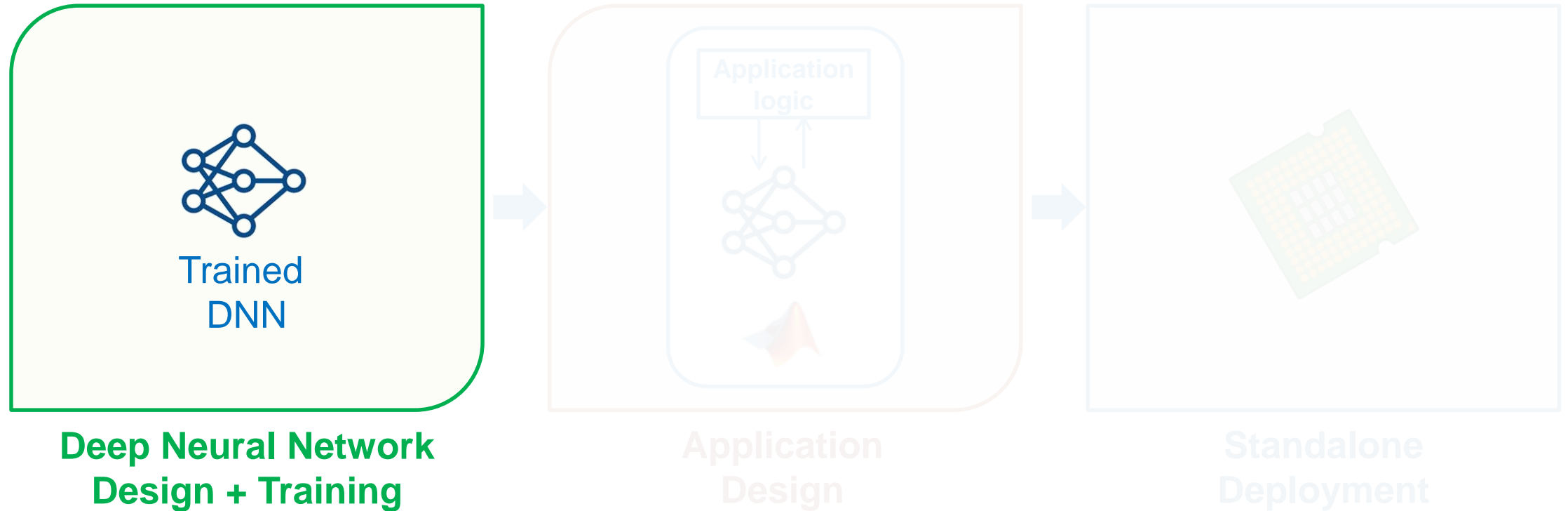# MATLAB EXPO 2019

## Deploying Deep Neural Networks to Embedded GPUs and CPUs
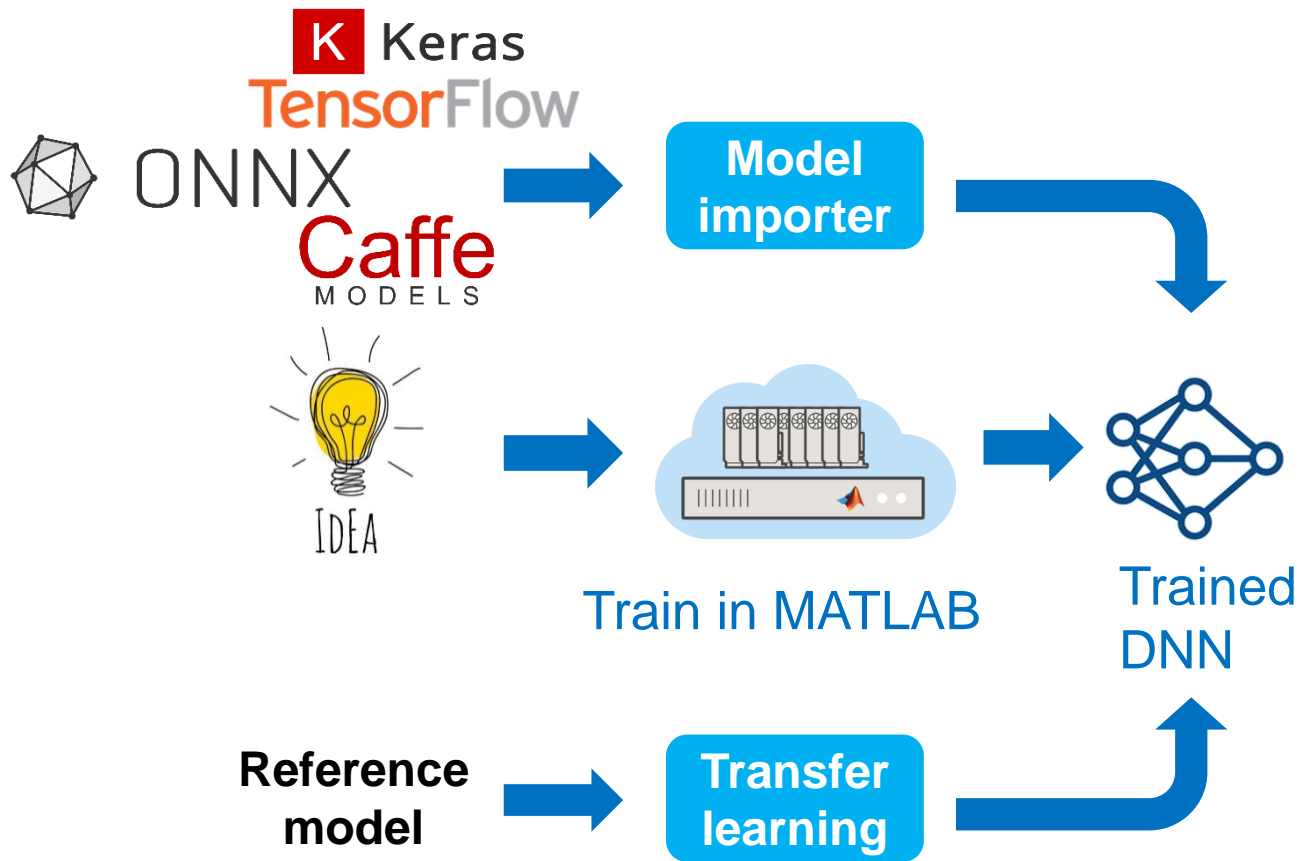
Abhijit Bhattacharjee

