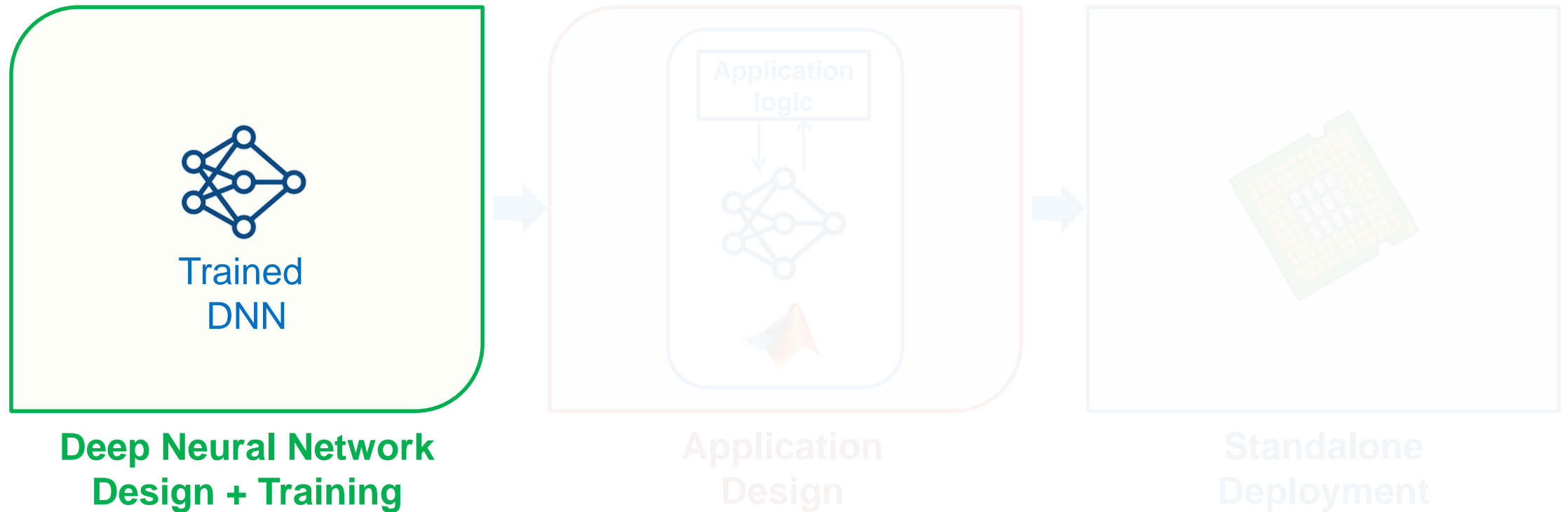# MATLAB EXPO 2019

# Deploying Deep Neural Networks to Embedded GPUs and CPUs
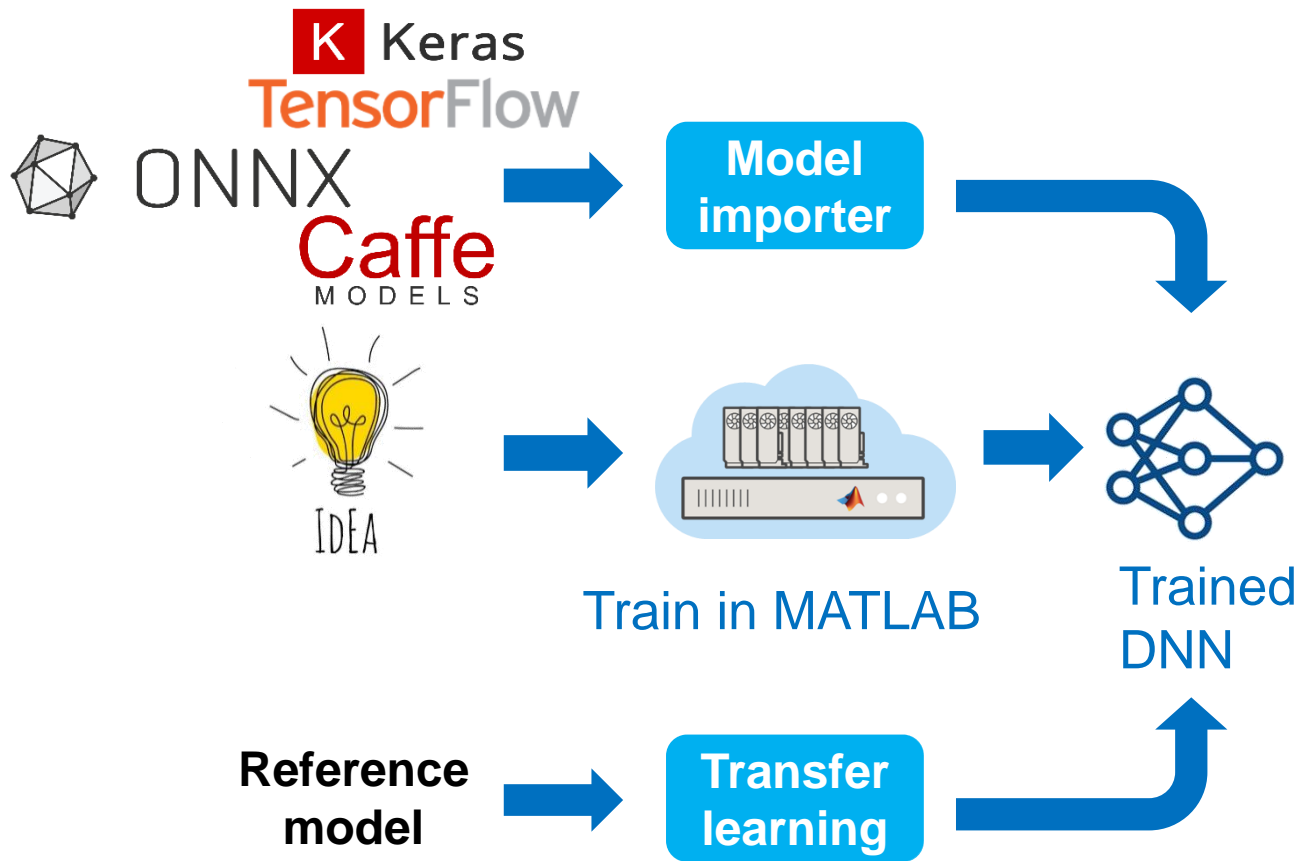
*Dr Rishu Gupta*
*Senior Application Engineer*

# Deep Learning Workflow in MATLAB



**Deep Neural Network
Design + Training**

Application
Design

Standalone
Deployment

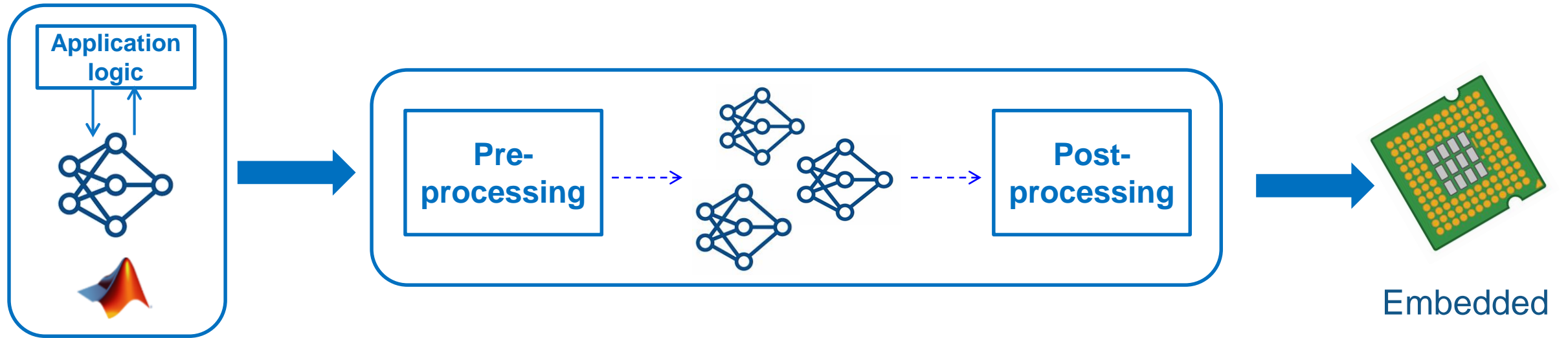# Deep Neural Network Design and Training



