

MATLAB EXPO 2019

Deploying Deep Neural Networks
to Embedded GPUs and CPUs

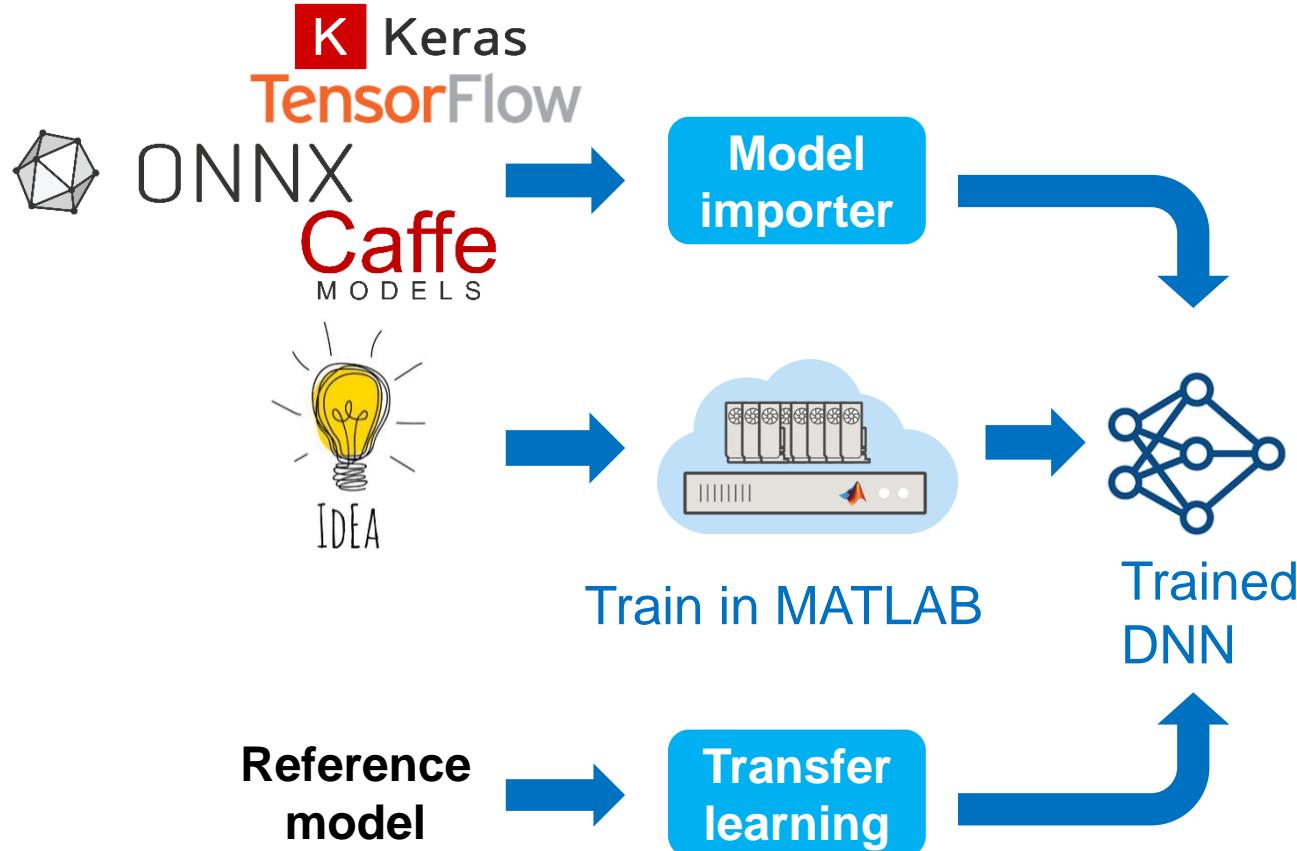
Steven Thomsett



Deep Learning Workflow in MATLAB

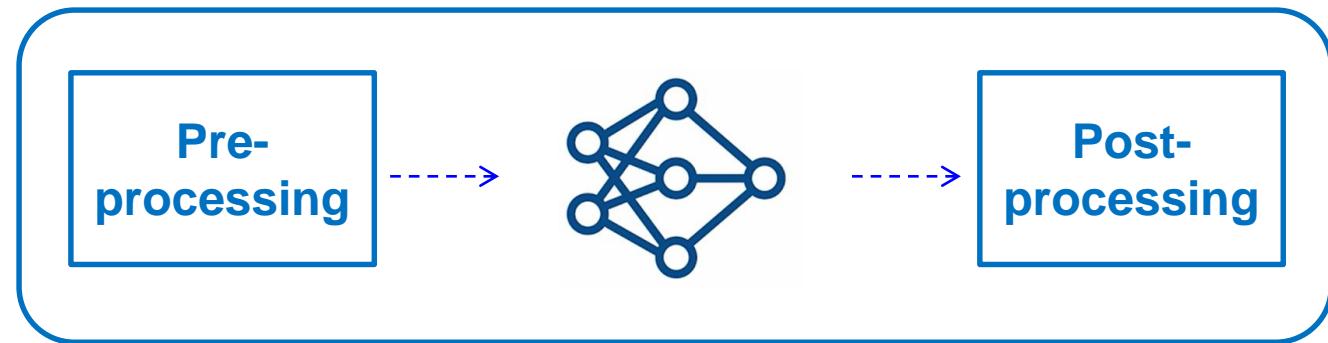


Deep Neural Network Design and Training

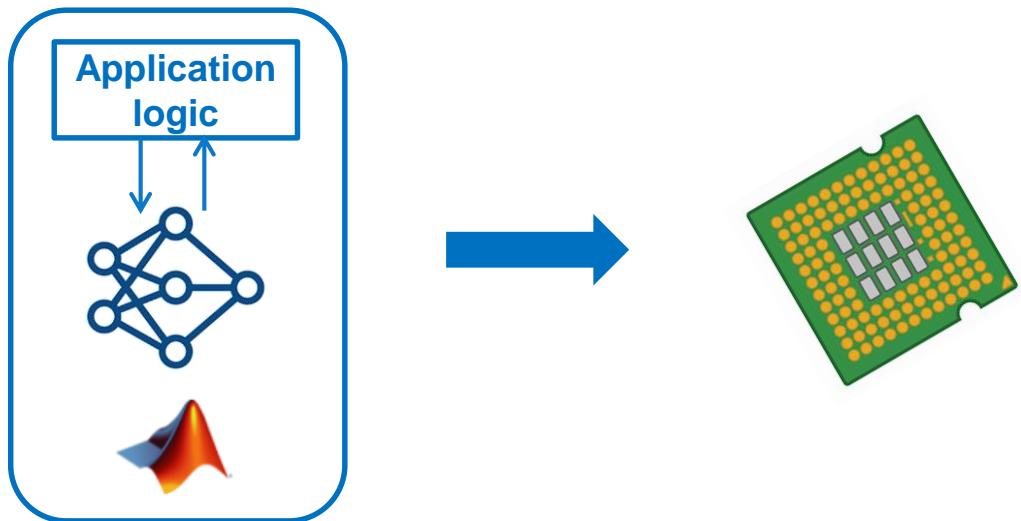
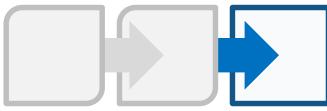


- **Design in MATLAB**
 - **Manage** large data sets
 - **Automate** data labeling
 - **Easy access** to models
- **Training in MATLAB**
 - **Acceleration** with GPU's
 - **Scale** to clusters

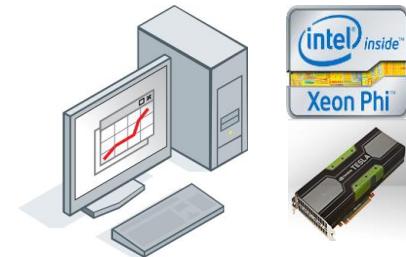
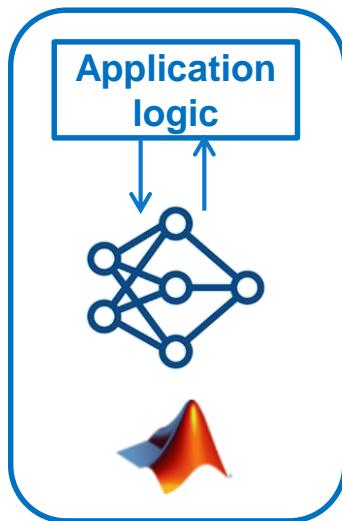
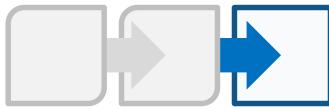
Application Design



Multi-Platform Deep Learning Deployment



Multi-Platform Deep Learning Deployment



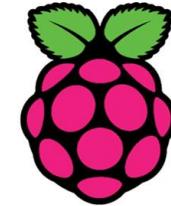
Desktop



Data Center



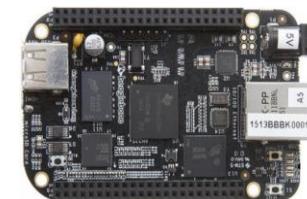
NVIDIA Jetson



Raspberry pi



Mobile



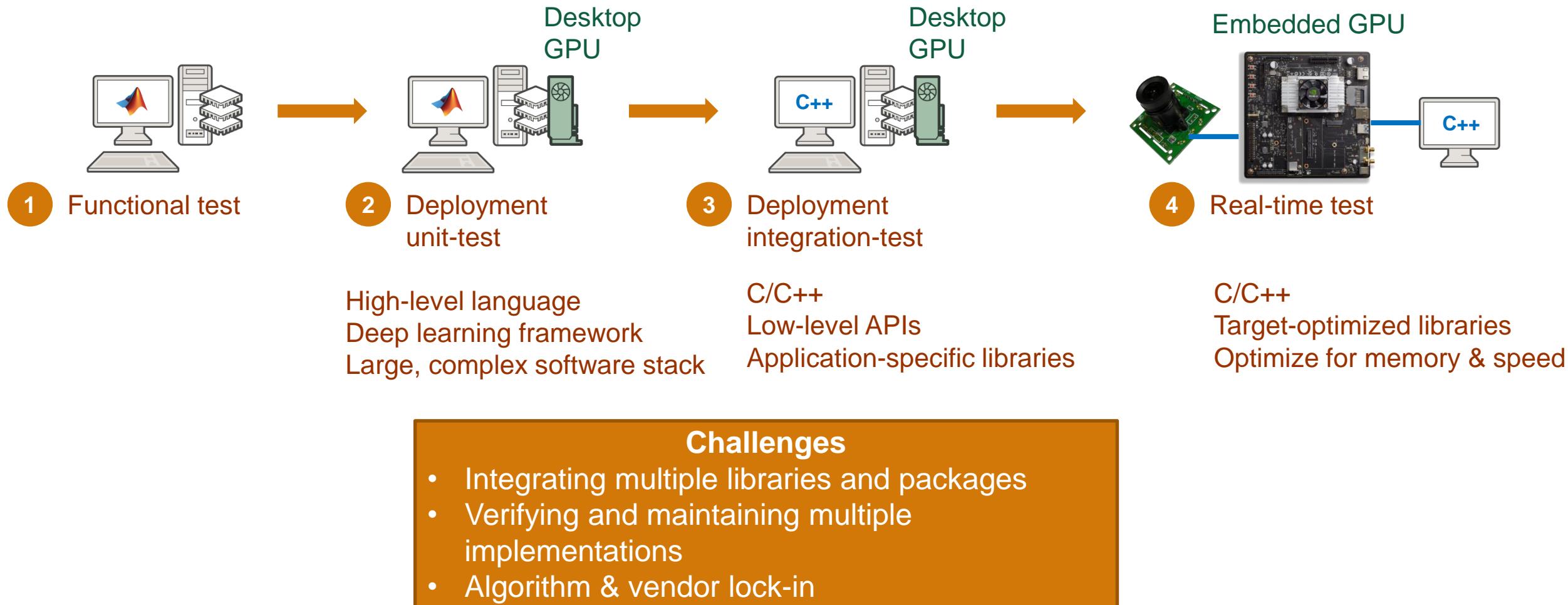
Beaglebone

...