# Deep Learning Workflow in MATLAB



**Deep Neural Network
Design + Training**

Application
Design

Standalone
Deployment

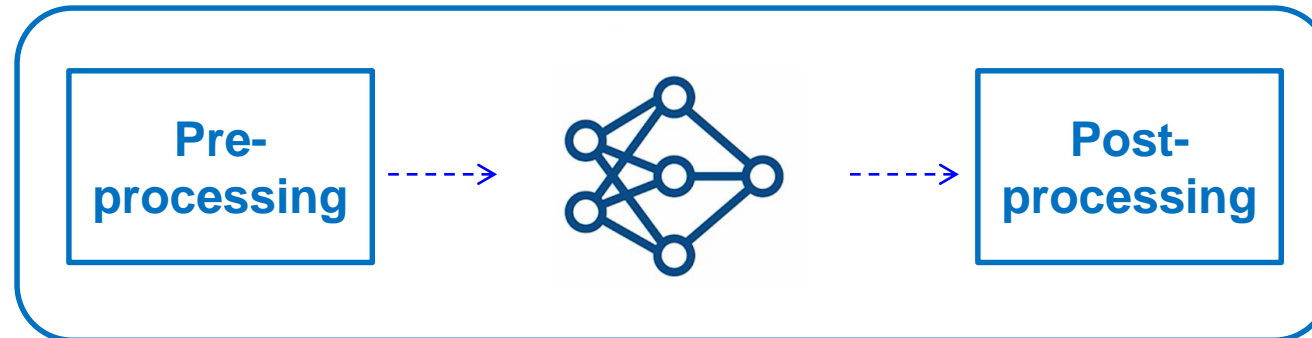# Deep Neural Network Design and Training



- **Design in MATLAB**
  - **Manage** large data sets
  - **Automate** data labeling
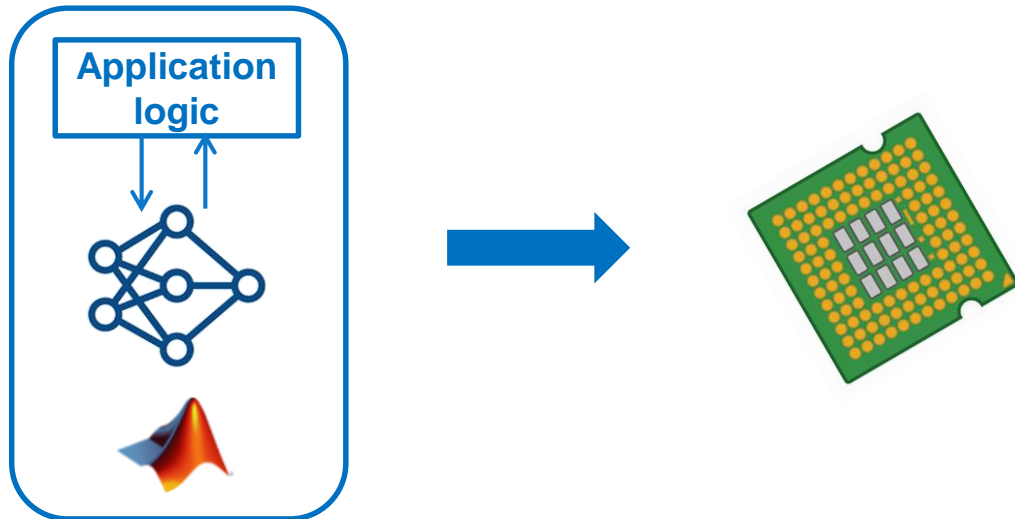  - **Easy access** to models

- **Training in MATLAB**
  - **Acceleration** with GPU's
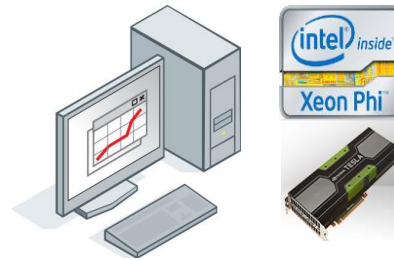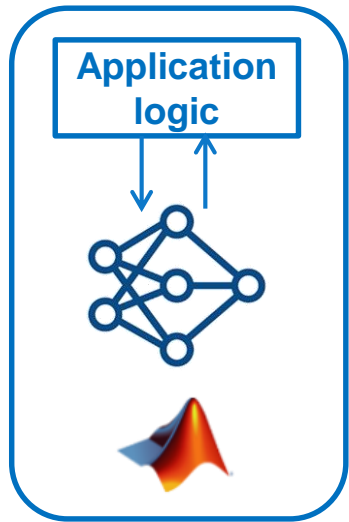  - **Scale** to clusters