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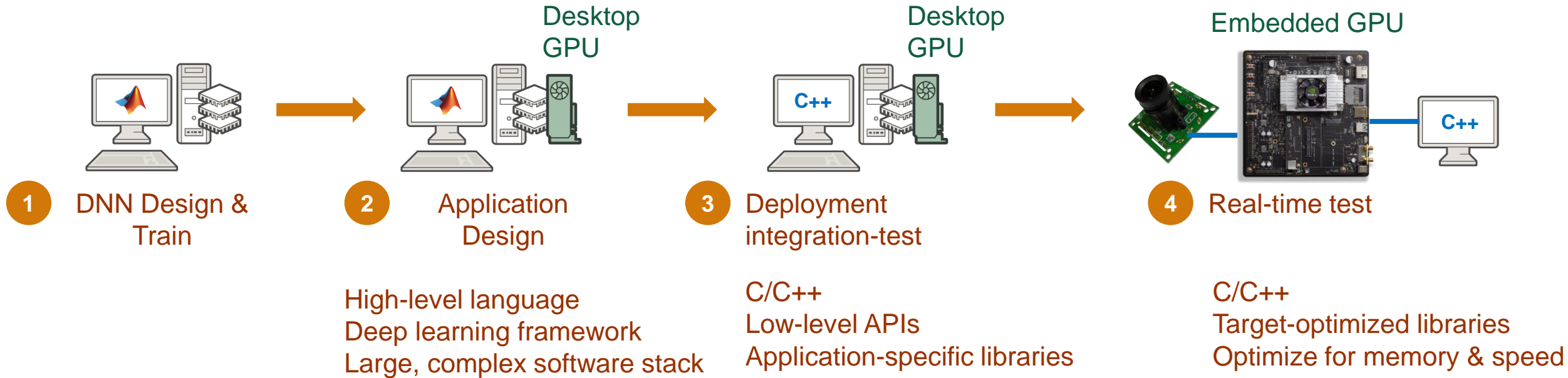
- **Training in MATLAB**
  - **Acceleration** with GPU's
  - **Scale** to clusters

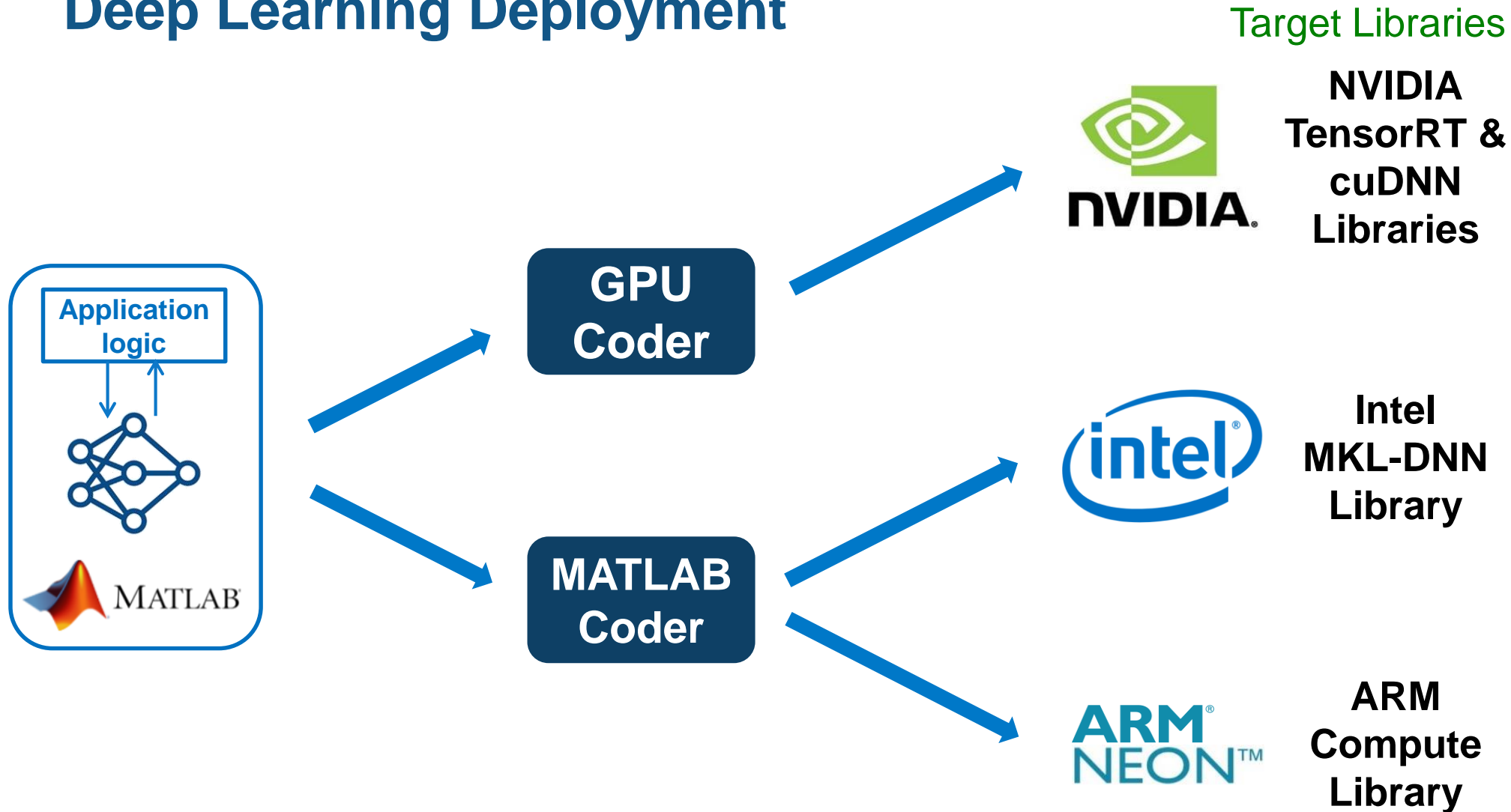# Application Design



**Multi-Platform Deep Learning Deployment**
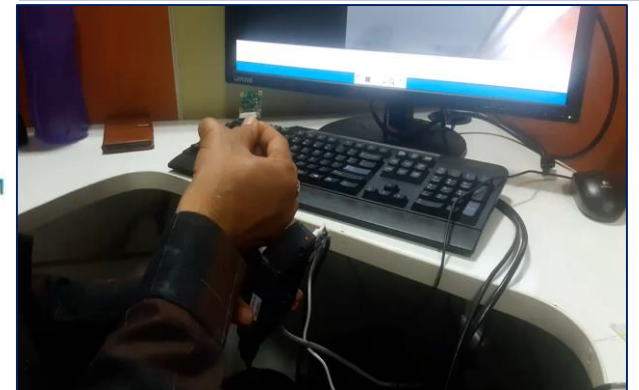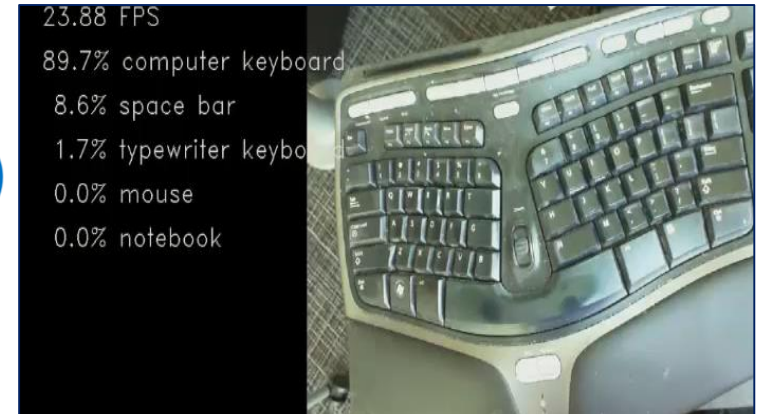
# Algorithm Design to Embedded Deployment Workflow

Desktop GPU

Desktop GPU

Embedded GPU

**1** DNN Design & Train

**2** Application Design

**3** Deployment integration-test

**4** Real-time test

High-level language
Deep learning framework
Large, complex software stack

C/C++
Low-level APIs
Application-specific libraries

C/C++
Target-optimized libraries
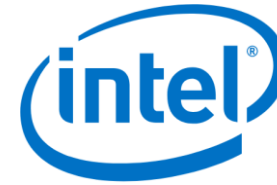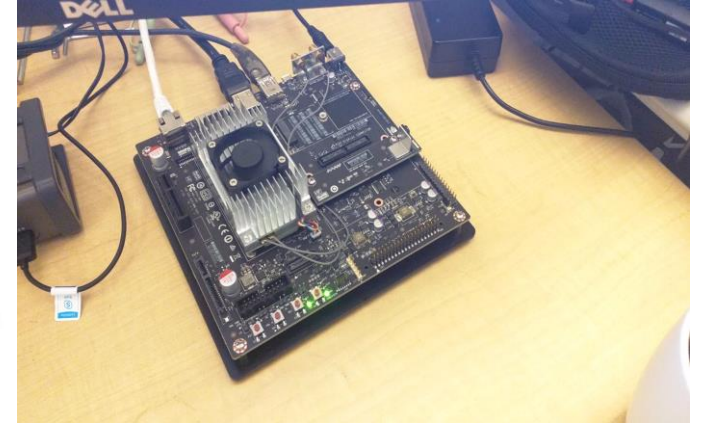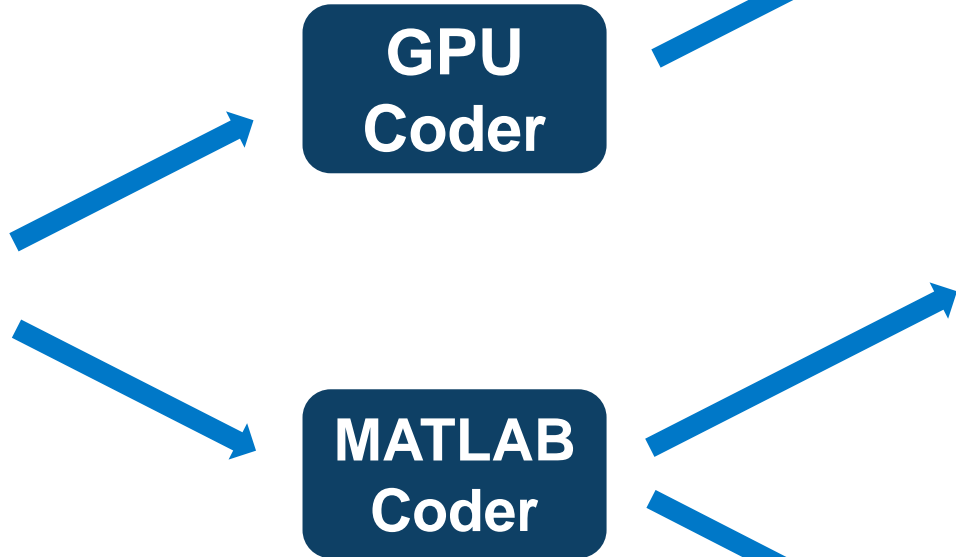Optimize for memory & speed