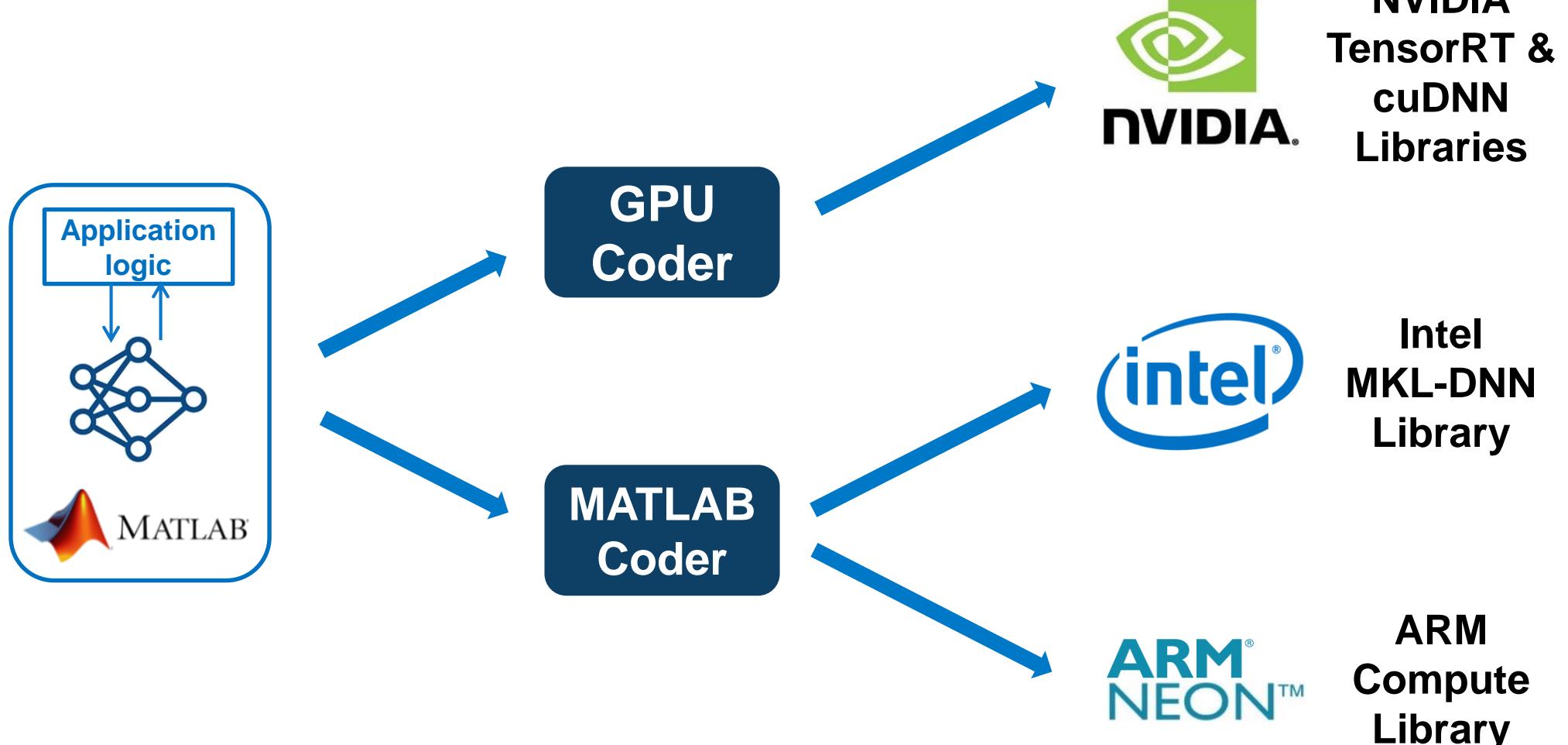
Embedded

Algorithm Design to Embedded Deployment Workflow

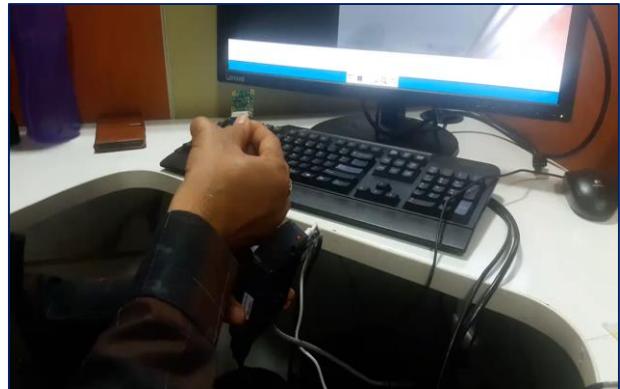
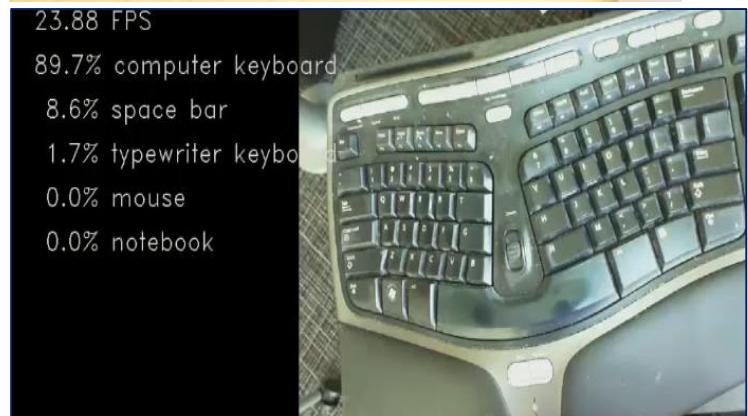
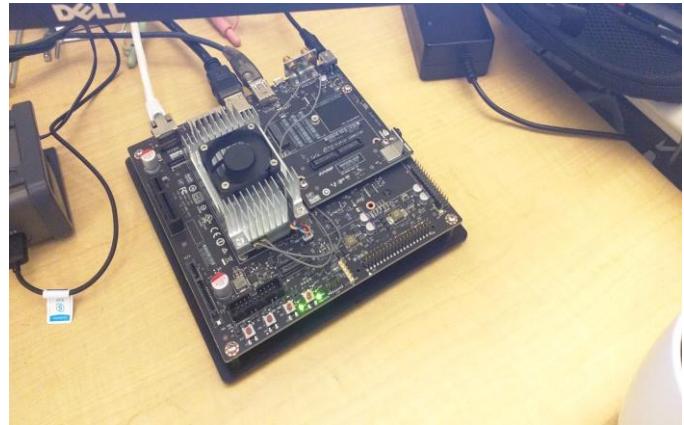
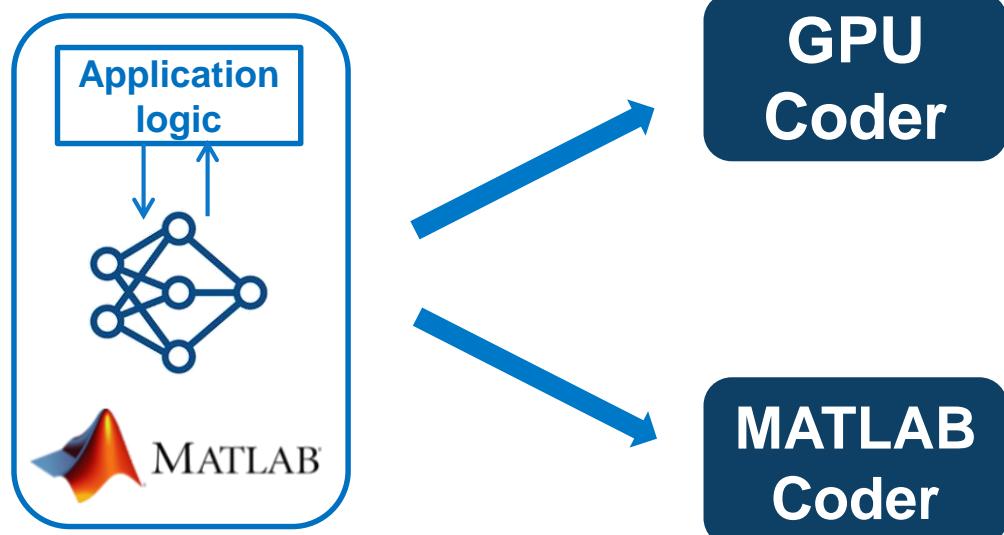
Conventional Approach