# Application Design

# Multi-Platform Deep Learning Deployment

# Multi-Platform Deep Learning Deployment
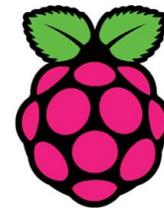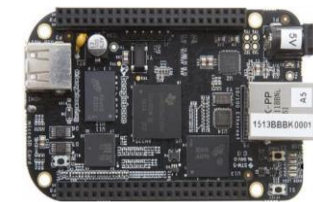
**Application logic**
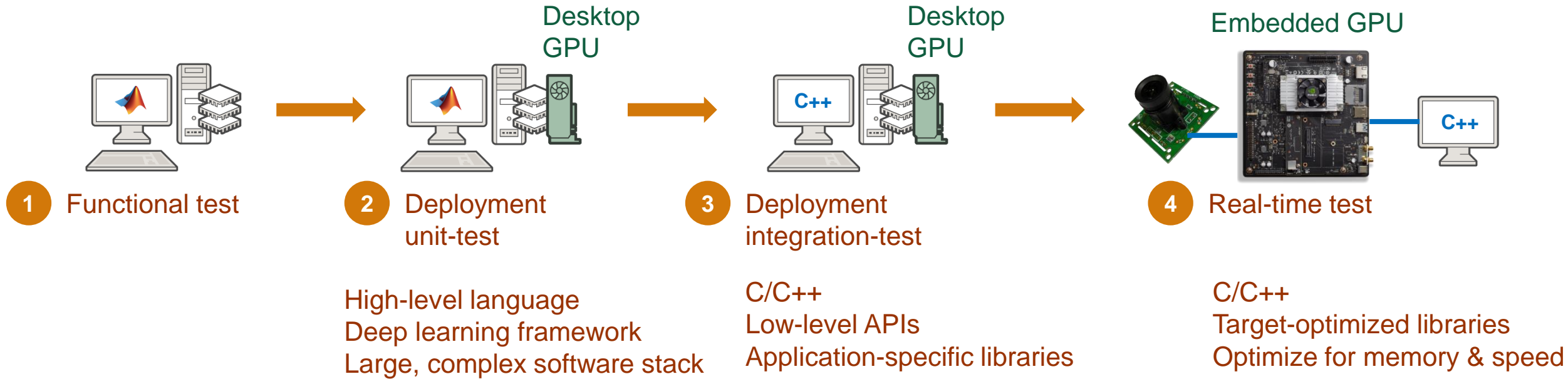
Desktop

Data Center

NVIDIA Jetson

Raspberry pi

Mobile

Beaglebone

. . .