## Challenges
- Integrating multiple libraries and packages
- Verifying and maintaining multiple implementations
- Algorithm & vendor lock-in

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

Target Libraries

**Application logic**

MATLAB

**GPU Coder**

**MATLAB Coder**

**NVIDIA TensorRT & cuDNN Libraries**

**Intel MKL-DNN Library**

**ARM Compute Library**

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

# Musashi Seimitsu Industry Co.,Ltd.
## Detect Abnormalities in Automotive Parts



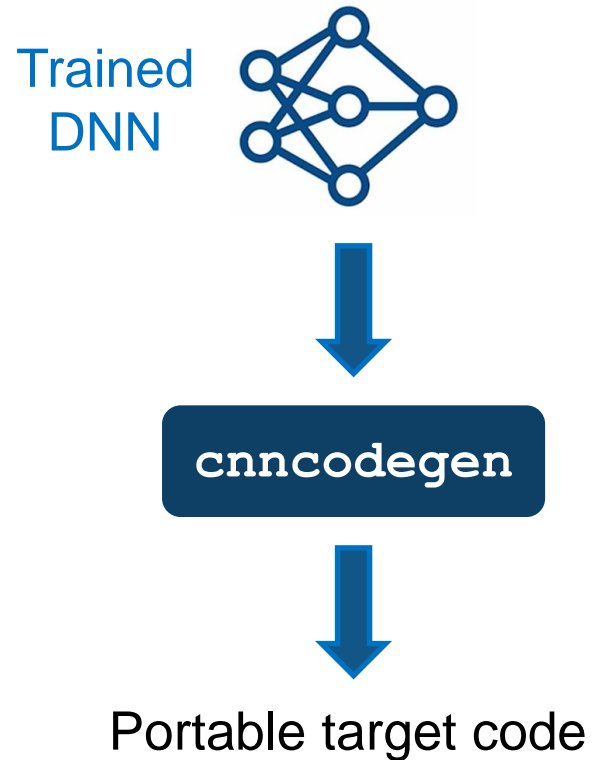Automated visual inspection of 1.3 million bevel gear per month

**MATLAB use in project:**

- Preprocessing of captured images

- Image annotation for training

- Deep learning based analysis
  - Various transfer learning methods (Combinations of CNN models, Classifiers)
  - Estimation of defect area using Class Activation Map (CAM)
  - Abnormality/defect classification

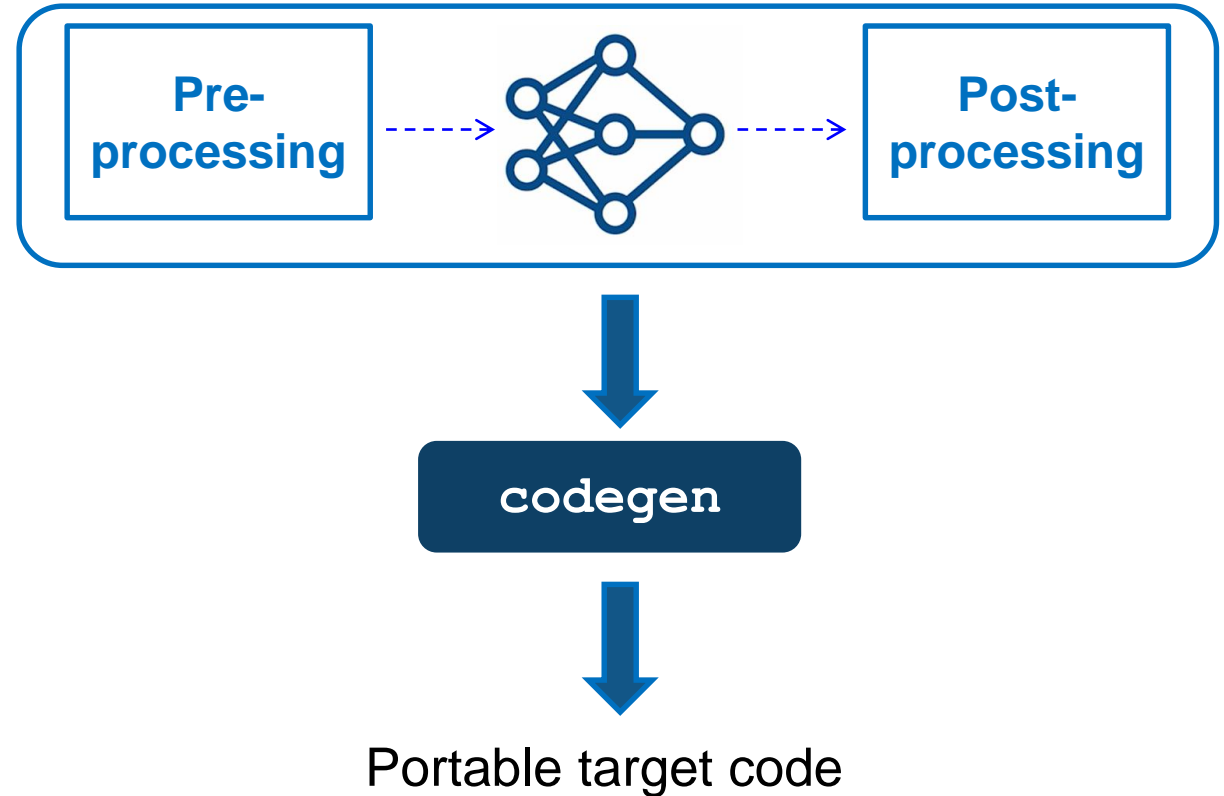- Deployment to NVIDIA Jetson using GPU Coder