Solution: Use MATLAB Coder & GPU Coder for Deep Learning Deployment

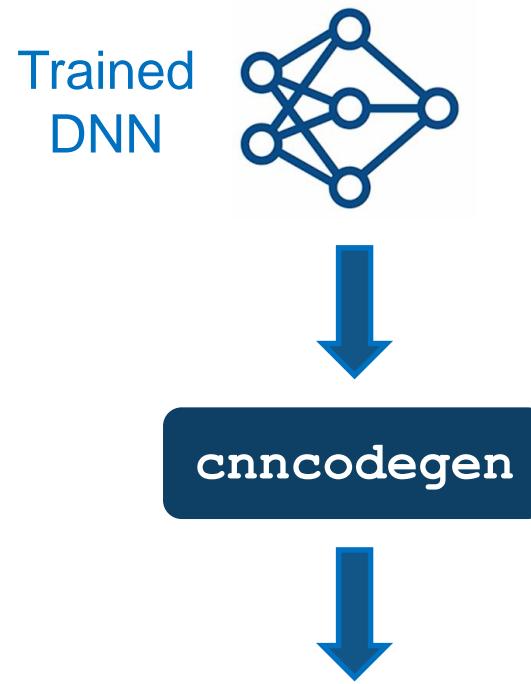


Solution: Use MATLAB Coder & GPU Coder for Deep Learning Deployment

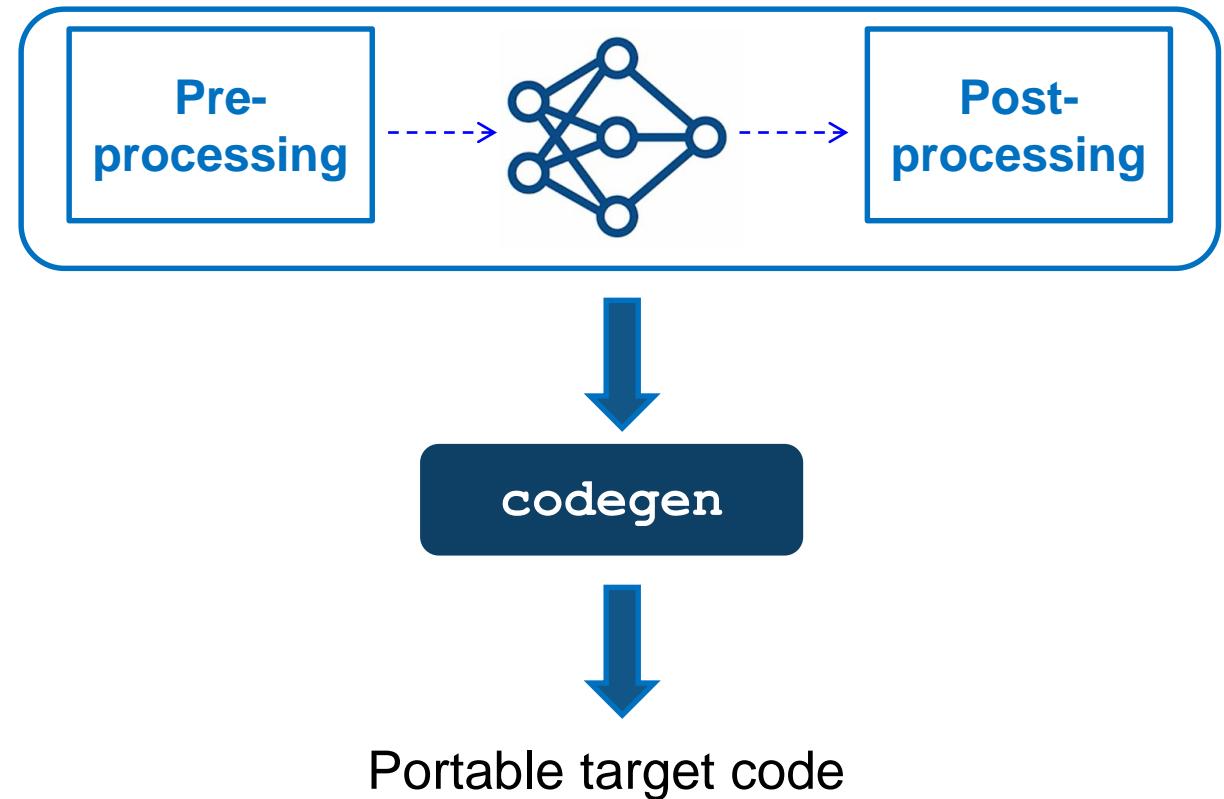


Deep Learning Deployment Workflows

INFERENCE ENGINE DEPLOYMENT

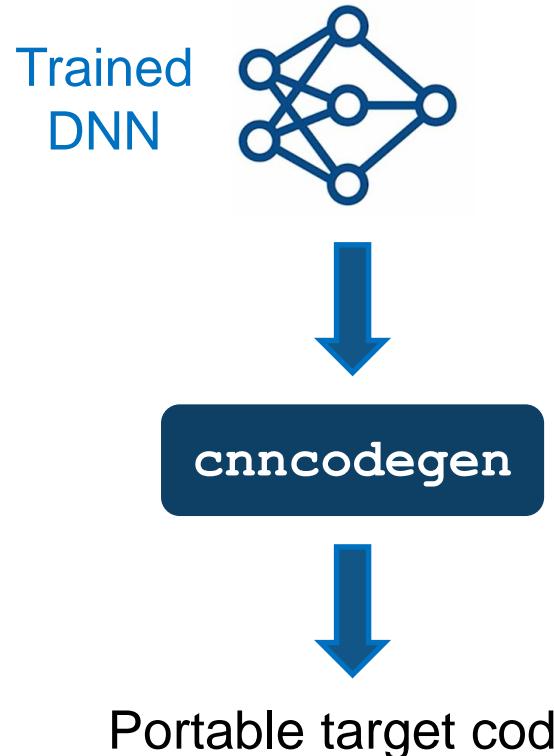


INTEGRATED APPLICATION DEPLOYMENT



Workflow for Inference Engine Deployment

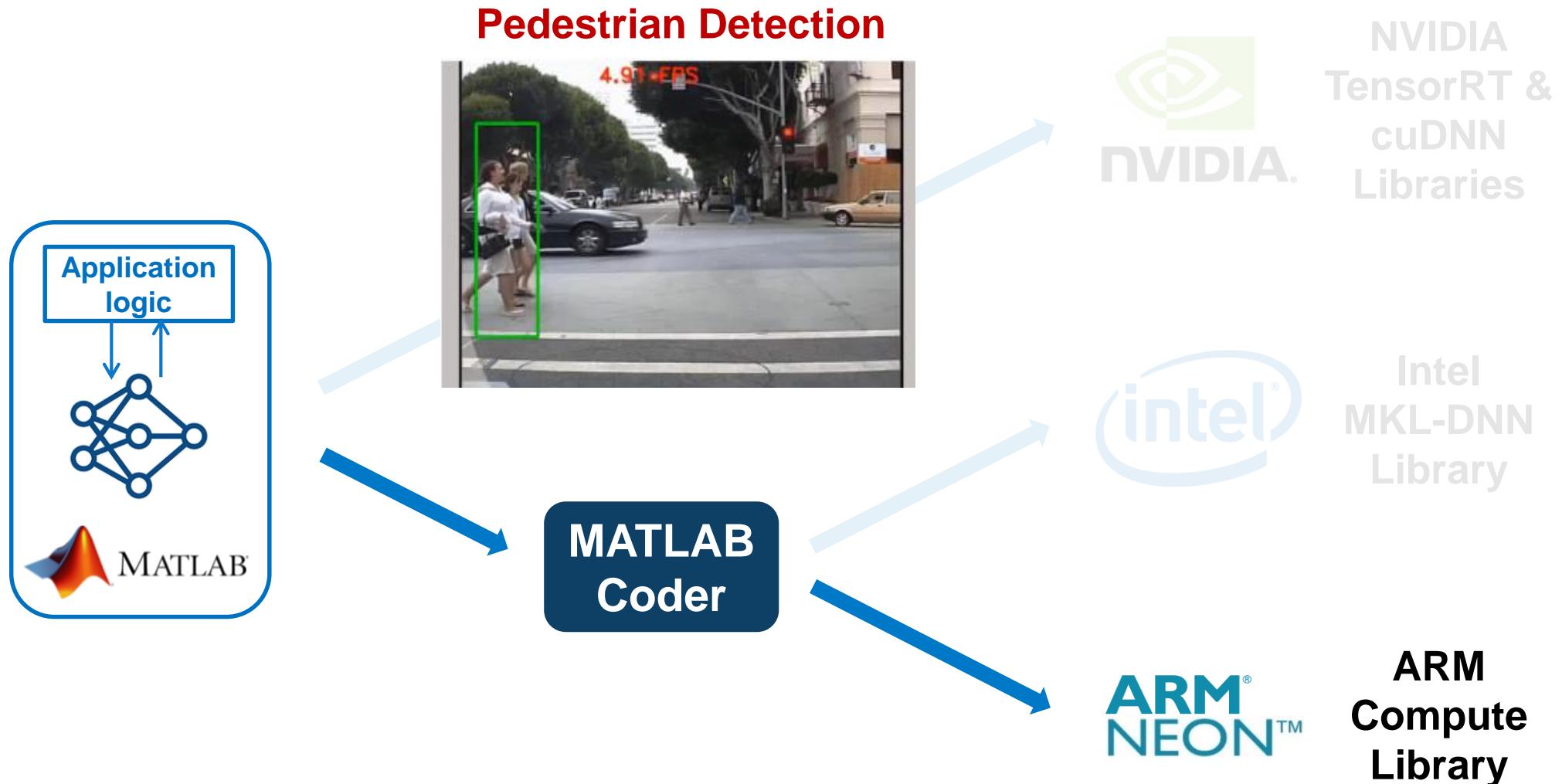
INFERENCE ENGINE DEPLOYMENT



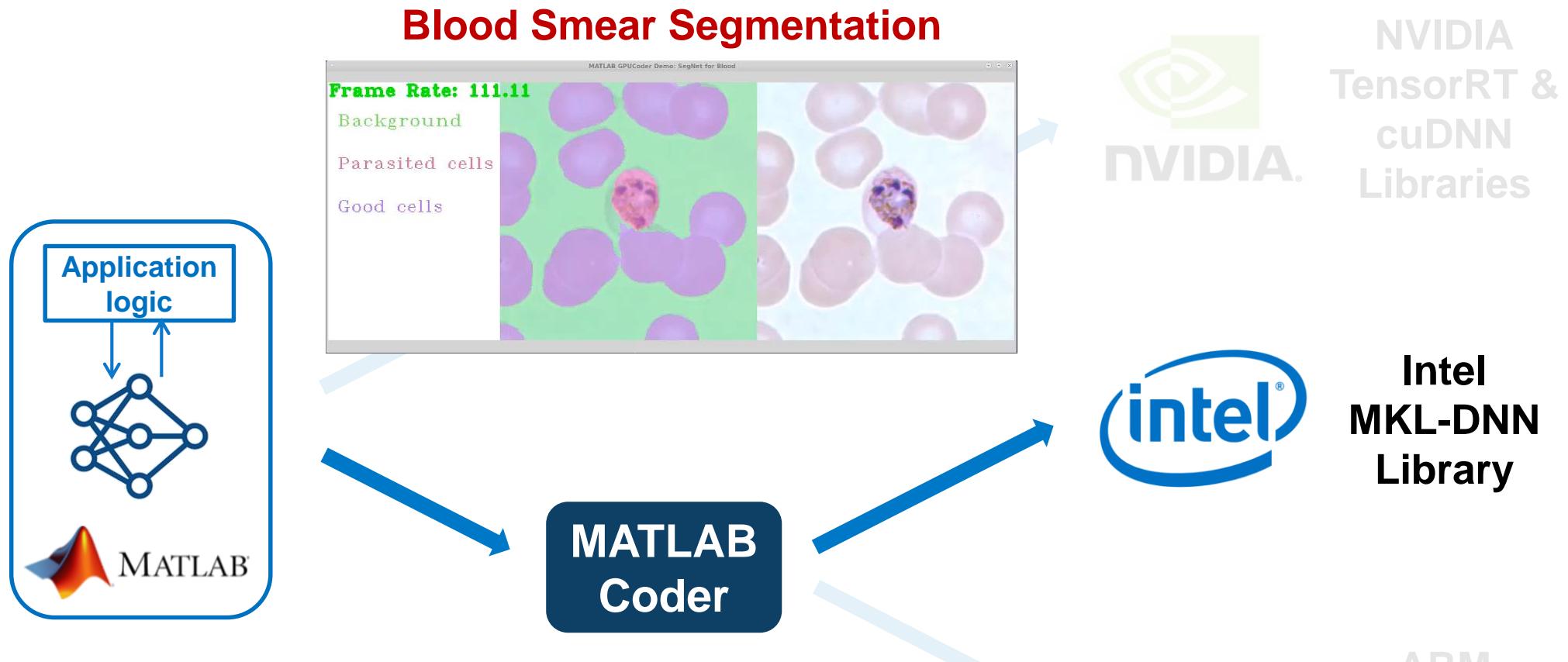
Steps for inference engine deployment

1. Generate the code for trained model
`>> cnncodegen (net, 'targetlib', 'arm-compute')`
2. Copy the generated code onto target board
3. Build the code for the inference engine
`>> make -C ./codegen -f ...mk`
4. Use hand written main function to call inference engine
5. Generate the exe and test the executable
`>> make -C ./`

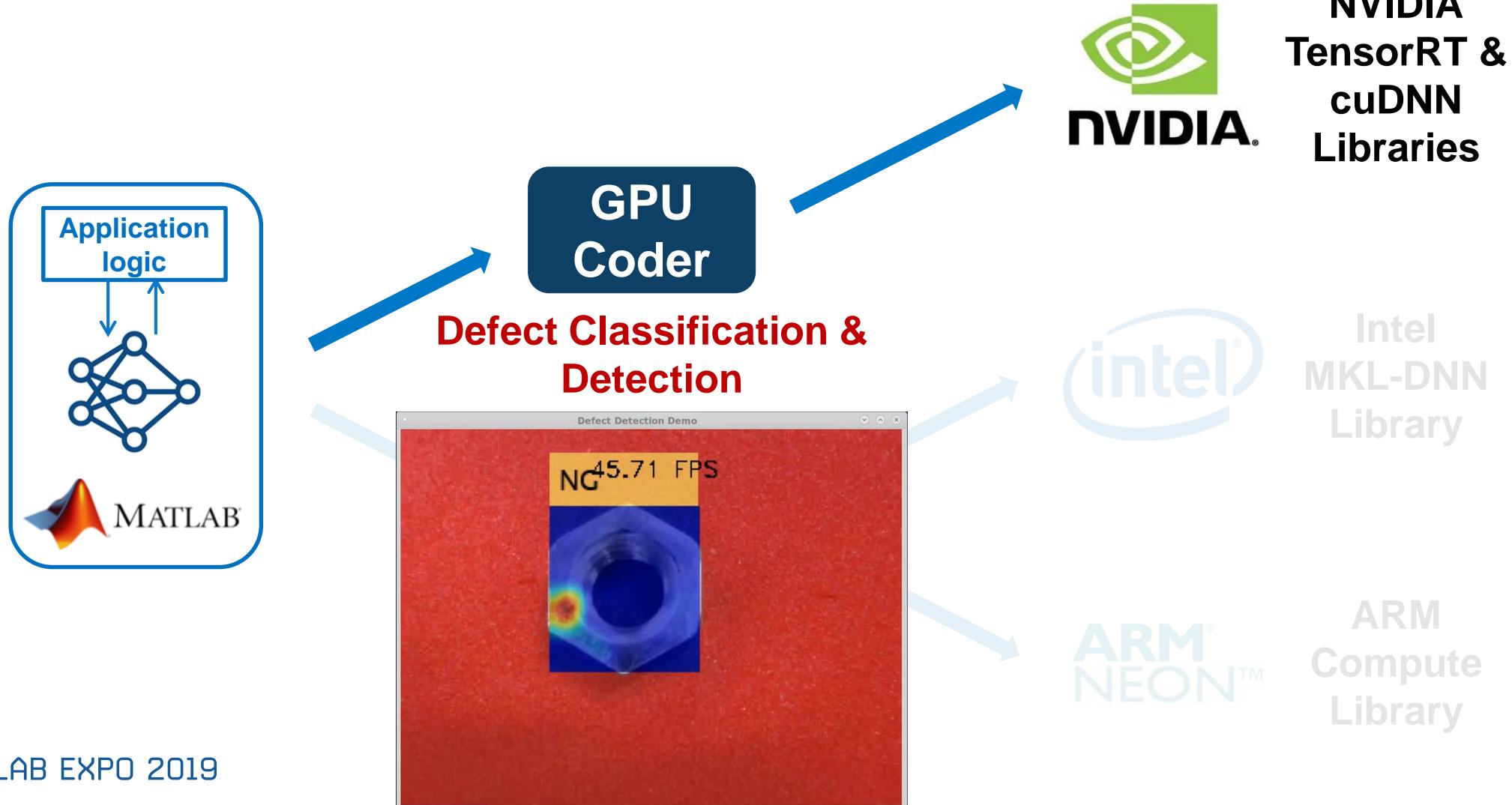
Deep Learning Inference Deployment



Deep Learning Inference Deployment



Deep Learning Inference Deployment





How is the Performance?