Embedded

# Algorithm Design to Embedded Deployment Workflow
## *Conventional Approach*

Desktop GPU

Desktop GPU

Embedded GPU

**①** Functional test

**②** Deployment unit-test

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

**③** Deployment integration-test

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

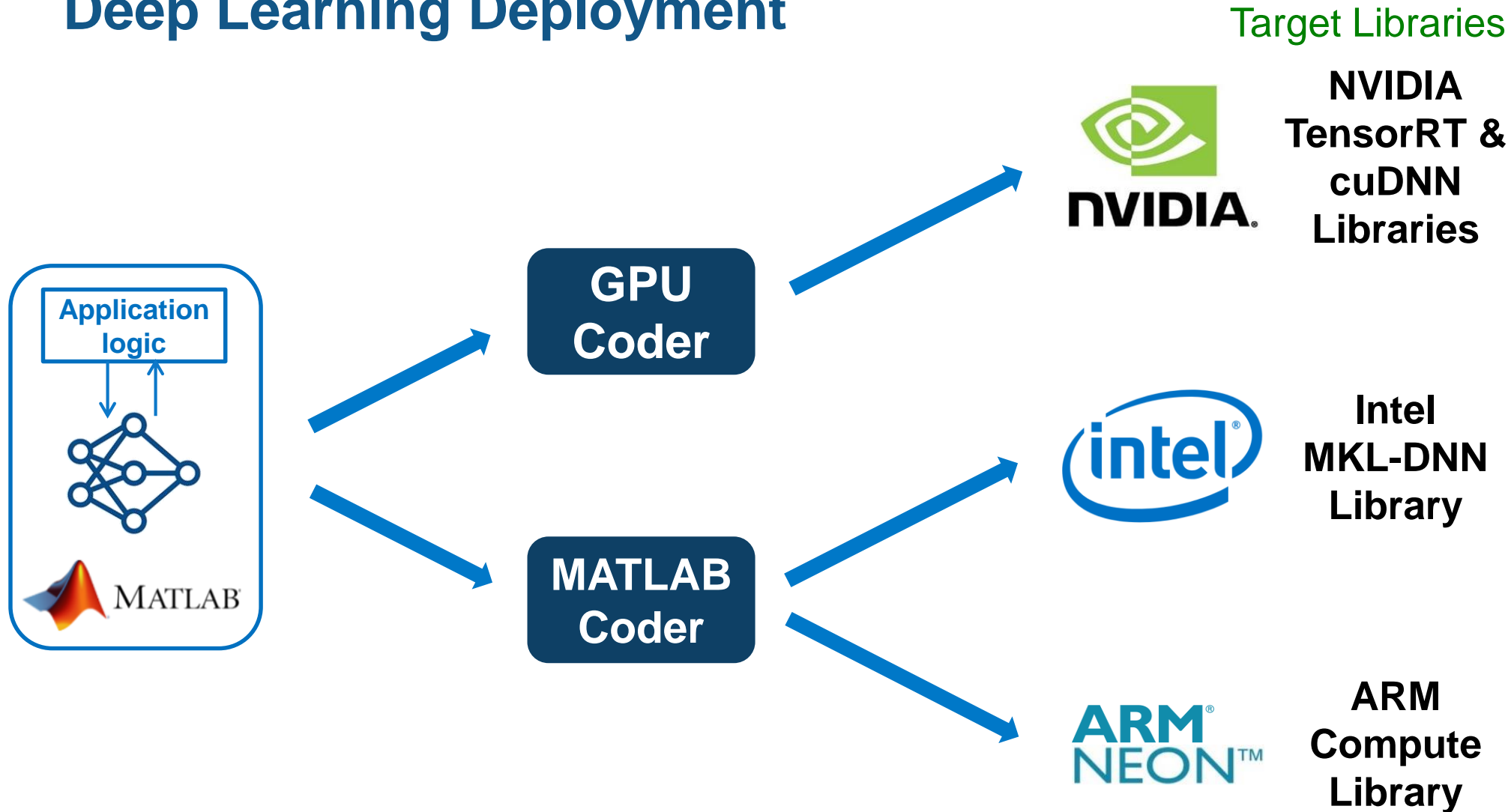**④** Real-time test

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

---

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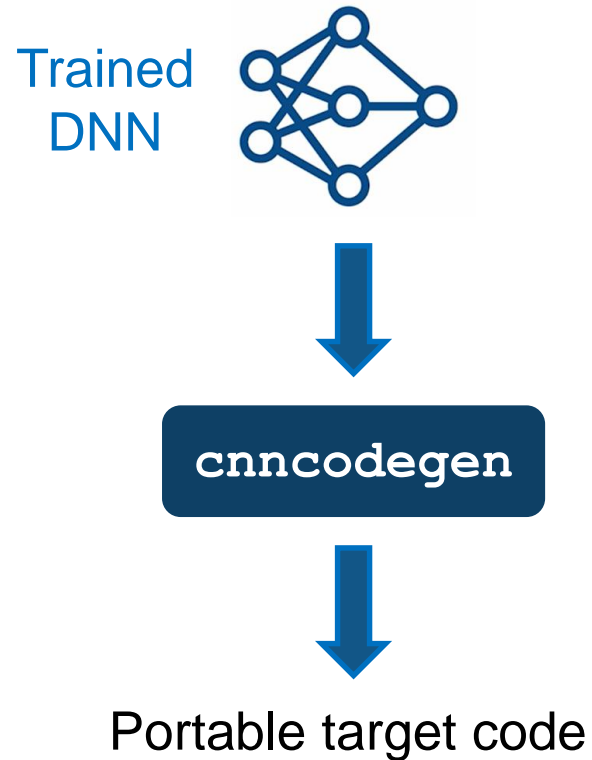- Verifying and maintaining multiple implementations
- Algorithm & vendor lock-in

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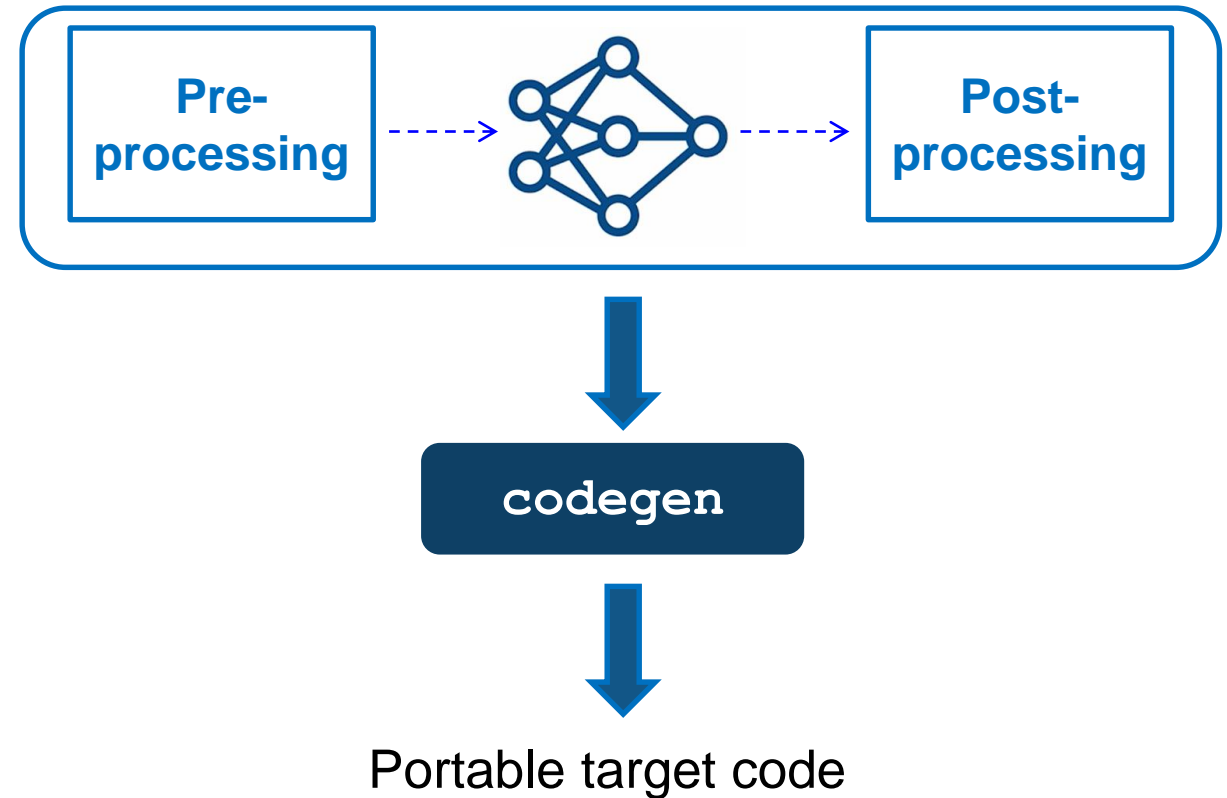
Target Libraries