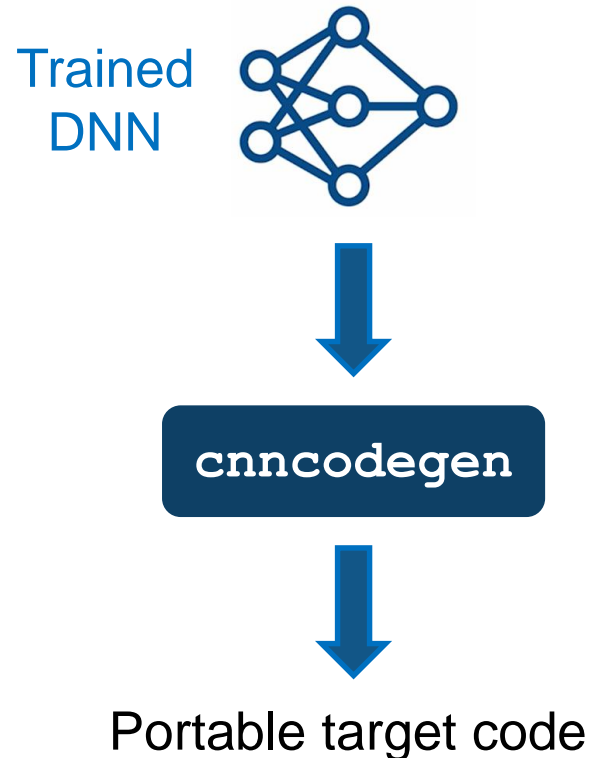# Deep Learning Deployment Workflows

**INFERENCE ENGINE DEPLOYMENT**

**INTEGRATED APPLICATION DEPLOYMENT**

Trained
DNN

Pre-
processing

Post-
processing

**cnncodegen**

**codegen**

Portable target code

Portable target code

# Workflow for Inference Engine Deployment

**INFERENCE ENGINE DEPLOYMENT**

Trained
DNN

⬇

**cnncodegen**

⬇

Portable target code

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
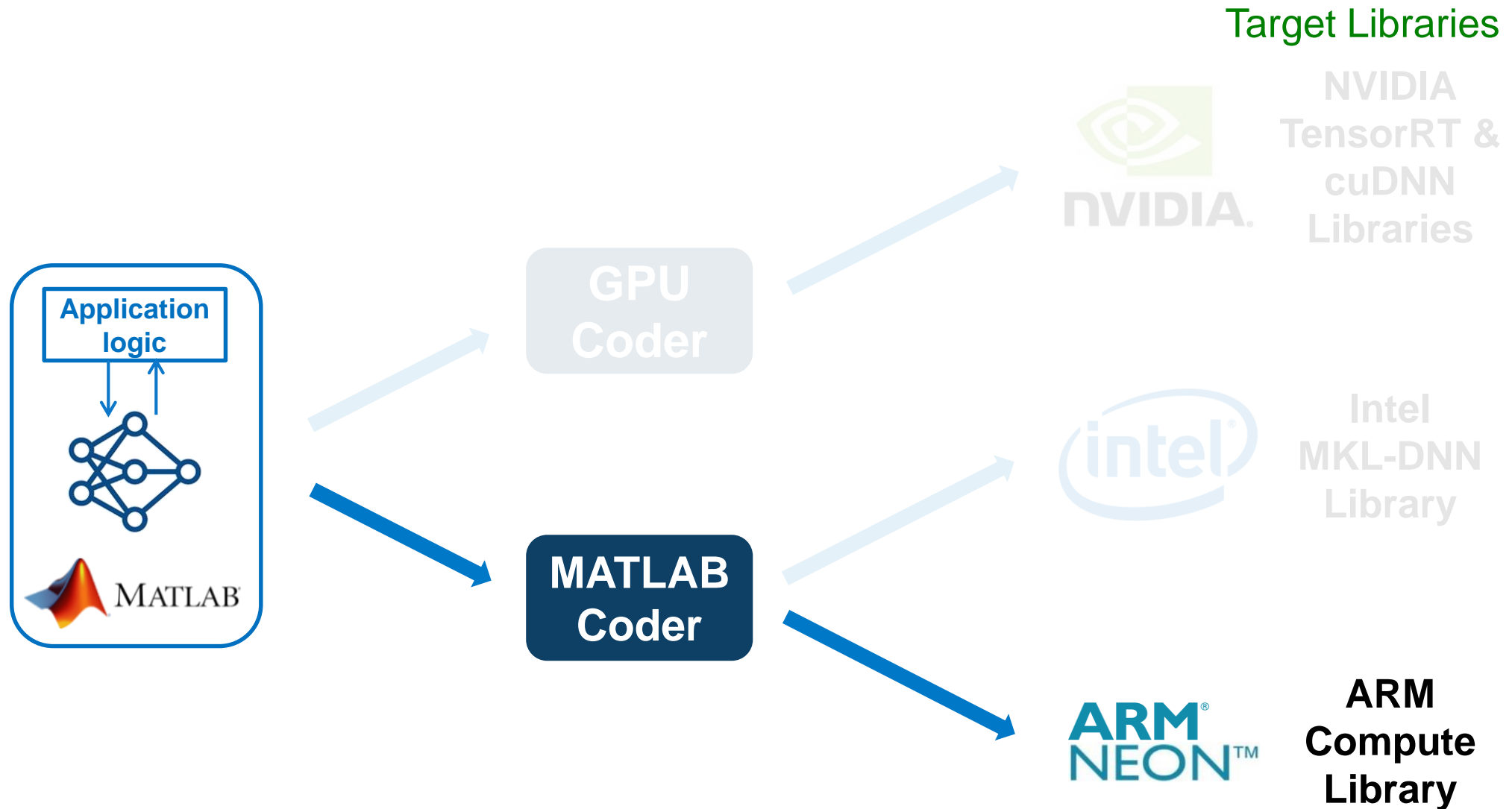
3. Build the code for the inference engine
```
>> make –C ./codegen –f …mk
```

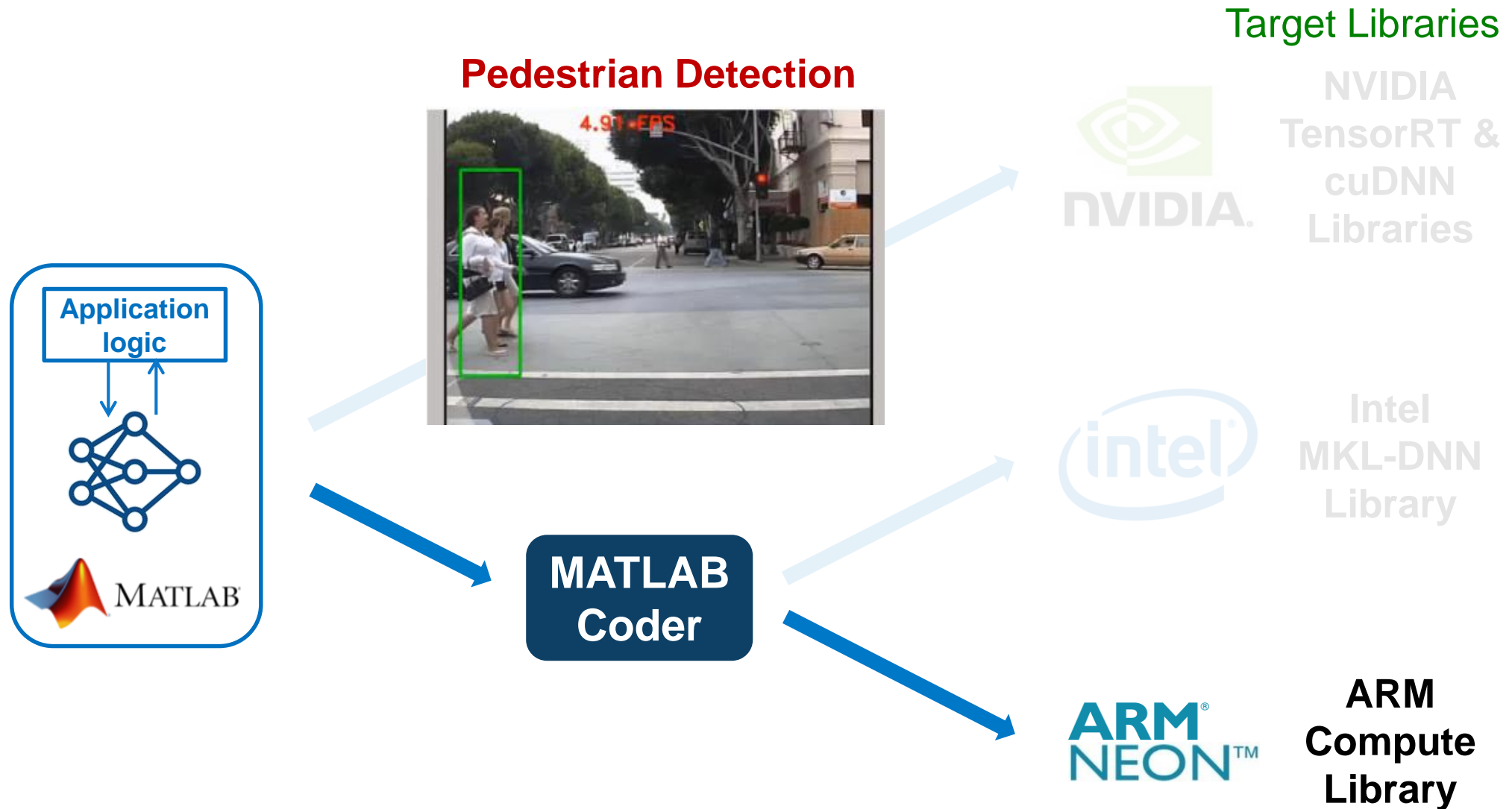4. Use hand written main function to call inference engine