Target Libraries
NVIDIA
TensorRT &
cuDNN
Libraries



Intel
MKL-DNN
Library

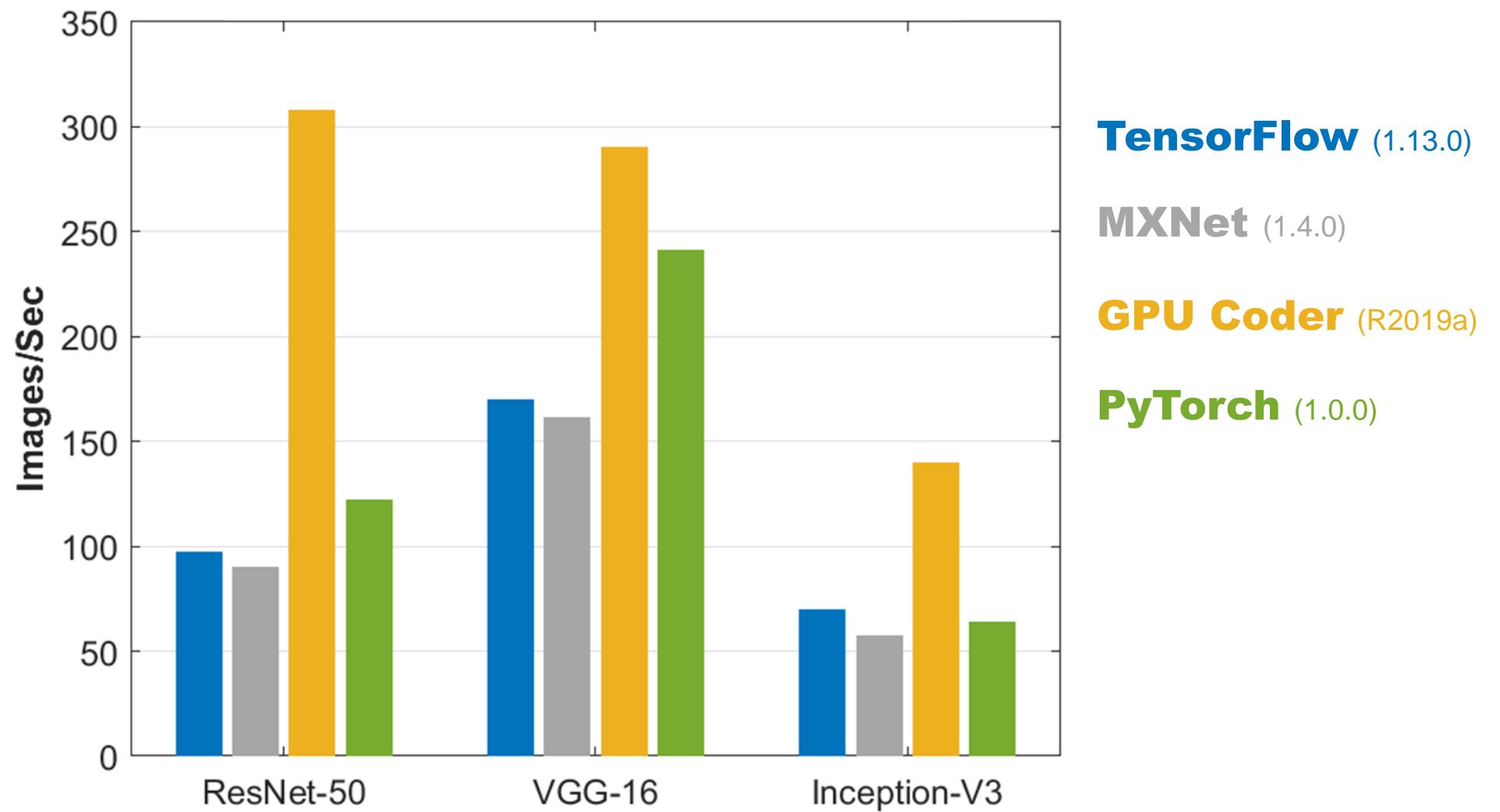


ARM
Compute
Library

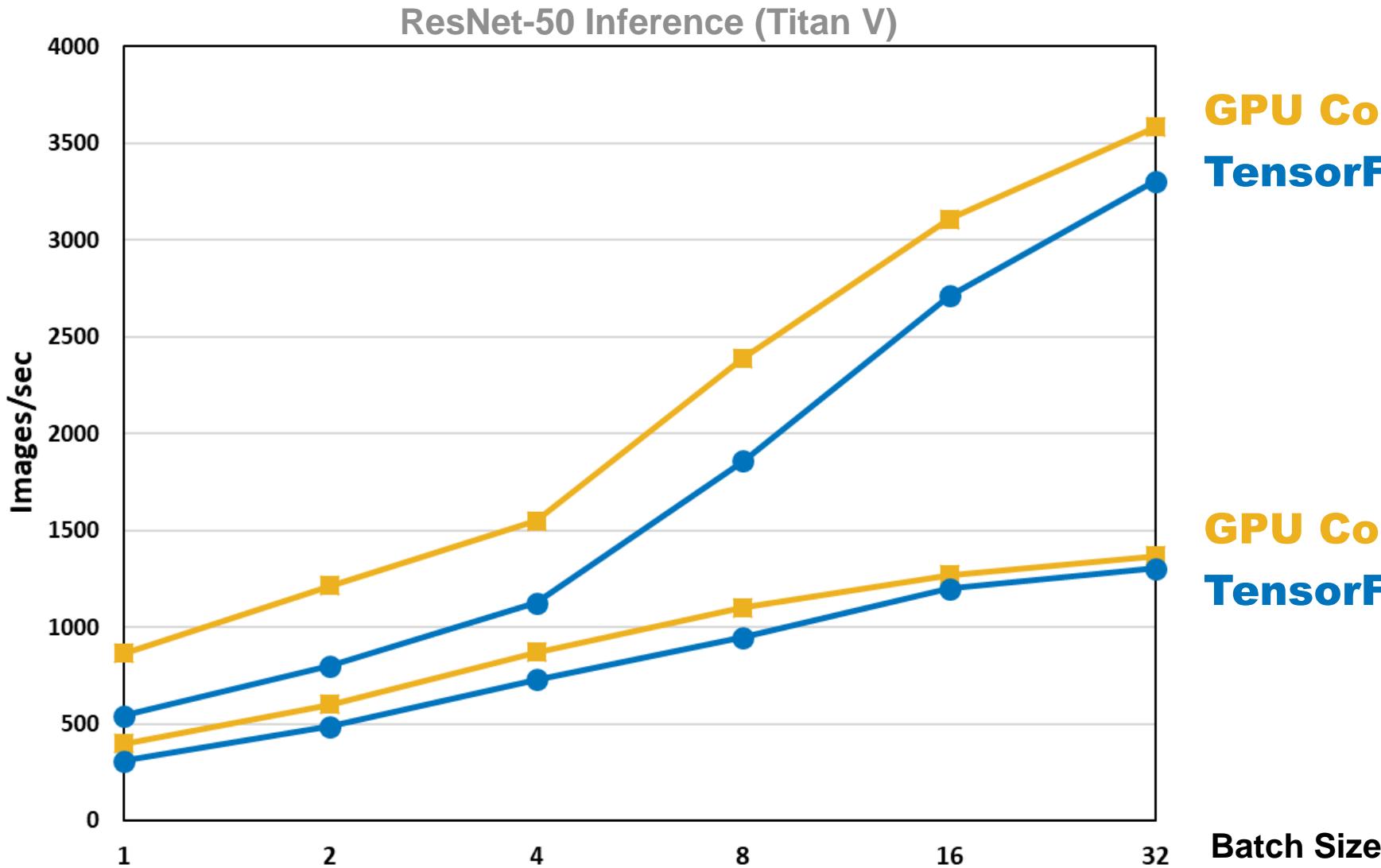
Performance of Generated Code

- CNN inference (ResNet-50, VGG-16, Inception V3) on Titan V GPU
- CNN inference (ResNet-50) on Jetson TX2
- CNN inference (ResNet-50 , VGG-16, Inception V3) on Intel Xeon CPU

Single Image Inference on Titan V using cuDNN



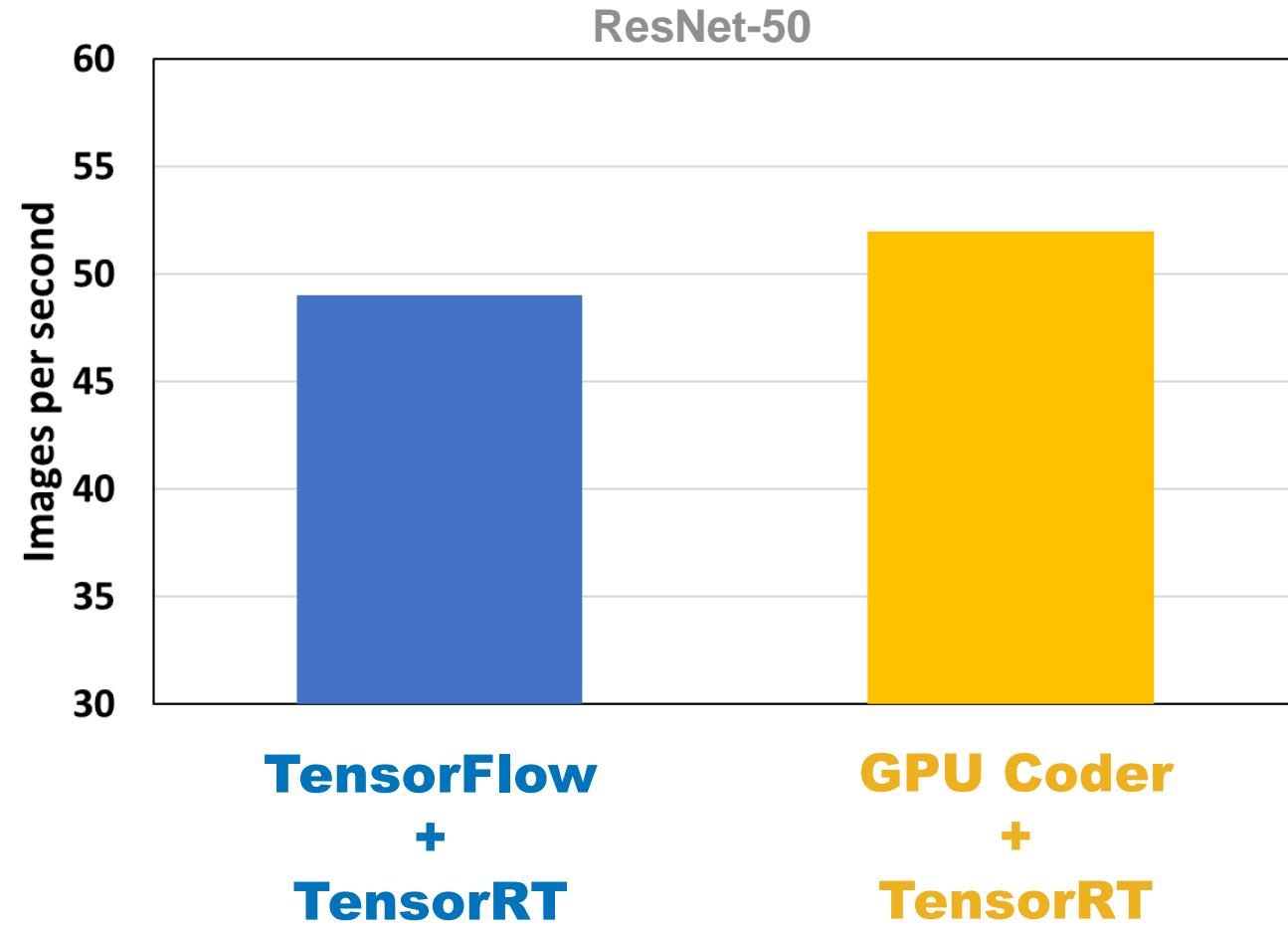
Even Stronger Performance with INT8 using TensorRT



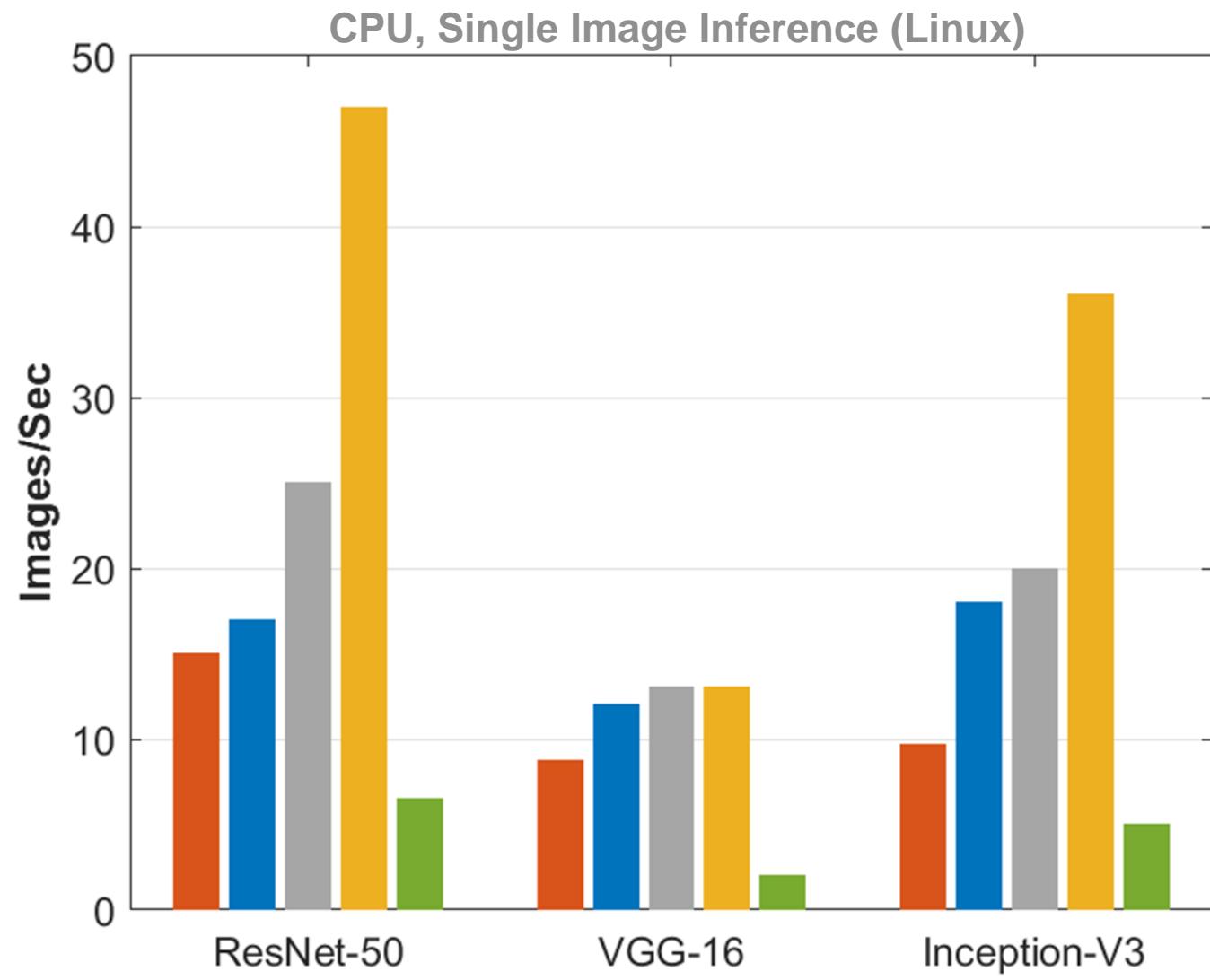
GPU Coder + TensorRT (INT8)
TensorFlow + TensorRT (INT8)

GPU Coder + TensorRT (FP32)
TensorFlow + TensorRT (FP32)

Single Image Inference on Jetson TX2



CPU Performance



MATLAB

Intel® Xeon® CPU 3.6 GHz - Frameworks: TensorFlow 1.6.0, MXNet 1.2.1, PyTorch 0.3.1

Brief Summary

DNN libraries are great for inference, ...

MATLAB Coder and GPU Coder generates code that takes advantage of:



NVIDIA® CUDA libraries, including TensorRT & cuDNN



Intel® Math Kernel Library for Deep Neural Networks
(MKL-DNN)



ARM® Compute libraries for mobile platforms

Brief Summary

DNN libraries are great for inference, ...