**Application logic**

MATLAB

**GPU Coder** → **NVIDIA** → **NVIDIA TensorRT & cuDNN Libraries**

**MATLAB Coder** → **intel** → **Intel MKL-DNN Library**

→ **ARM NEON™** → **ARM Compute Library**

# 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**
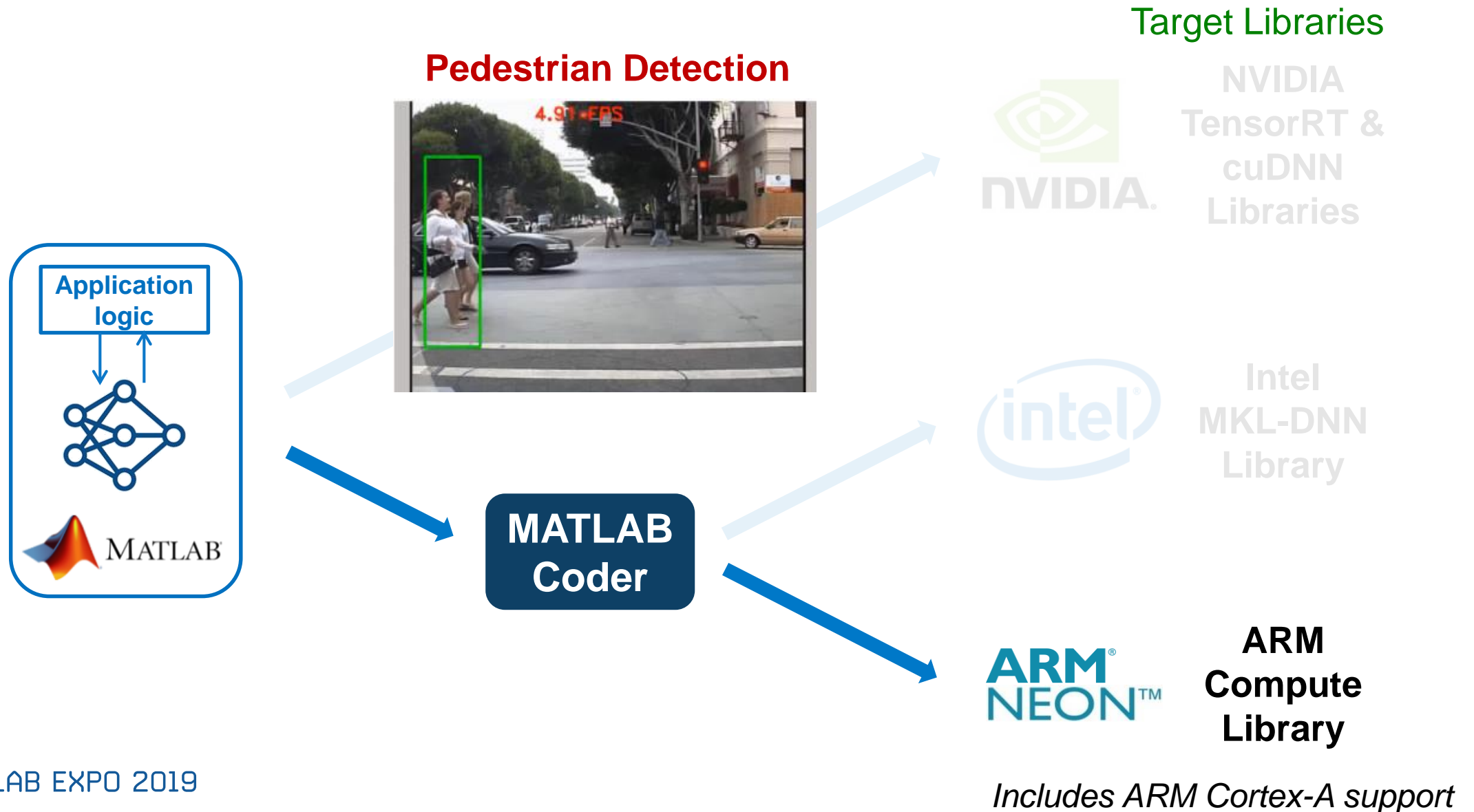
Trained
DNN