5. Generate the exe and test the executable
```
>> make –C ./ ……
```

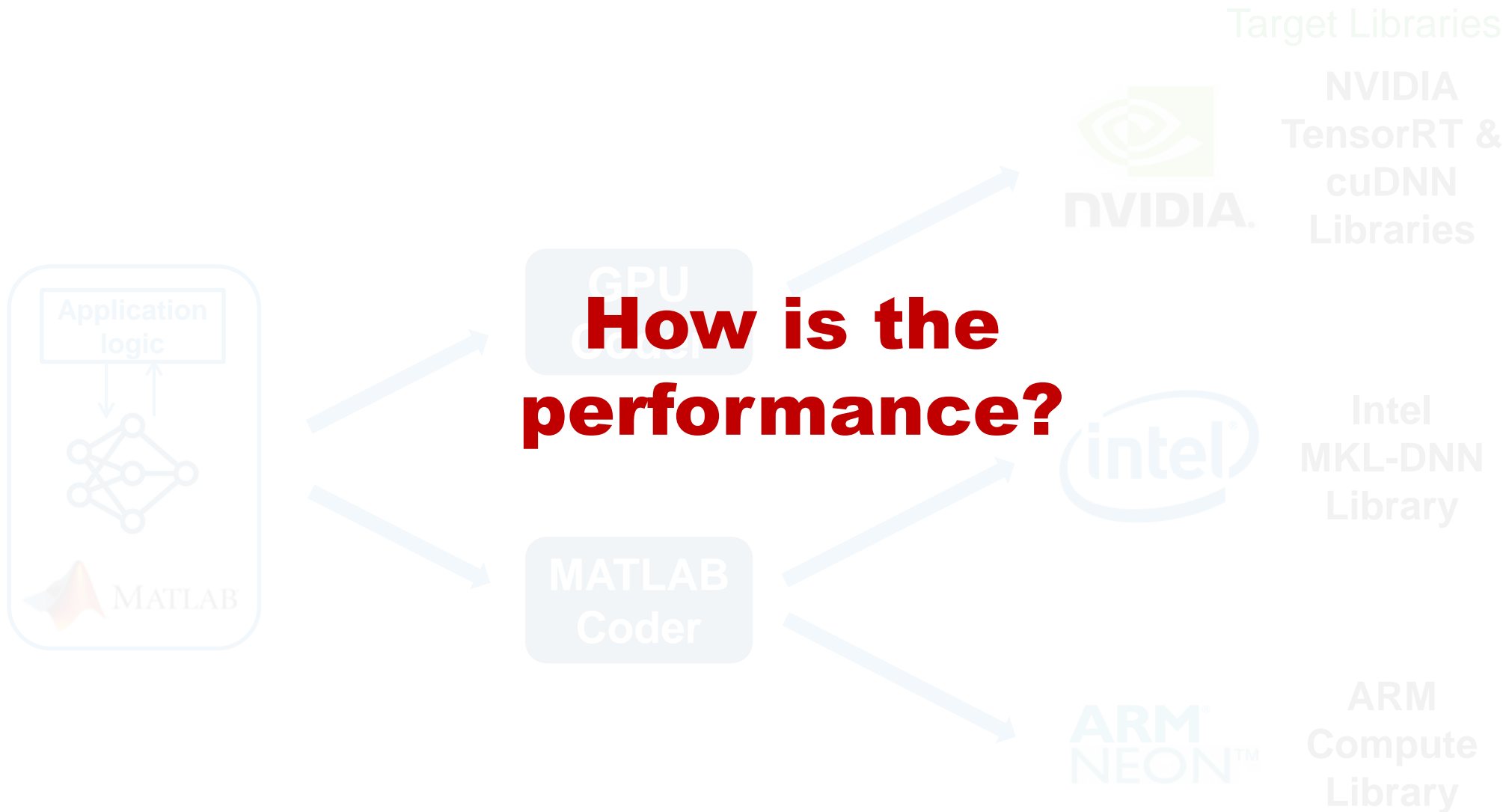# Deep Learning Inference Deployment

**Target Libraries**

**Pedestrian Detection**



**Application logic**

MATLAB

**MATLAB Coder**

NVIDIA TensorRT & cuDNN Libraries

Intel MKL-DNN Library

ARM NEON™ ARM Compute Library
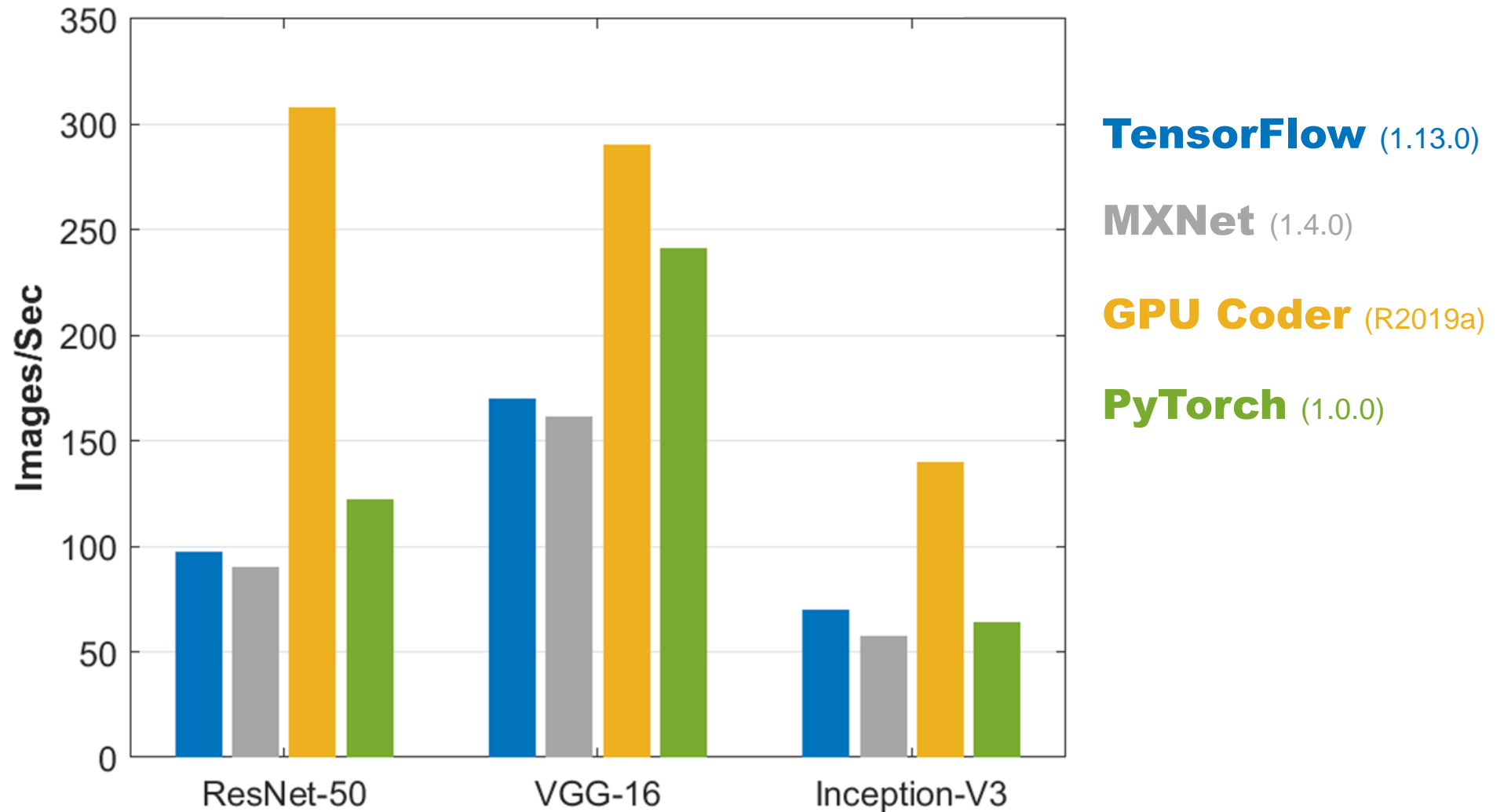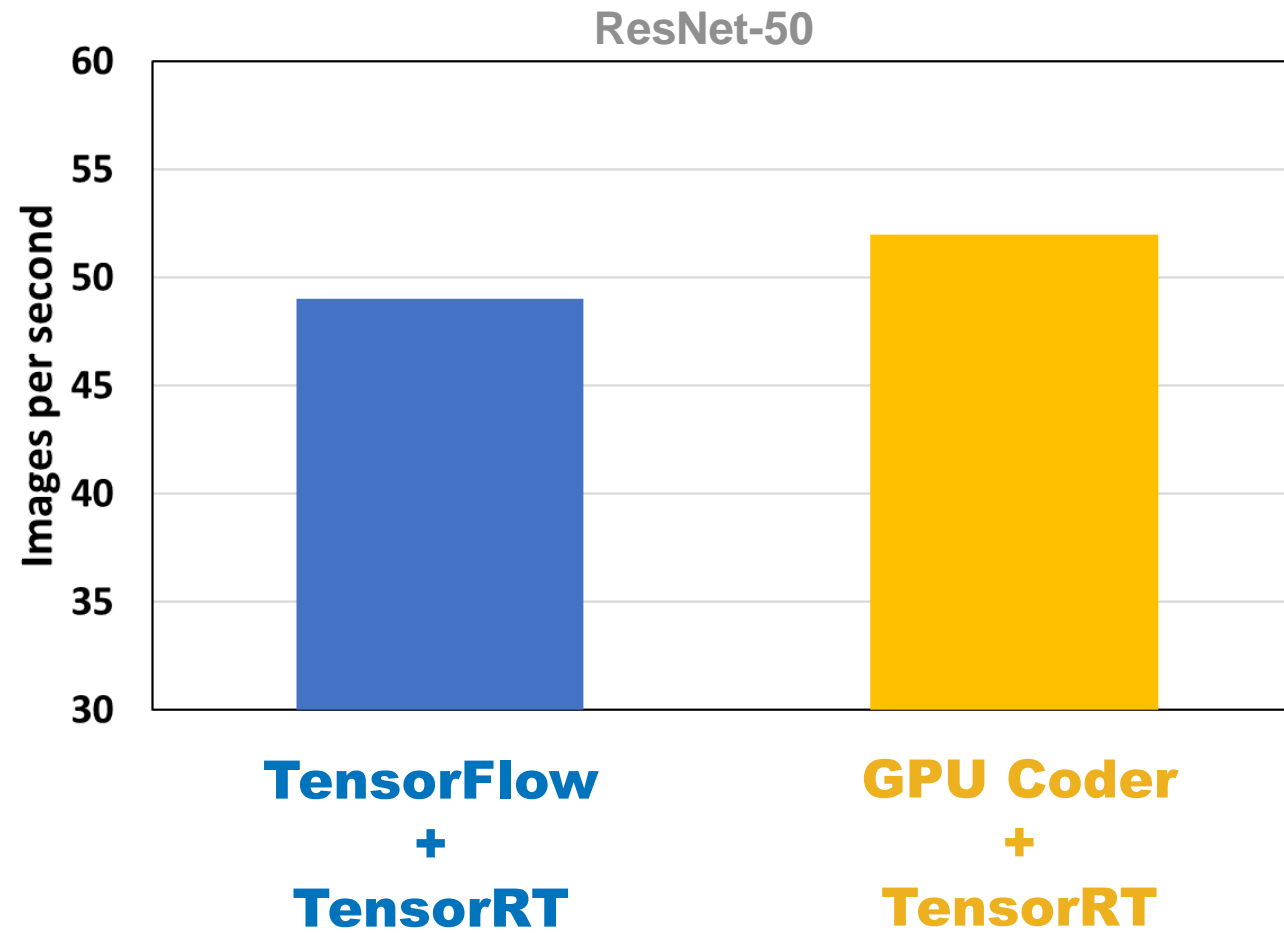
How is the performance?