MATLAB Coder and GPU Coder generates code that takes advantage of:

**But, Applications
Require More than
just Inference**



NVIDIA® CUDA libraries, including TensorRT & cuDNN

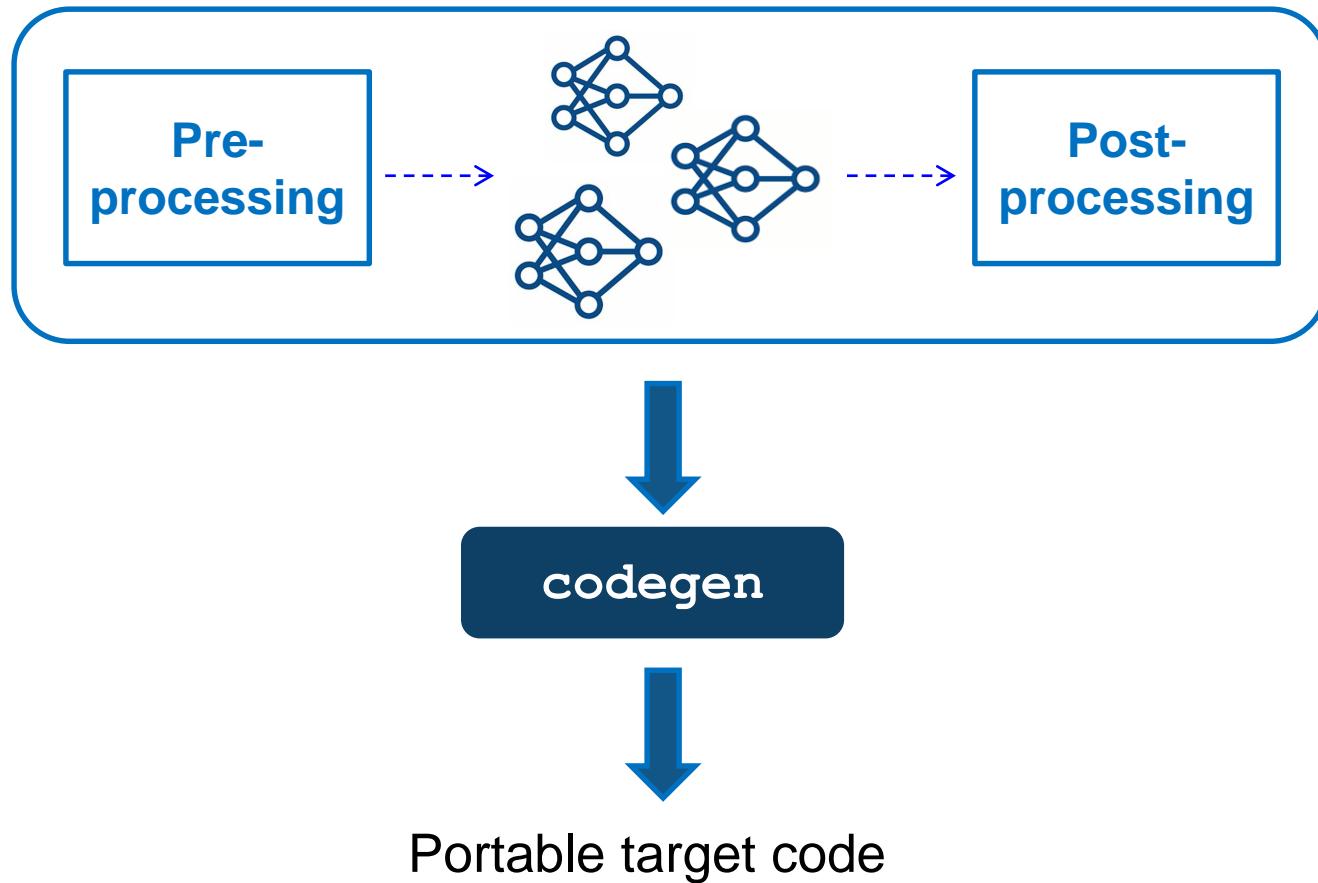


Intel® Math Kernel Library for Deep Neural Networks
(MKL-DNN)

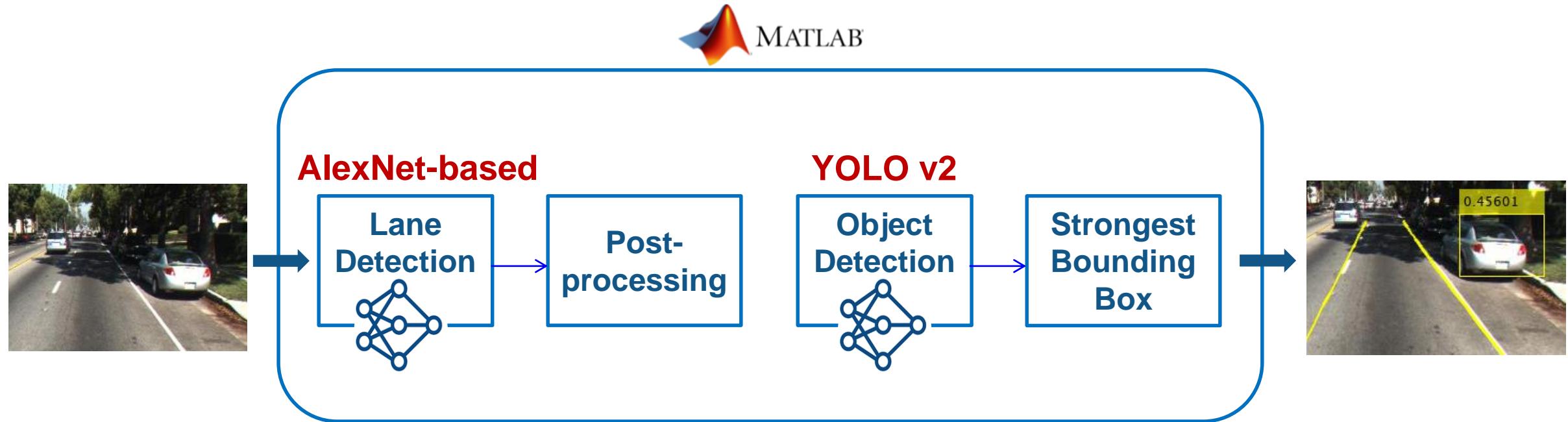


ARM® Compute libraries for mobile platforms

Deep Learning Workflows: Integrated Application Deployment



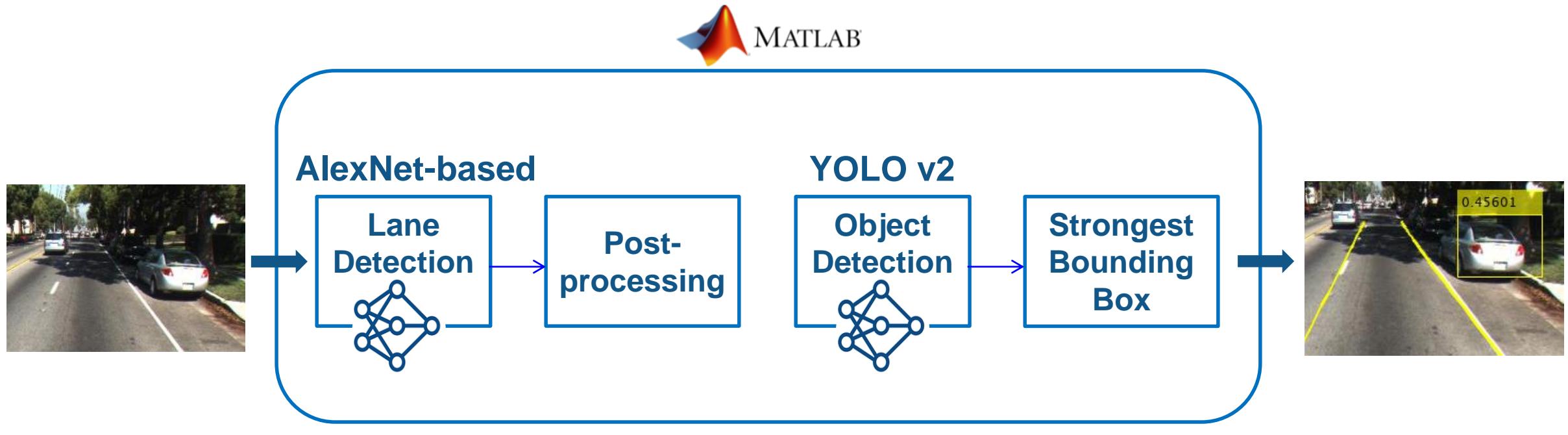
Lane and Object Detection using YOLO v2



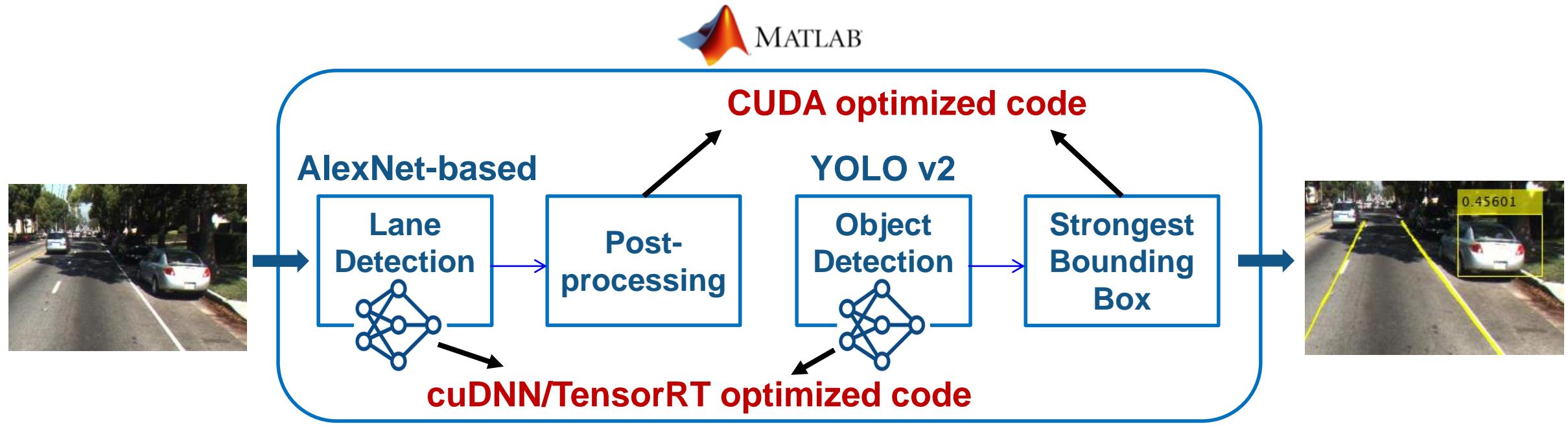
Workflow:

- 1) Test in MATLAB on CPU
- 2) Generate code and test on desktop GPU
- 3) Generate code and test on Jetson AGX Xavier GPU