**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. Use hand written main function to call inference engine

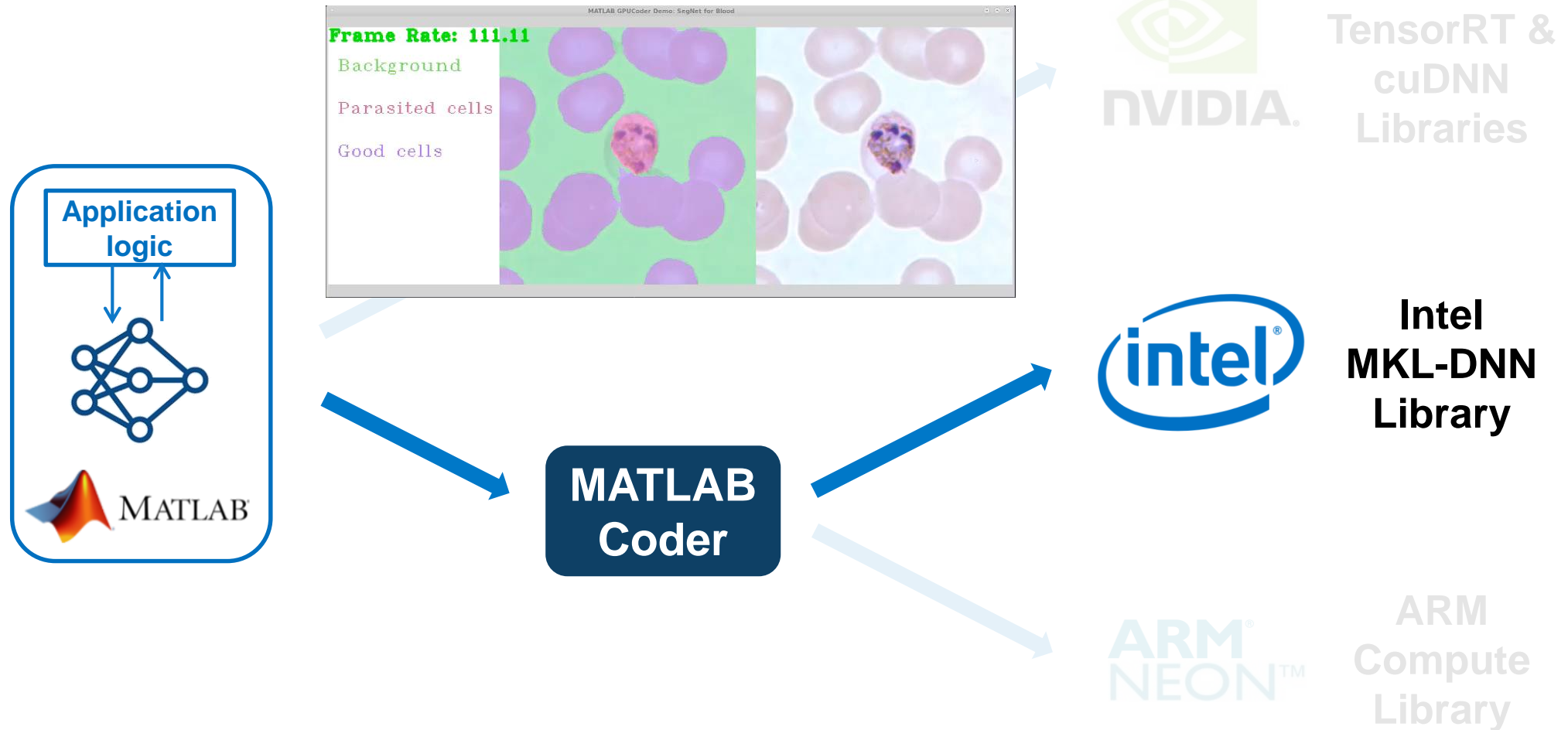4. Generate the exe and test the executable
```
>> make -C ./ ......
```
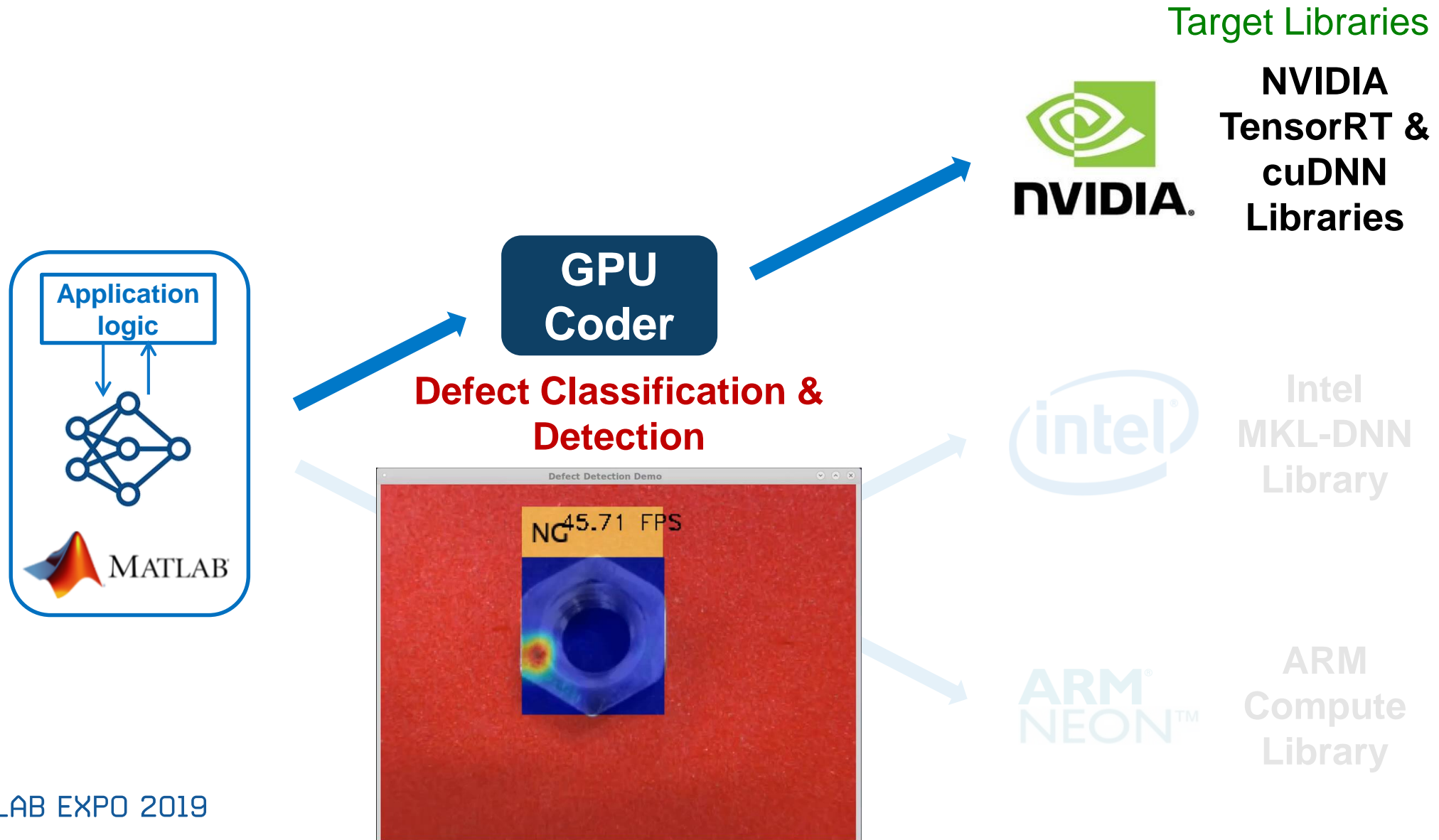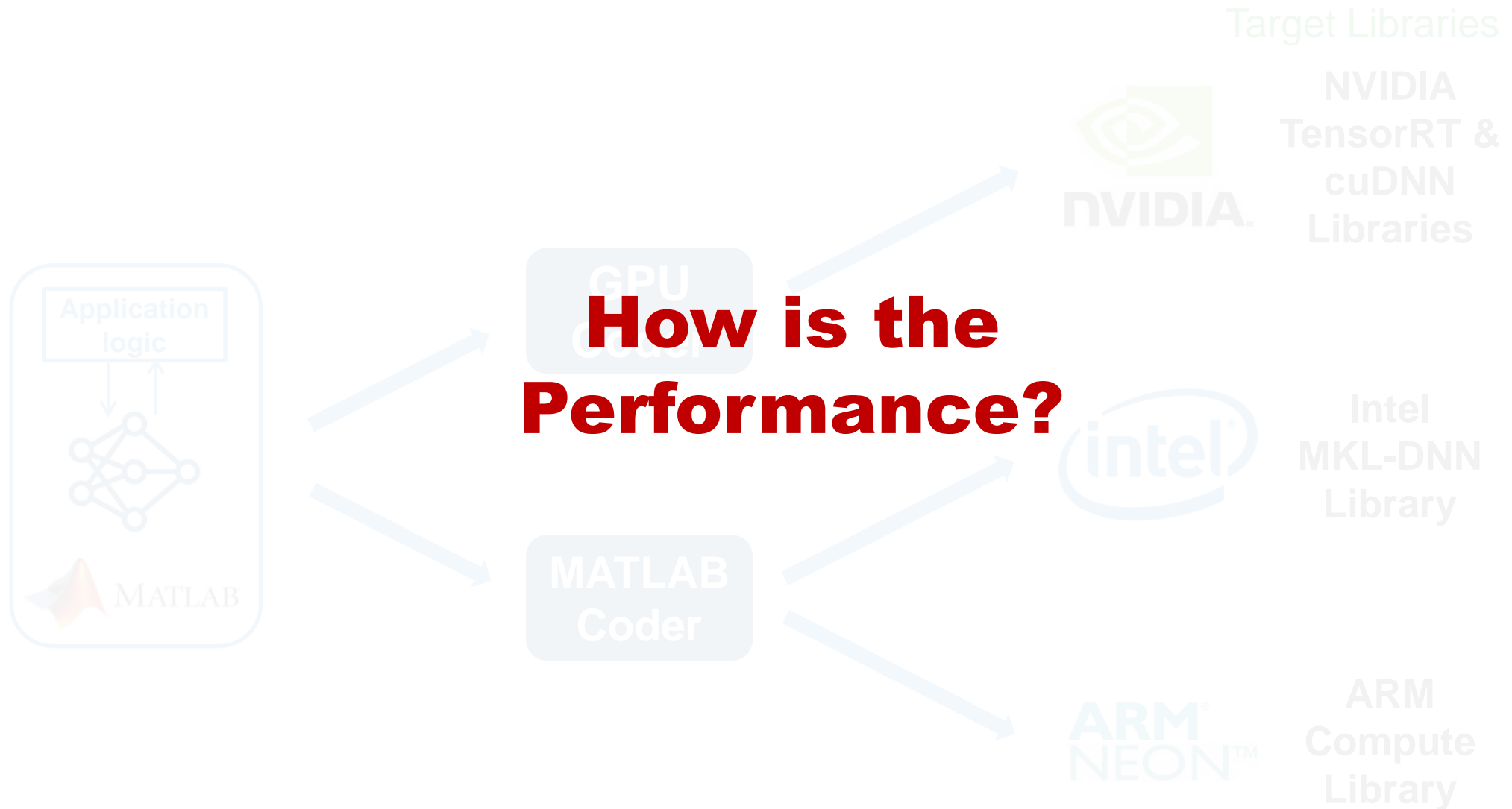
# Deep Learning Inference Deployment