# 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

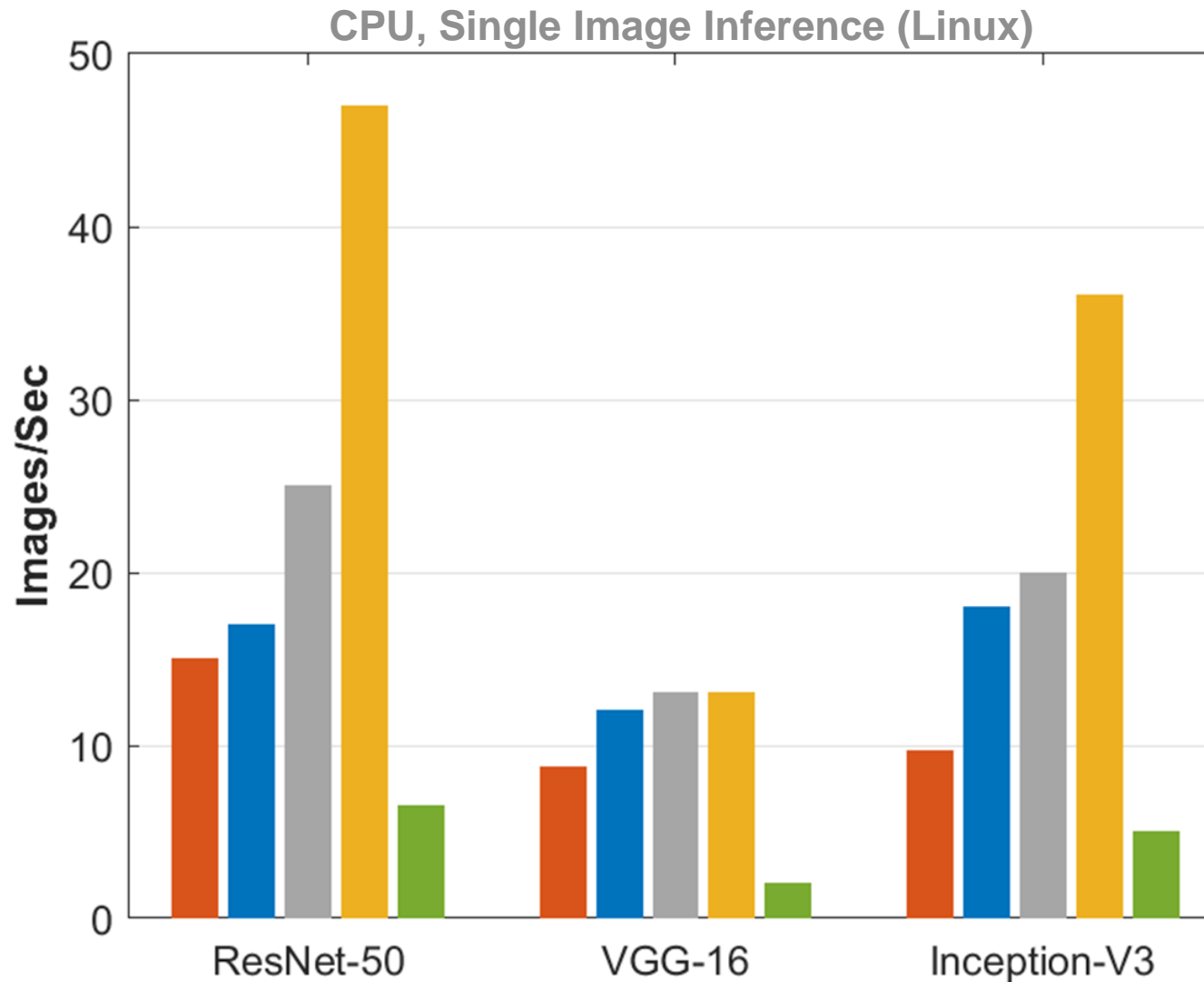# Single Image Inference on Titan V using cuDNN



**TensorFlow** (1.13.0)

**MXNet** (1.4.0)

**GPU Coder** (R2019a)

**PyTorch** (1.0.0)

*Intel® Xeon® CPU 3.6 GHz - NVIDIA libraries: CUDA10 - cuDNN 7 - Frameworks: TensorFlow 1.13.0, MXNet 1.4.0 PyTorch 1.0.0*

# CPU Performance

**CPU, Single Image Inference (Linux)**



**MATLAB**
**TensorFlow**
**MXNet**
**MATLAB Coder**
**PyTorch**

*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:

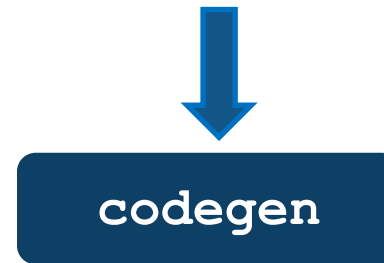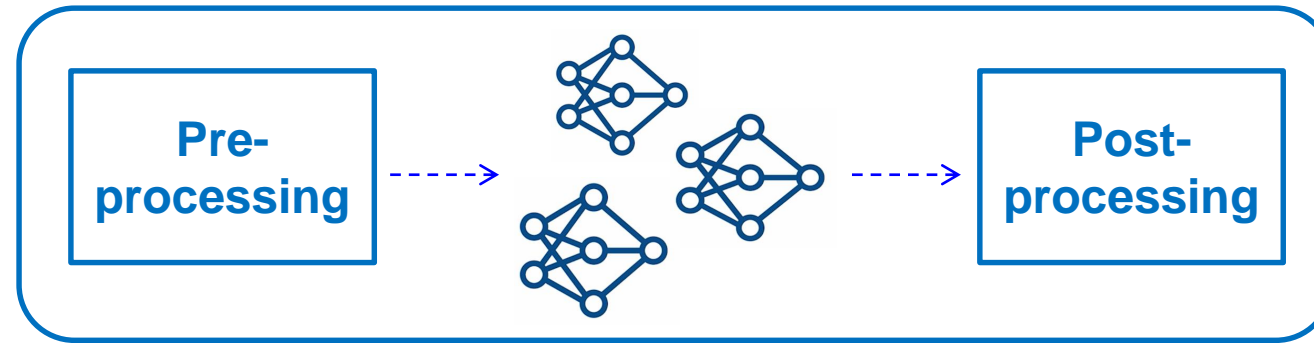**But, applications require more than just inference**