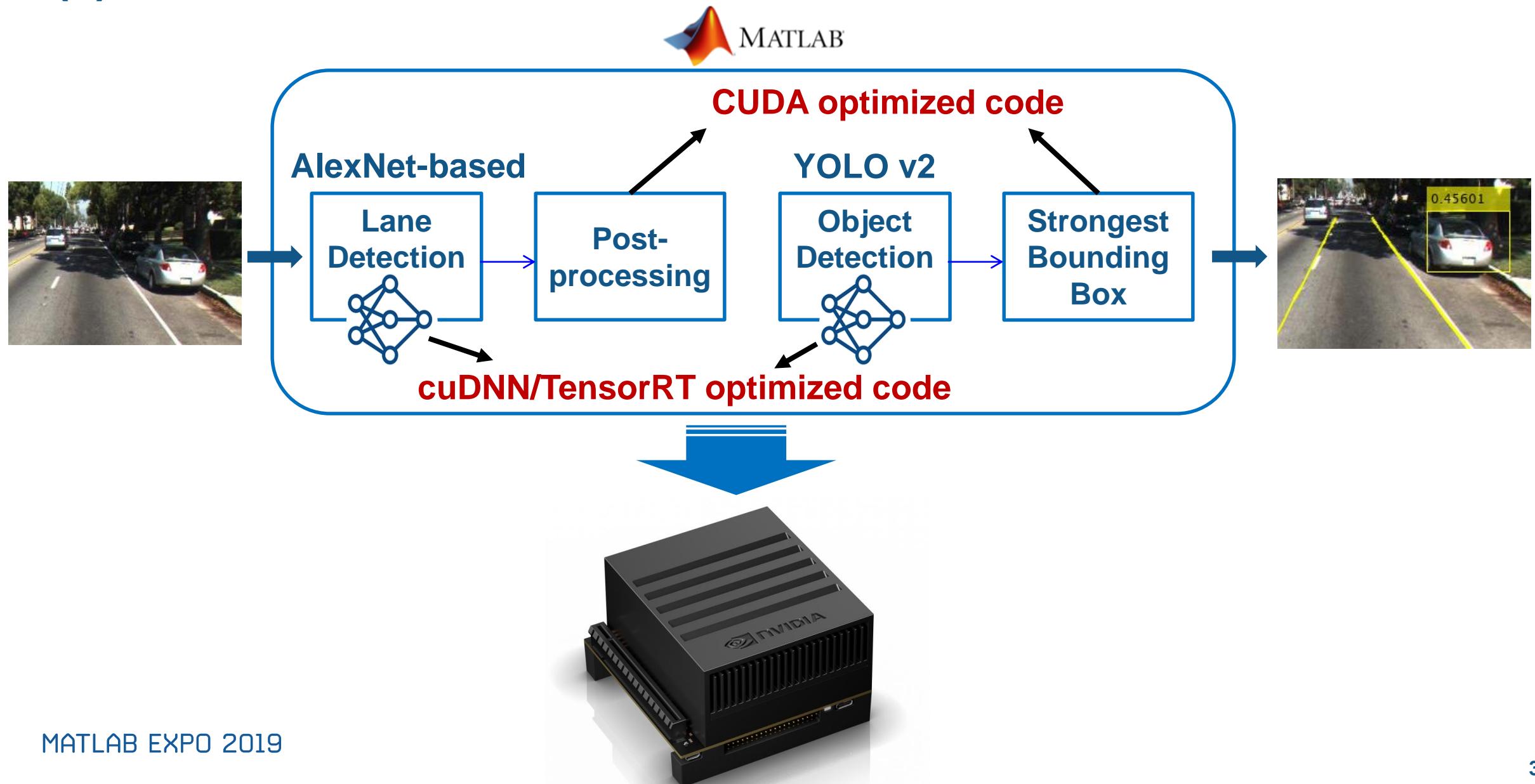
(1) Test in MATLAB on CPU



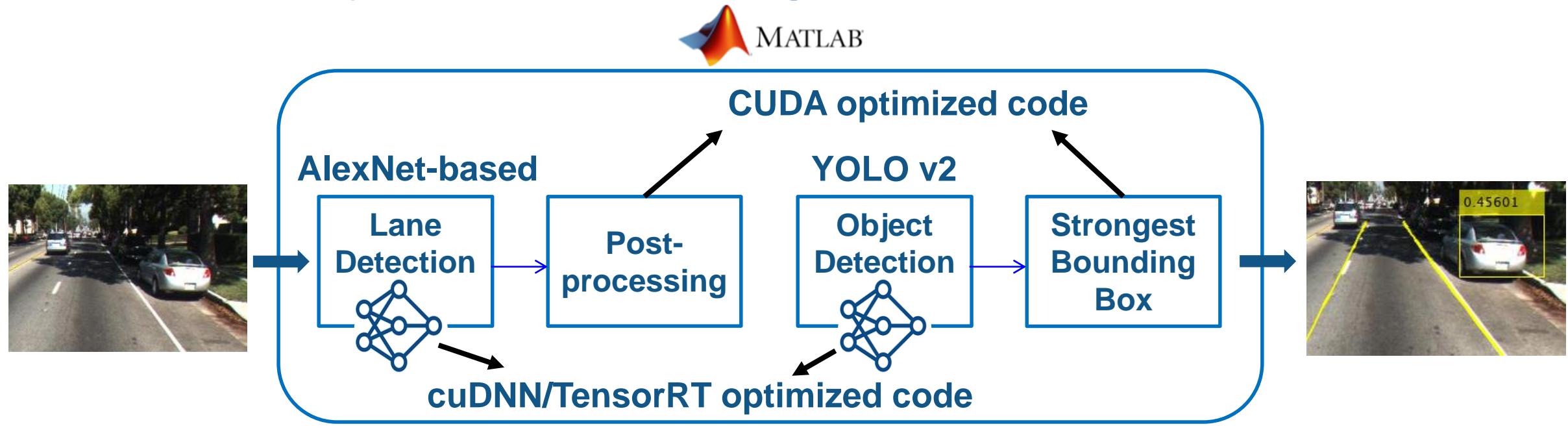
(2) Generate Code and Test on Desktop GPU



(3) Generate Code and Test on Jetson AGX Xavier GPU

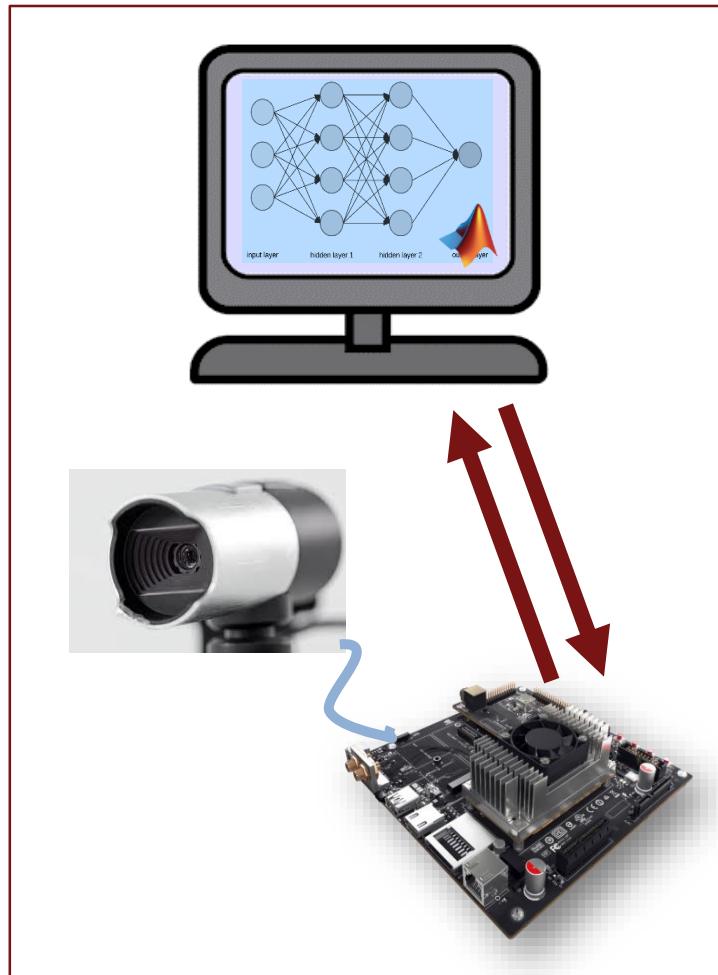


Lane and Object Detection using YOLO v2



- 1) Running on CPU
- 2) 7X faster running generate code on desktop GPU
- 3) Generate code and test on Jetson AGX Xavier GPU