**Target Libraries**

**Blood Smear Segmentation**



**Application logic**

MATLAB

**MATLAB Coder**

NVIDIA TensorRT & cuDNN Libraries

**Intel MKL-DNN Library**

ARM Compute Library

# Deep Learning Inference Deployment

Target Libraries

**NVIDIA TensorRT & cuDNN Libraries**

**Application logic**

MATLAB

**GPU Coder**

**Defect Classification & Detection**

Defect Detection Demo

NG 45.71 FPS

Intel MKL-DNN Library

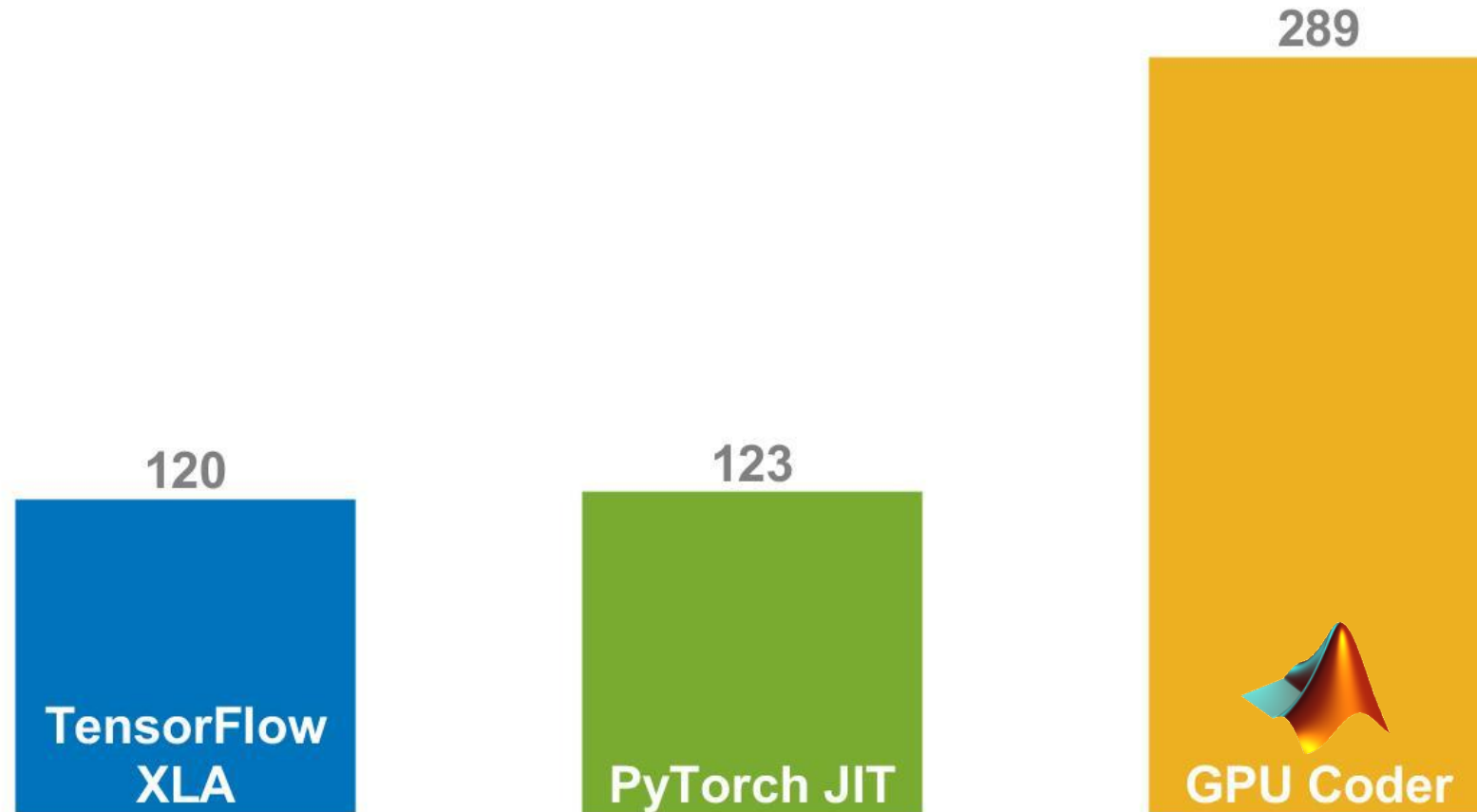ARM NEON™ ARM Compute Library

# 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

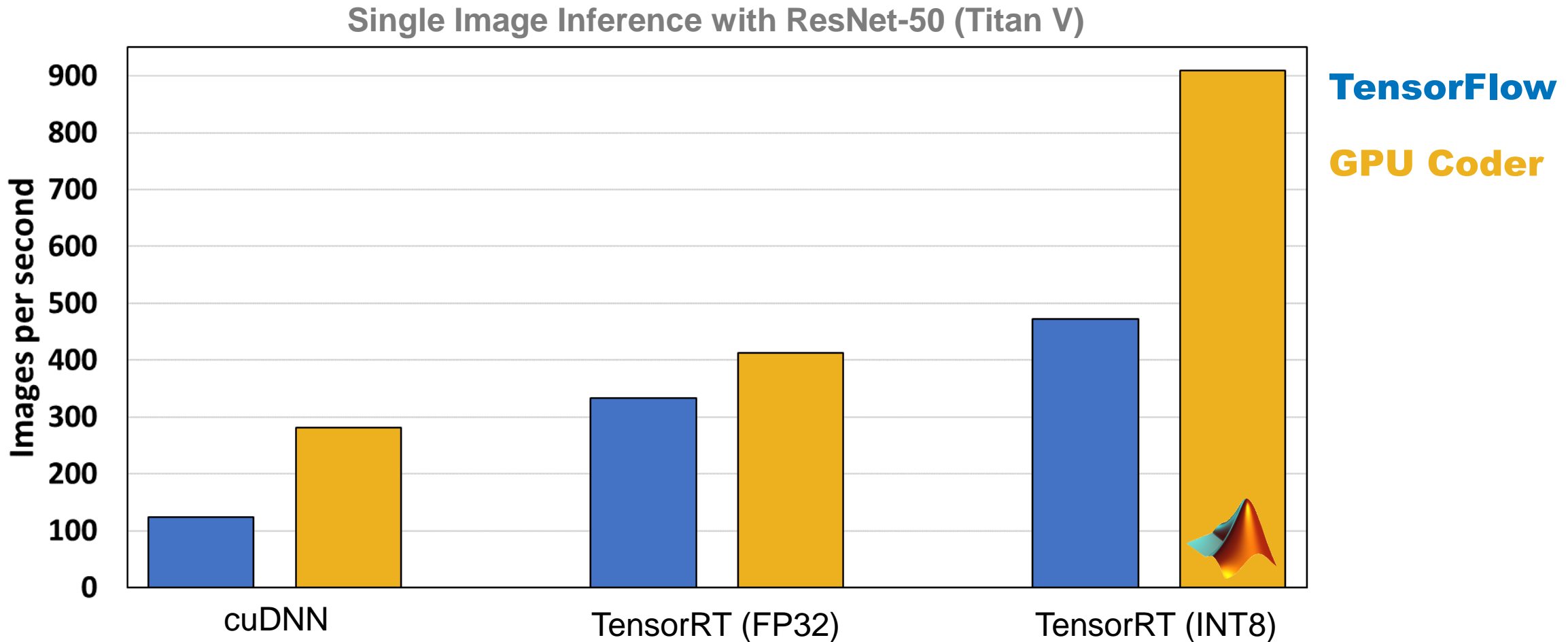- CNN inference (ResNet-50 , VGG-16, Inception V3) on Intel Xeon CPU

# TensorRT Accelerates Inference Performance on Titan V

**Single Image Inference with ResNet-50 (Titan V)**

**TensorFlow**

**GPU Coder**

*Intel® Xeon® CPU 3.6 GHz - NVIDIA libraries: CUDA10.0/1 - cuDNN 7.5.0 - TensorRT 5.1.2 - Frameworks: TensorFlow 1.13.1*