NVIDIA® CUDA libraries, including TensorRT & cuDNN

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

ARM NEON™ ARM® Compute libraries for mobile platforms

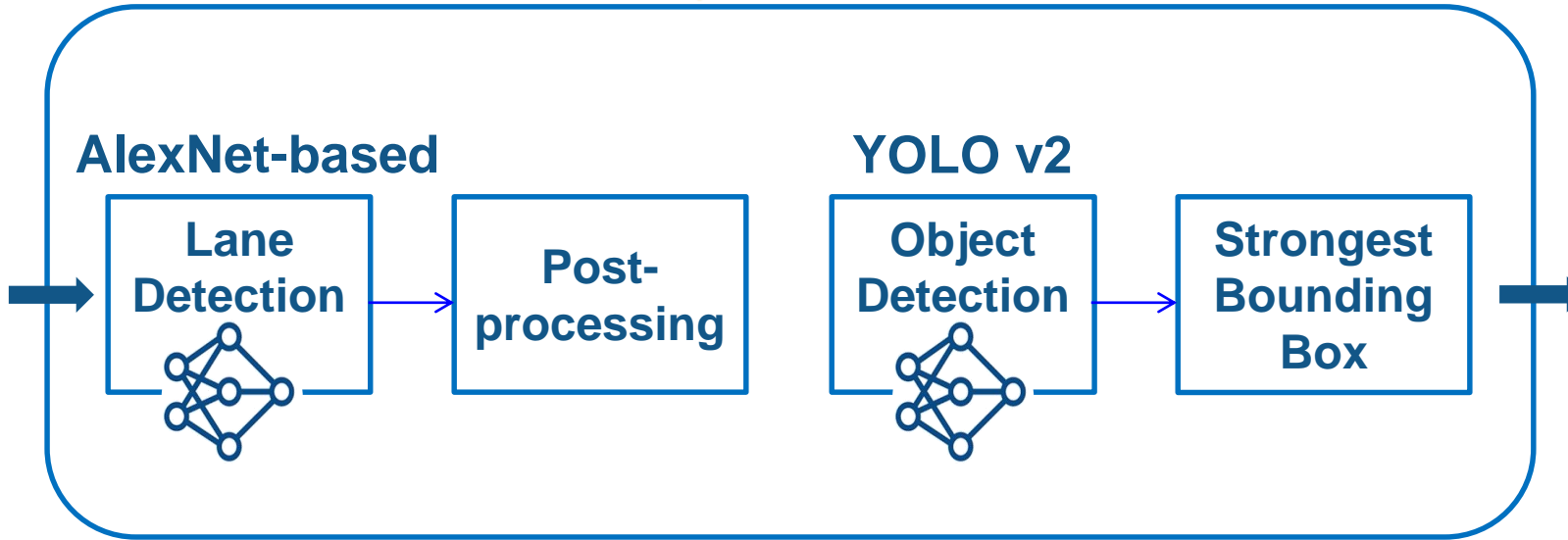# Deep Learning Workflows: Integrated Application Deployment



Pre-processing

Post-processing

codegen

Portable target code

# Lane and Object Detection using YOLO v2



**AlexNet-based**
**YOLO v2**

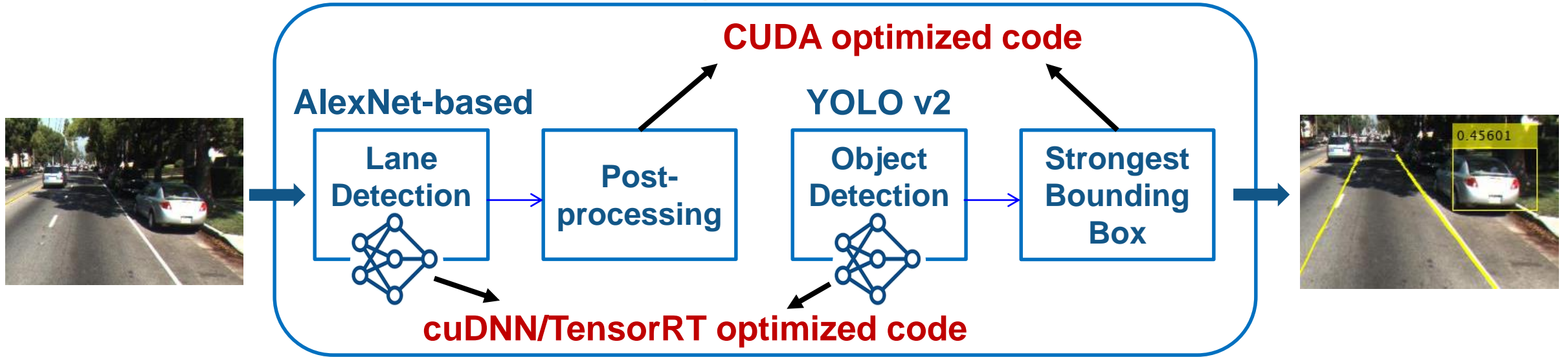Lane Detection → Post-processing → Object Detection → Strongest Bounding Box

Workflow:

1) Test in MATLAB

2) Generate code and test on desktop

3) Generate code and test on Jetson AGX Xavier GPU

# (1) Test in MATLAB

# (2) Generate Code and Test on Desktop GPU

MATLAB

**CUDA optimized code**

**AlexNet-based**

**YOLO v2**

Lane Detection → Post-processing → Object Detection → Strongest Bounding Box

**cuDNN/TensorRT optimized code**

0.43404

0.45601

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

MATLAB

**CUDA optimized code**

**AlexNet-based**

**YOLO v2**

| **Lane Detection** | **Post-processing** | **Object Detection** | **Strongest Bounding Box** |

**cuDNN/TensorRT optimized code**

0.45601

NVIDIA

# Lane and Object Detection using YOLO v2



Workflow:

1) Test in MATLAB

2) Generate code and test on desktop
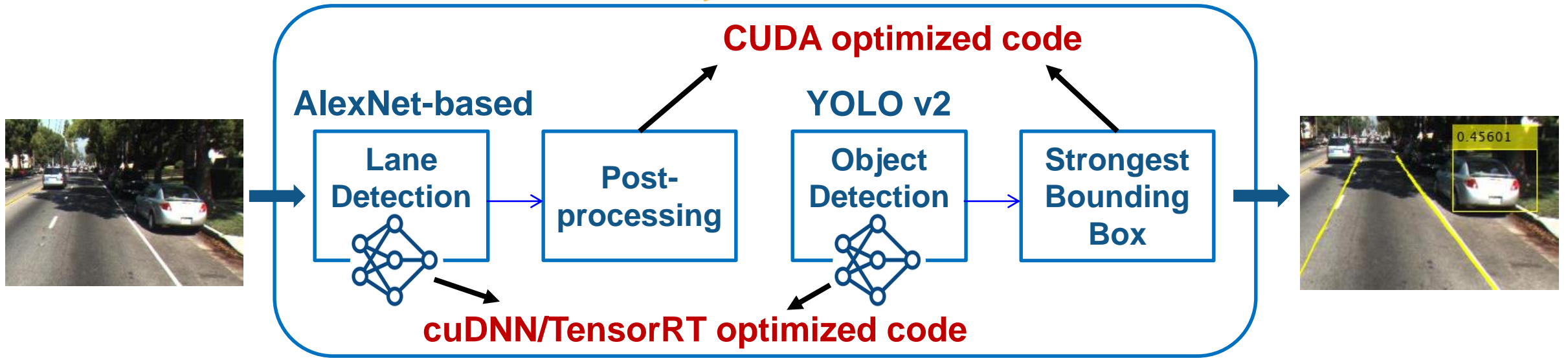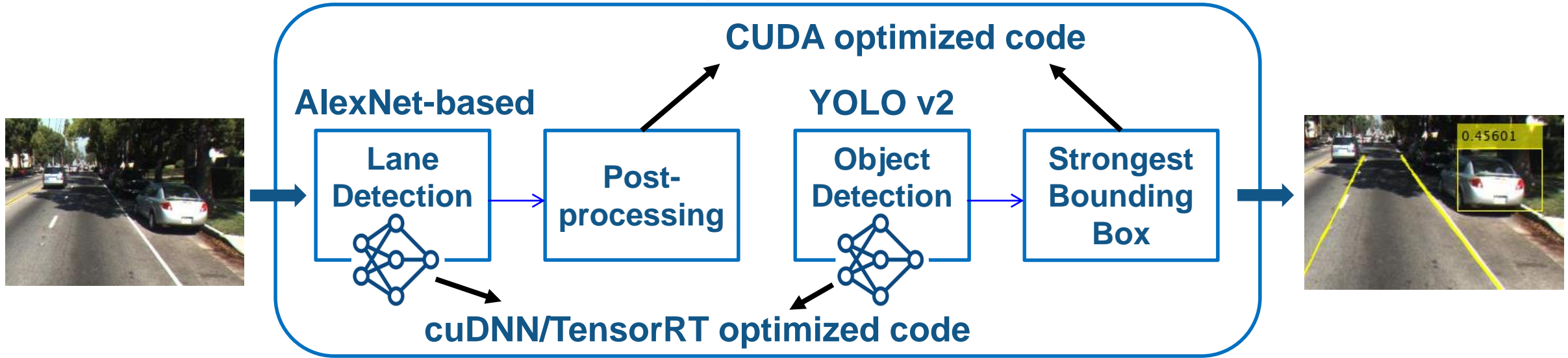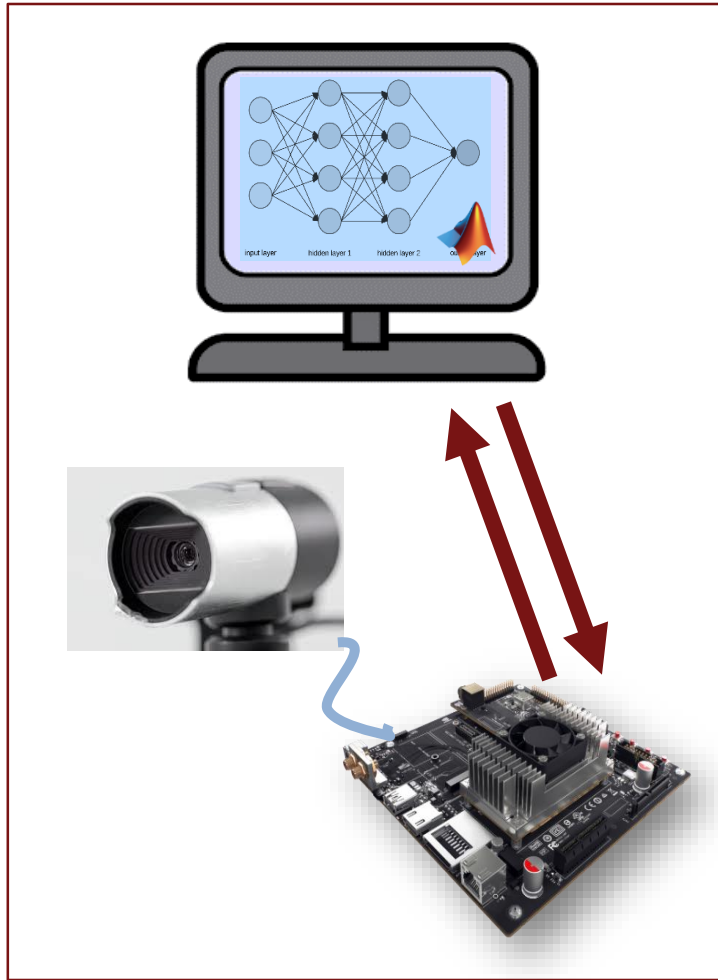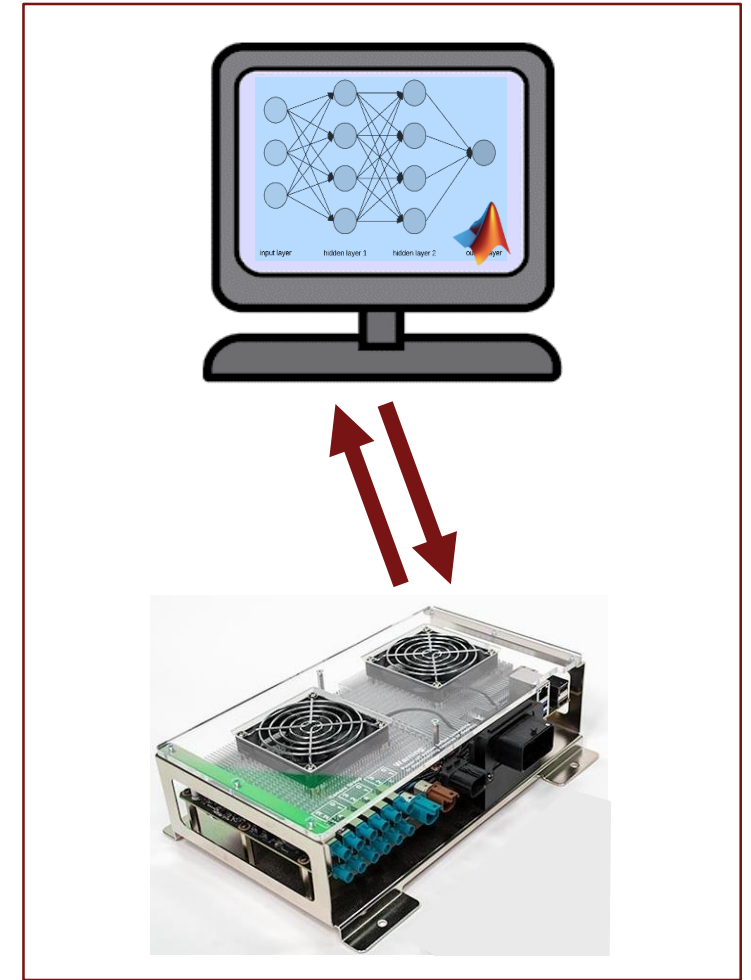
3) Generate code and test on Jetson AGX Xavier GPU