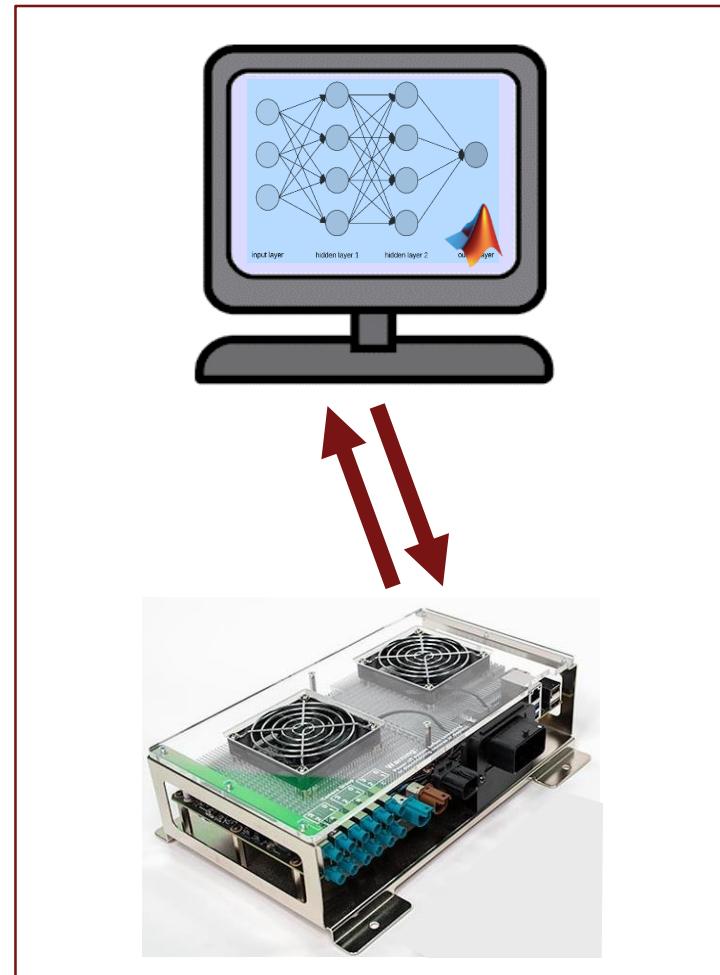
Accessing Hardware



Access Peripheral
from MATLAB

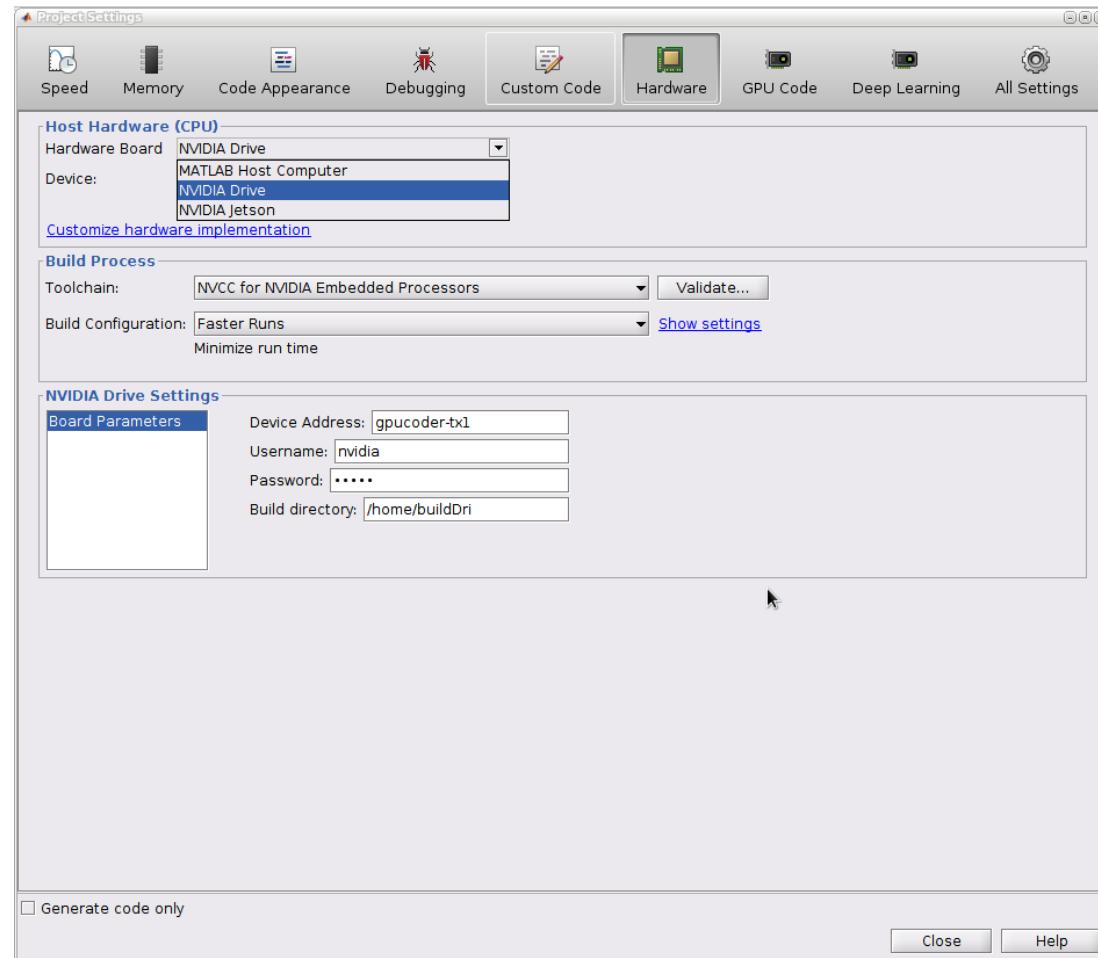


Deploy Standalone
Application



Processor-in-Loop
Verification

Deploy to Target Hardware via Apps and Command Line



```

%% Deploy and launch through NVIDIA HSP

%% setup hardware object
% create jetson/drive hardware object with IP or hostname of jeston/drive
% also pass credentials for login
hwObj = jetson('gpcoder-tx2-2','ubuntu','ubuntu');
hwObj.setupCodegenContext;

%% setup codegen config object
% create congen config and connect to hardware object.
cfg_hsp = coder.gpuConfig('exe');
cfg_hsp.Hardware = coder.hardware(hwObj.BoardPref);
buildDir ='~/buildDir';
cfg_hsp.Hardware.BuildDir = buildDir;

%% add user written main files for building executable
% and generate/build the code.
cfg_hsp.CustomSource = 'driver_files_alexnet/main.cu';
cfg_hsp.CustomInclude = 'driver_files_alexnet/';

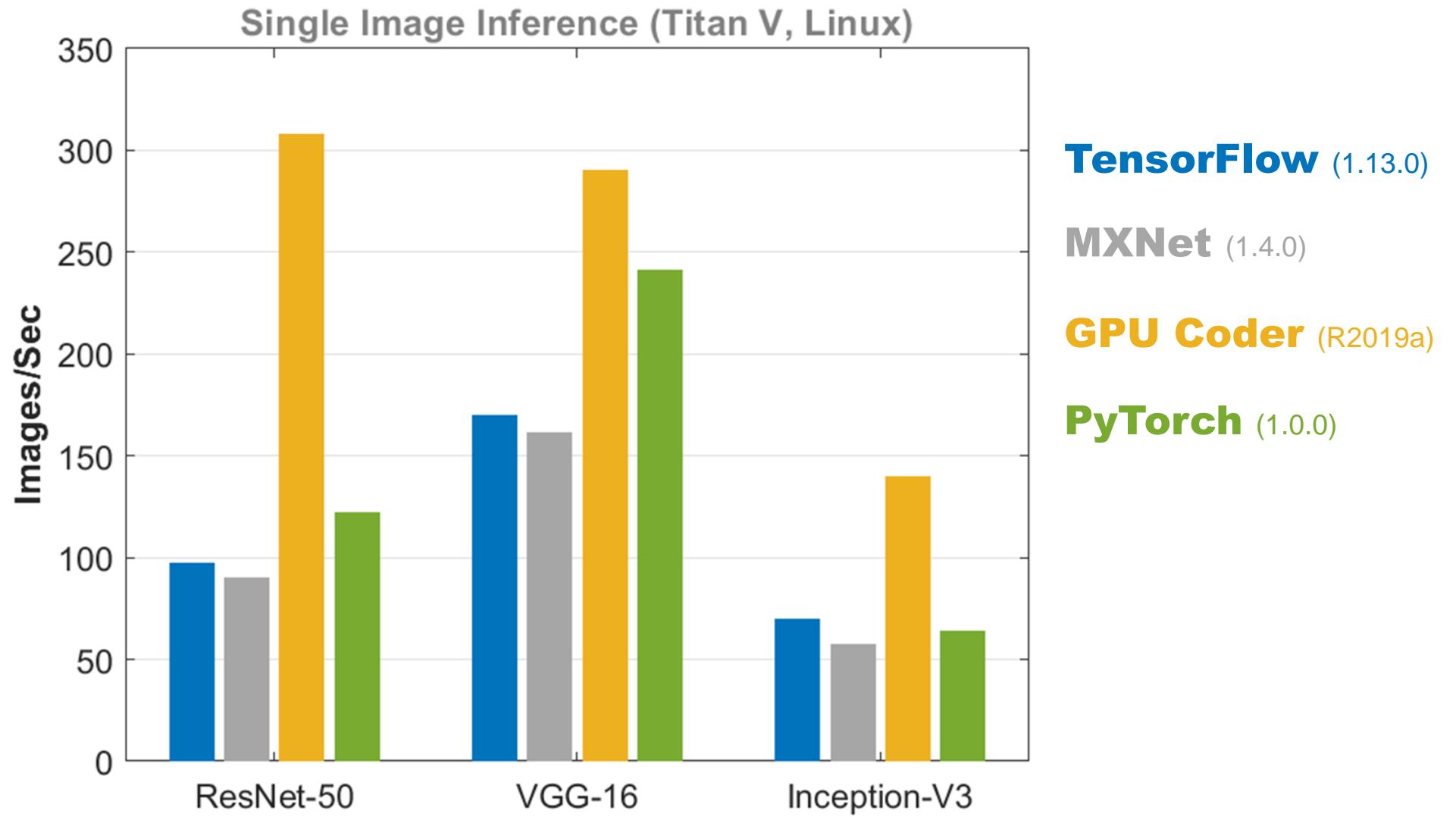
codegen -config cfg_hsp -args {im, coder.Constant(cnnMatFile)} alexnet_test

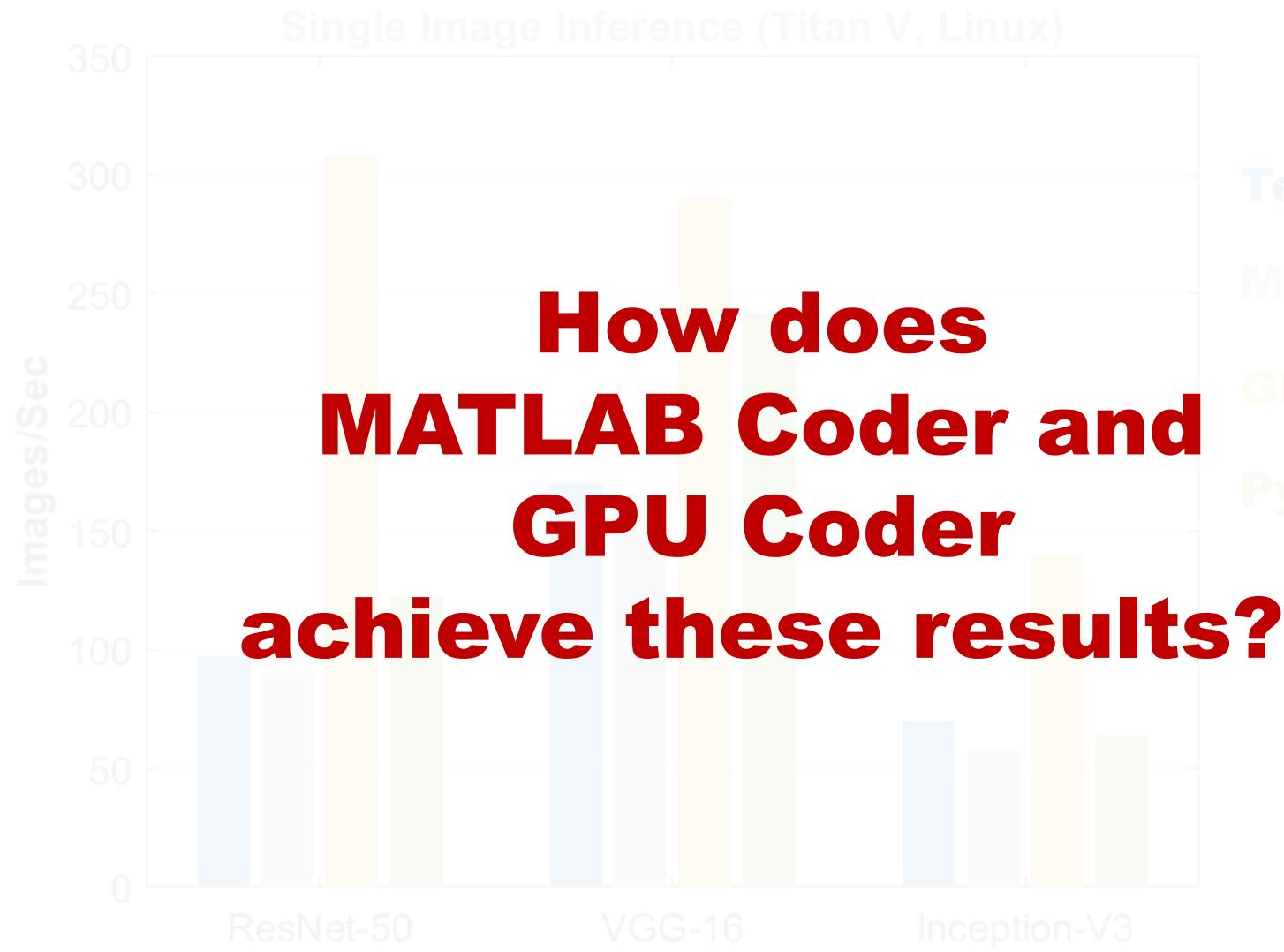
%% copy input and run the executable
hwObj.putFile('input2.txt', buildDir);
hwObj.putFile('synsetWords.txt', buildDir);

%execute on Jetson
hwObj.runExecutable([buildDir '/alexnet_test.elf'], 'input2.txt')

%% copy the output file back to host machine
hwObj.getFile([buildDir '/tOut.txt']);

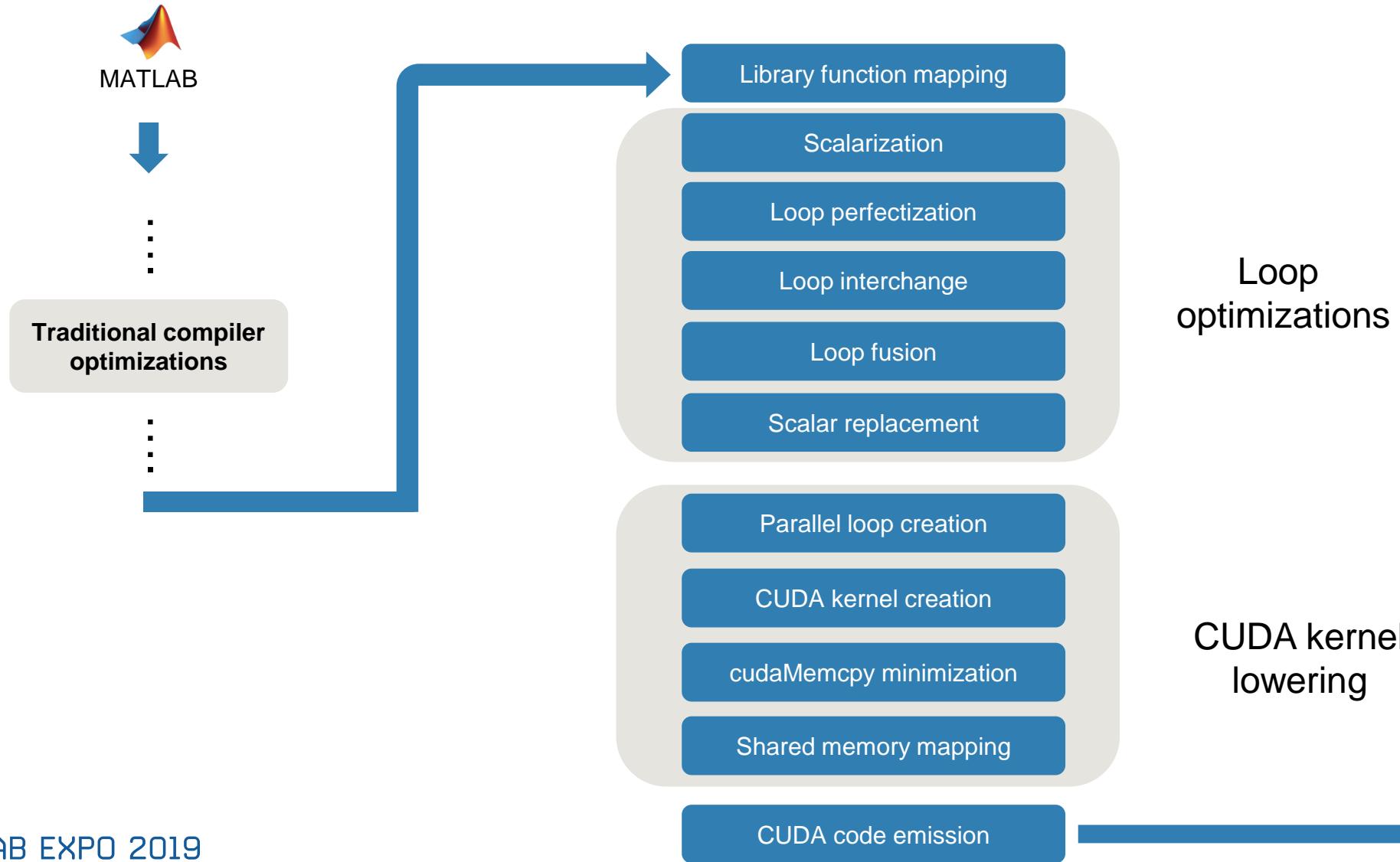
```





