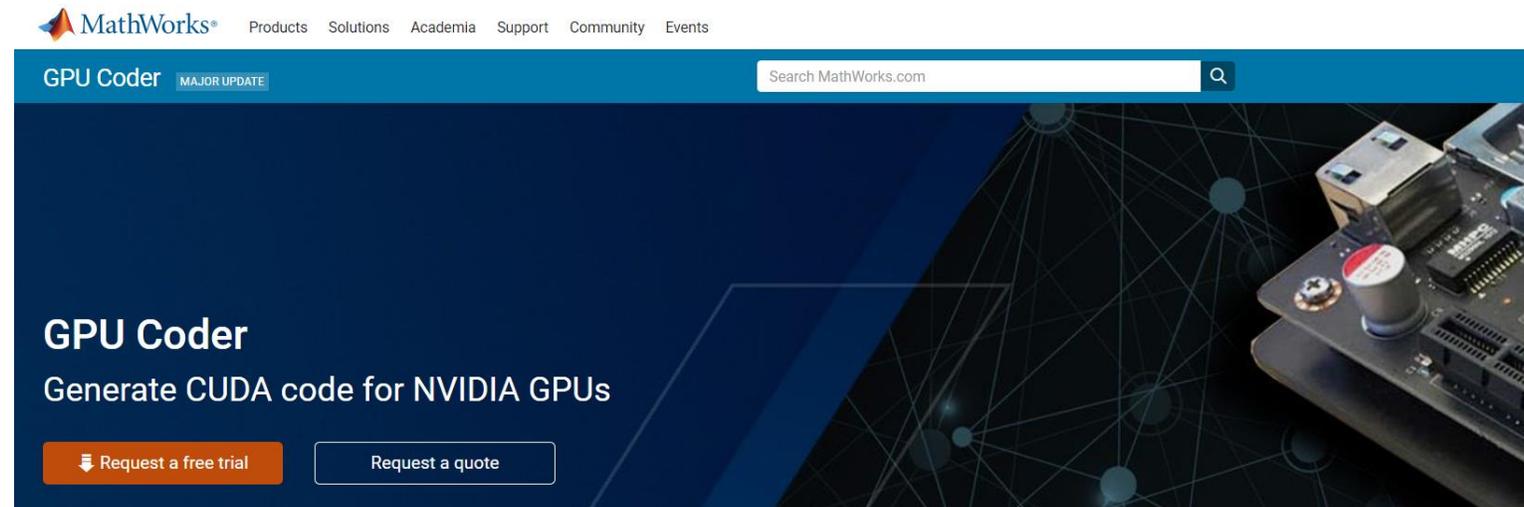
#### Lane Detection on the GPU



#### Edge Detection with Sobel

# Next Steps

- To learn more, including playing with a trial license, visit:
  - [www.mathworks.com/products/gpu-coder](http://www.mathworks.com/products/gpu-coder)
- Read [white paper](#) on generating CUDA code from MATLAB
- Attend 4-hour hands-on workshop
  - Arrange with your account manager



The screenshot shows the MathWorks website for the GPU Coder product. The header includes the MathWorks logo and navigation links for Products, Solutions, Academia, Support, Community, and Events. The main heading is "GPU Coder" with a "MAJOR UPDATE" badge. Below the heading is the text "Generate CUDA code for NVIDIA GPUs". There are two buttons: "Request a free trial" and "Request a quote". The background features a dark blue network diagram and a close-up of a GPU circuit board.



The screenshot shows a white paper titled "Generating CUDA Code from MATLAB: Accelerating Embedded Vision and Deep Learning Algorithms on GPUs". The background is a stylized illustration of a factory or industrial setting with purple and blue tones. There are three speed limit signs with the number "70" on them. A green button labeled "Read white paper" is visible in the bottom left corner.

# Q&A