# Accessing Hardware
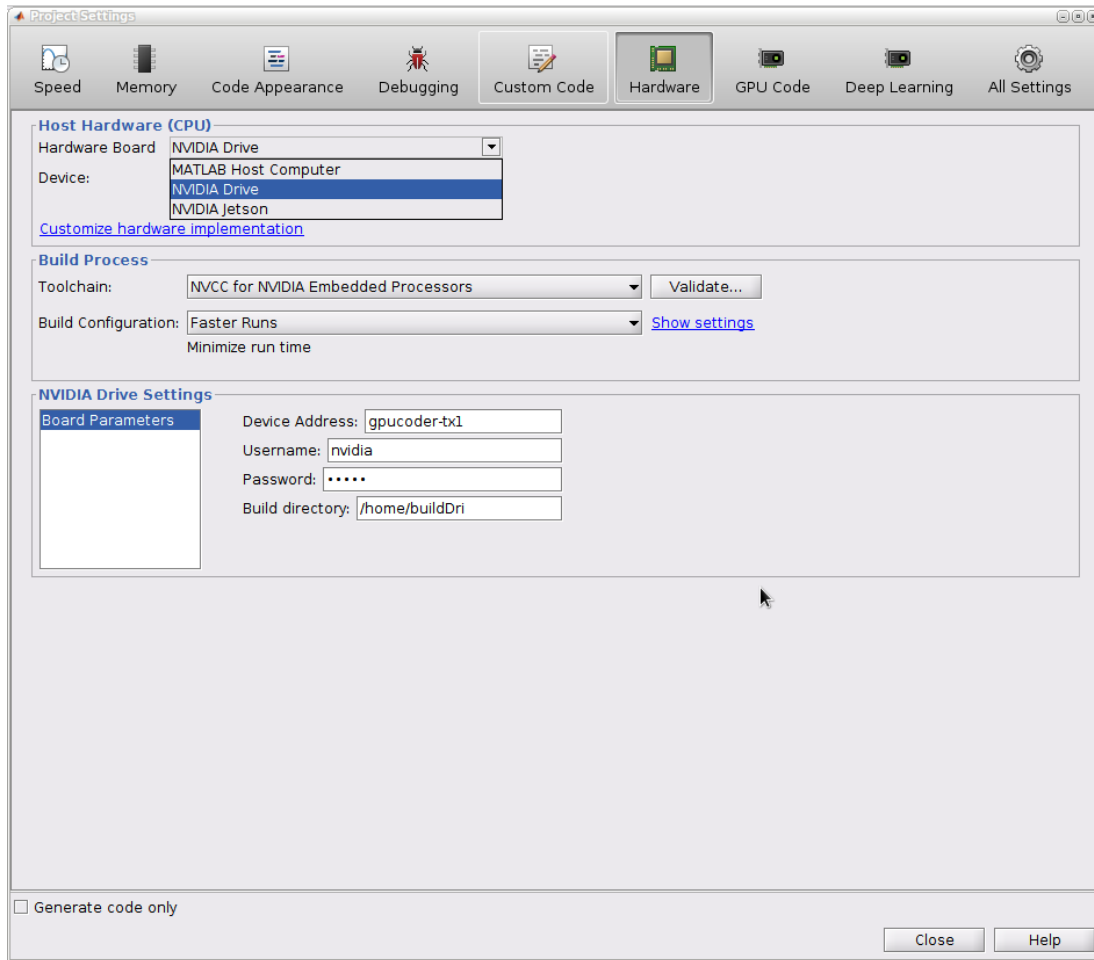


## Access Peripheral from MATLAB

## Deploy Standalone Application

## Processor-in-Loop Verification

# Deploy to Target Hardware via Apps and Command Line

Single Image Inference (Titan V, Linux)

**TensorFlow** (1.13.0)

**MXNet** (1.4.0)

**GPU Coder** (R2019a)

**PyTorch** (1.0.0)

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

# Coders Apply Various Optimizations

MATLAB

**Traditional compiler optimizations**

Library function mapping

Scalarization

Loop perfectization

Loop interchange

Loop fusion

Scalar replacement

Loop optimizations

Parallel loop creation

CUDA kernel creation

cudaMemcpy minimization

Shared memory mapping

CUDA kernel lowering

CUDA code emission

Deep Learning Workflow in MATLAB

## Deep Learning with MATLAB

This two-day course provides a comprehensive introduction to practical deep learning using MATLAB®.

**Topics include:**

- Importing image and sequence data
- Using convolutional neural networks for image classification, regression, and object detection
- Using long short-term memory networks for sequence classification and forecasting
- Modifying common network architectures to solve custom problems
- Improving the performance of a network by modifying training options



Transfer Learning



Classifying Sequence Data

# MATLAB EXPO 2019

Email: rishu.g@mathworks.com,

LinkedIn: https://www.linkedin.com/in/rishu-gupta-72148914/

# Please provide feedback for this block of sessions



- Scan this QR Code or log onto link below (link also sent to your phone and email)

- http://bit.ly/expo19-feedback

-  Enter the registration id number displayed on your badge

- Provide feedback for this session

Email: rishu.g@mathworks.com,

LinkedIn: https://www.linkedin.com/in/rishu-gupta-72148914/