**How does
MATLAB Coder and
GPU Coder
achieve these results?**

Coders Apply Various Optimizations



Coders Apply Various Optimizations



Traditional compiler optimizations

Network Optimization

Optimized Libraries

Coding Patterns

String processing
Memory management
Loop unrolling
Loop exchange

Loop fusion

Scalar replacement

Parallel loop

CUDA kernel

cudaMemcpy minimization

Shared memory mapping

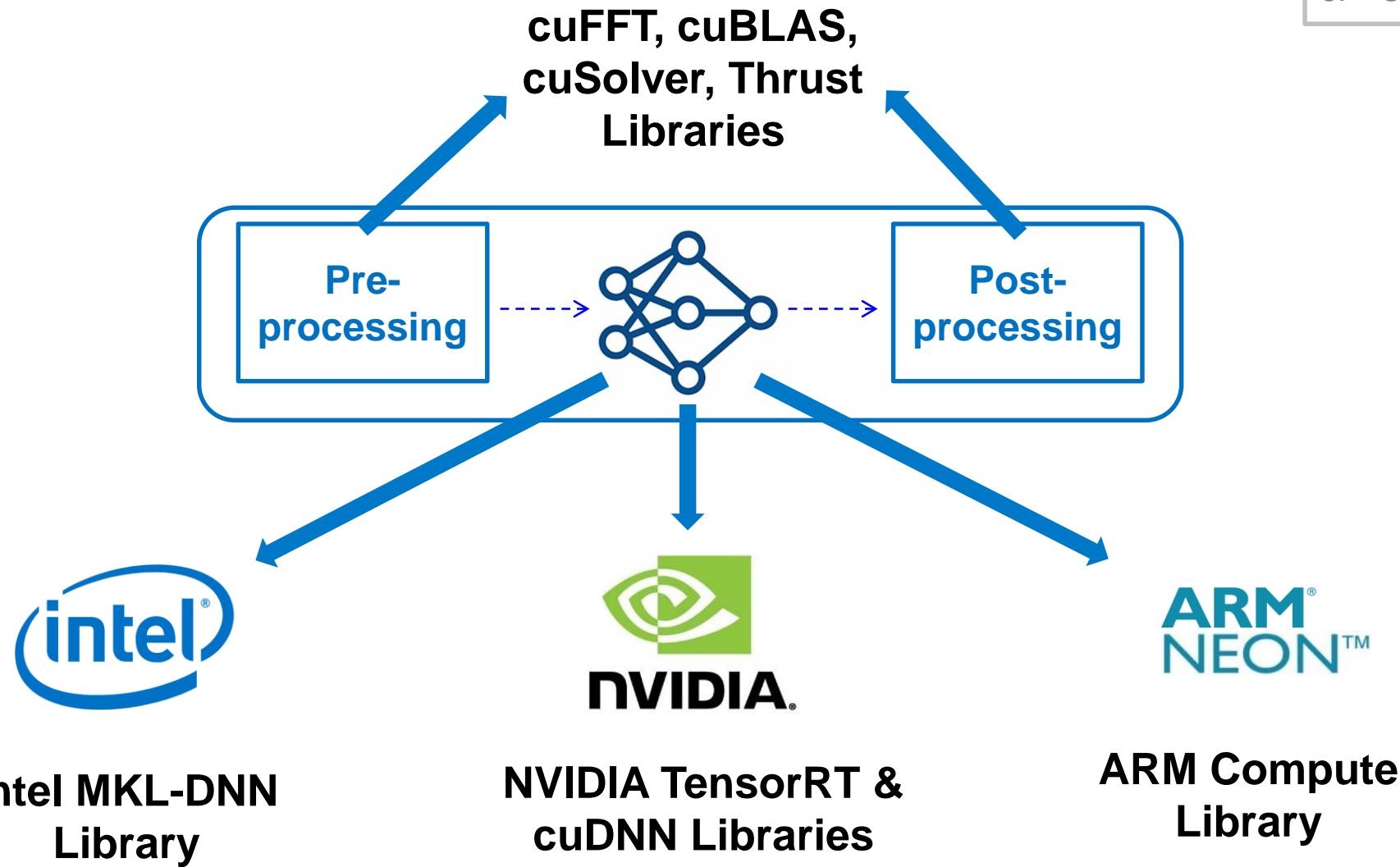
CUDA code emission

Loop optimizations

Kernel reordering

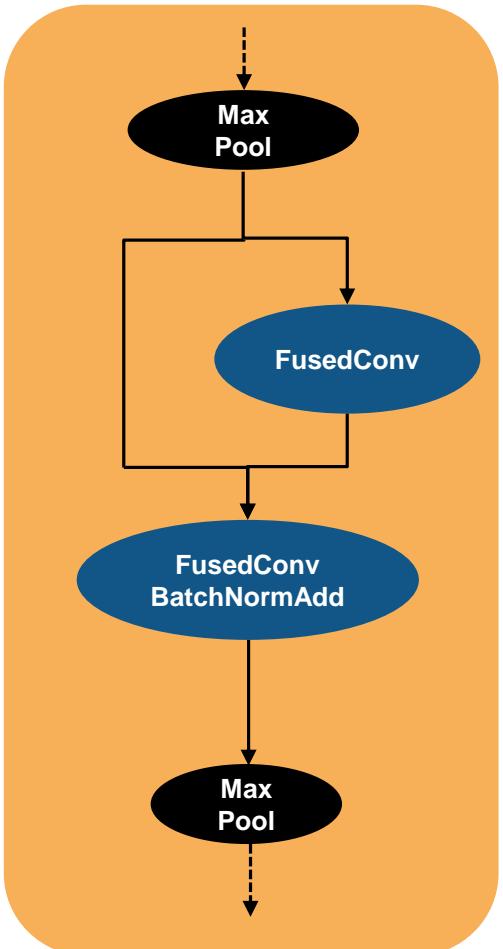
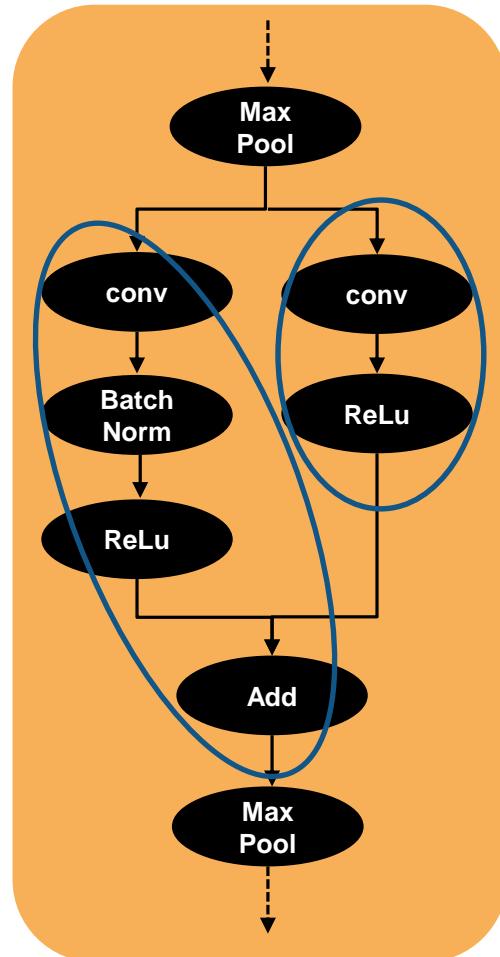
Generated Code Calls Optimized Libraries

- Performance
1. Optimized Libraries
 2. Network Optimizations
 3. Coding Patterns



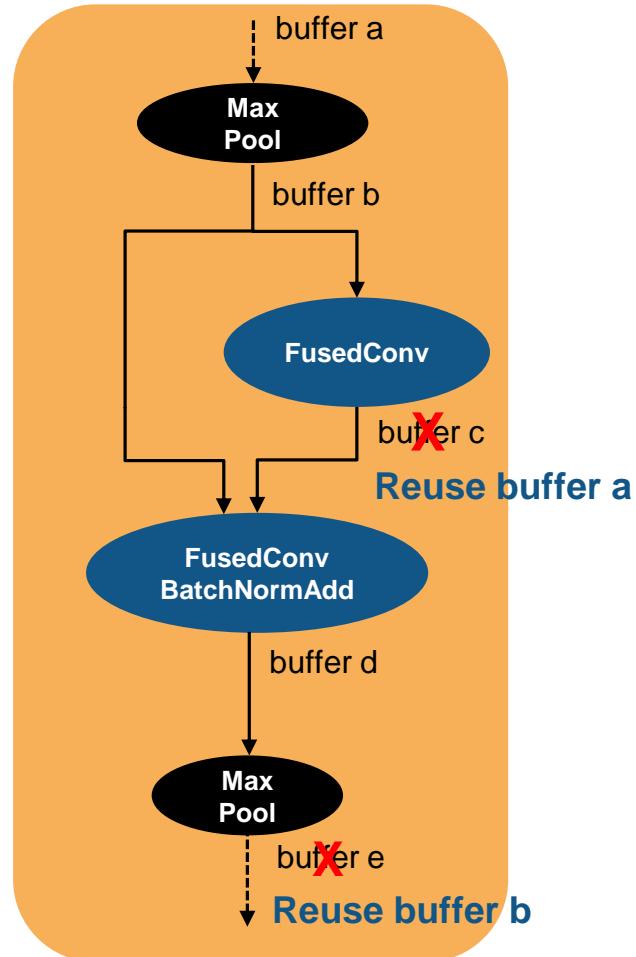
Deep Learning Network Optimization

- Performance
1. Optimized Libraries
 - 2. Network Optimizations**
 3. Coding Patterns



MATLAB EXPO 2025
Network

Layer fusion
Optimized computation

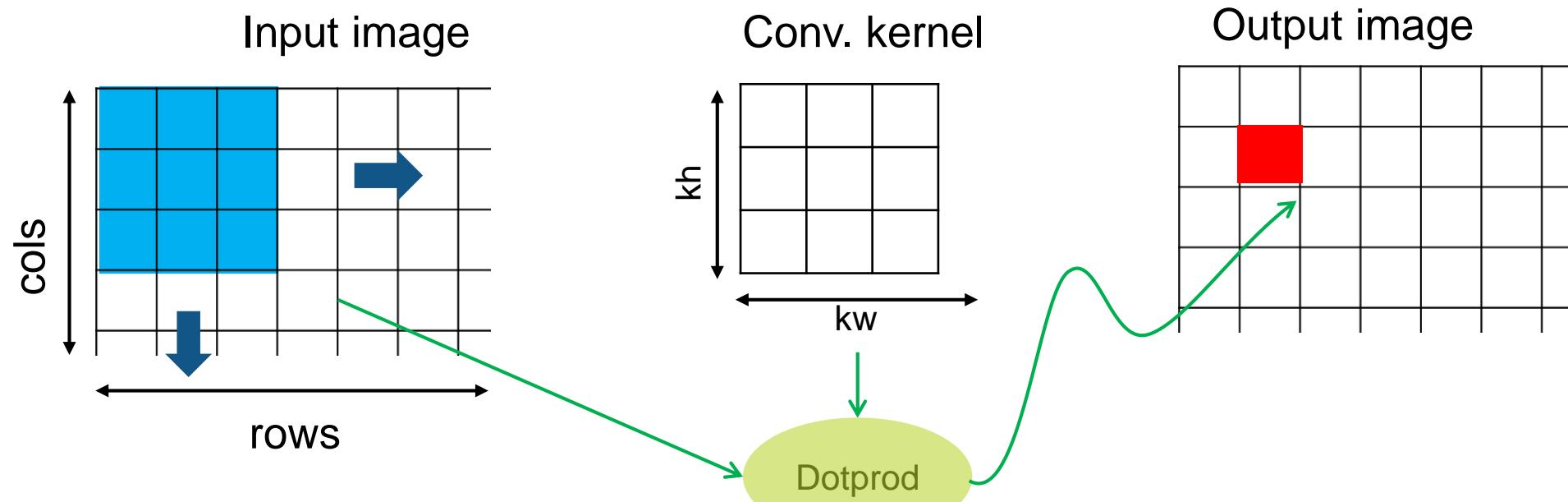


Buffer minimization
Optimized memory

Coding Patterns: Stencil Kernels

- Performance
1. Optimized Libraries
 2. Network Optimizations
 3. **Coding Patterns**

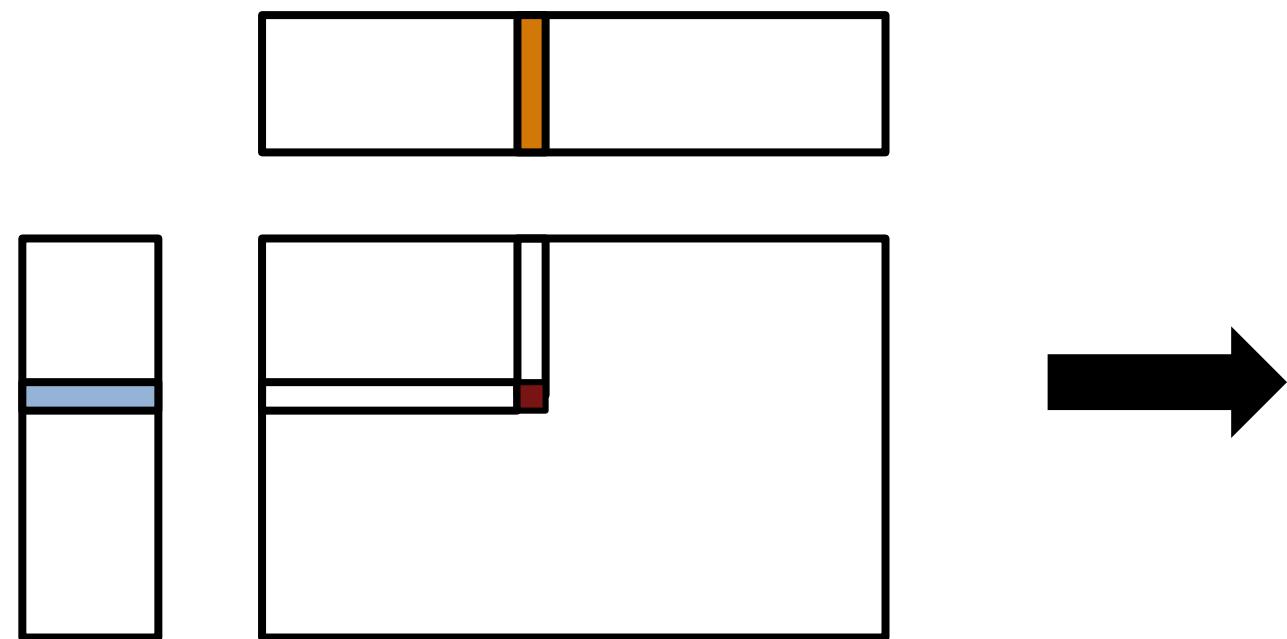
- Automatically applied for image processing functions (e.g. `imfilter`, `imerode`, `imdilate`, `conv2`, ...)
- Manually apply using `gpuCoder.stencilKernel()`



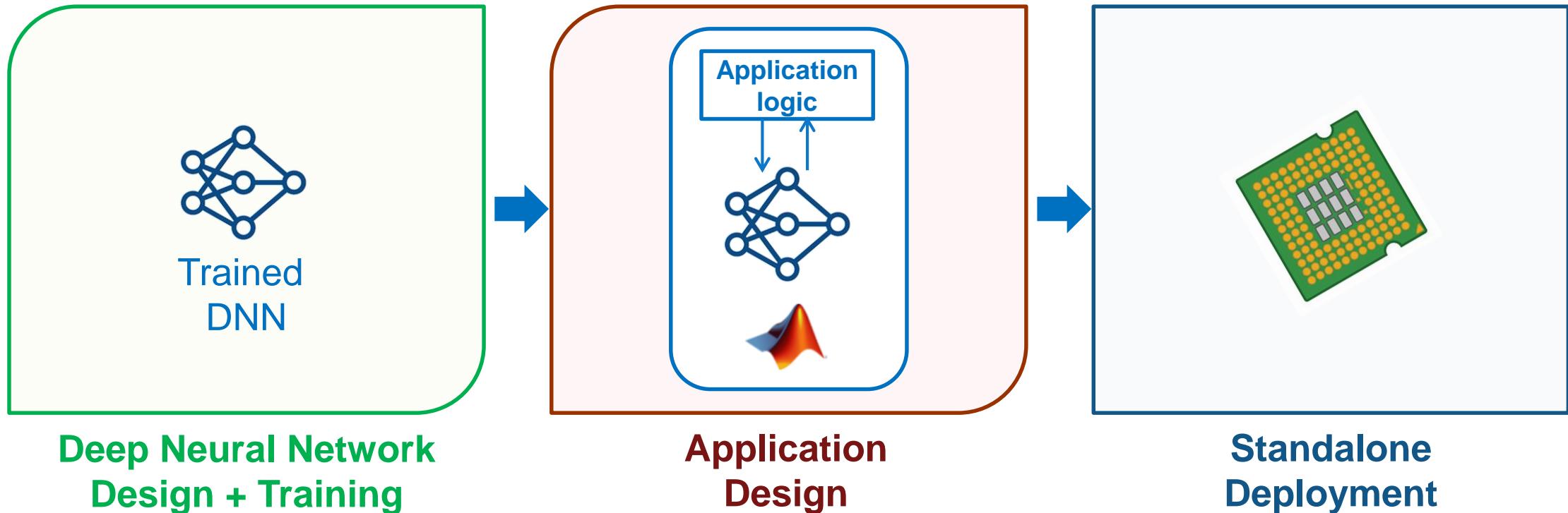
Coding Patterns: Matrix-Matrix Kernels

- Performance
1. Optimized Libraries
 2. Network Optimizations
 3. **Coding Patterns**

- Automatically applied for many MATLAB functions (e.g. matchFeatures SAD, SSD, pdist, ...)
- Manually apply using *gpuCoder.matrixMatrixKernel()*



Deep Learning Workflow in MATLAB



Deep Learning Workflow in MATLAB