Single Image Inference on Jetson TX2

ResNet-50

*NVIDIA libraries: CUDA10.0/1 - cuDNN 7.5.0 - TensorRT 5.1.2 - Frameworks: TensorFlow 1.13.1*

# CPU Performance

CPU, Single Image Inference (Linux)

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

MATLA

*Intel® Xeon® CPU 3.6 GHz - NVIDIA libraries: CUDA10.0/1 - cuDNN 7.5.0 - Frameworks: TensorFlow 1.13.1, MXNet 1.4.1 PyTorch 1.1*     **24**

# Brief Summary

**DNN libraries are great for inference,   …**

MATLAB Coder and GPU Coder generates code that takes advantage of:
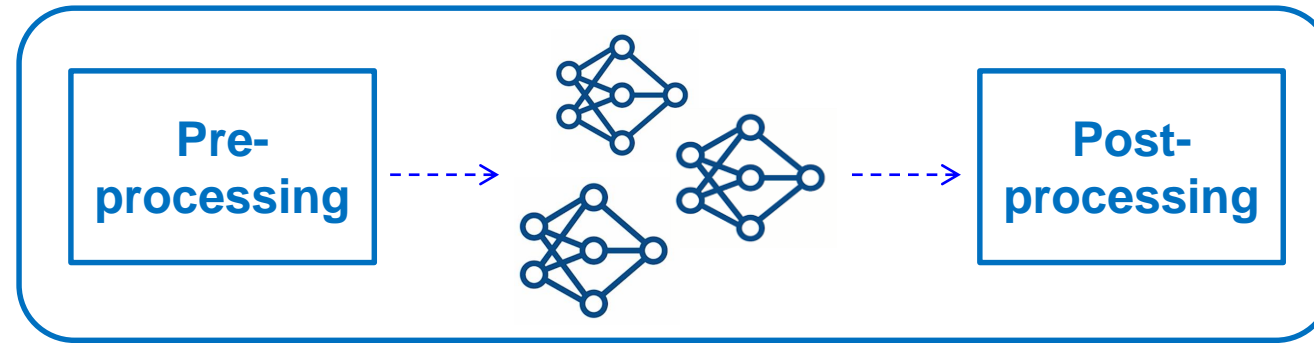
NVIDIA® CUDA libraries, including TensorRT & cuDNN

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:

NVIDIA® CUDA libraries, including TensorRT & cuDNN

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

ARM® Compute libraries for mobile platforms

## But, Applications Require More than just Inference

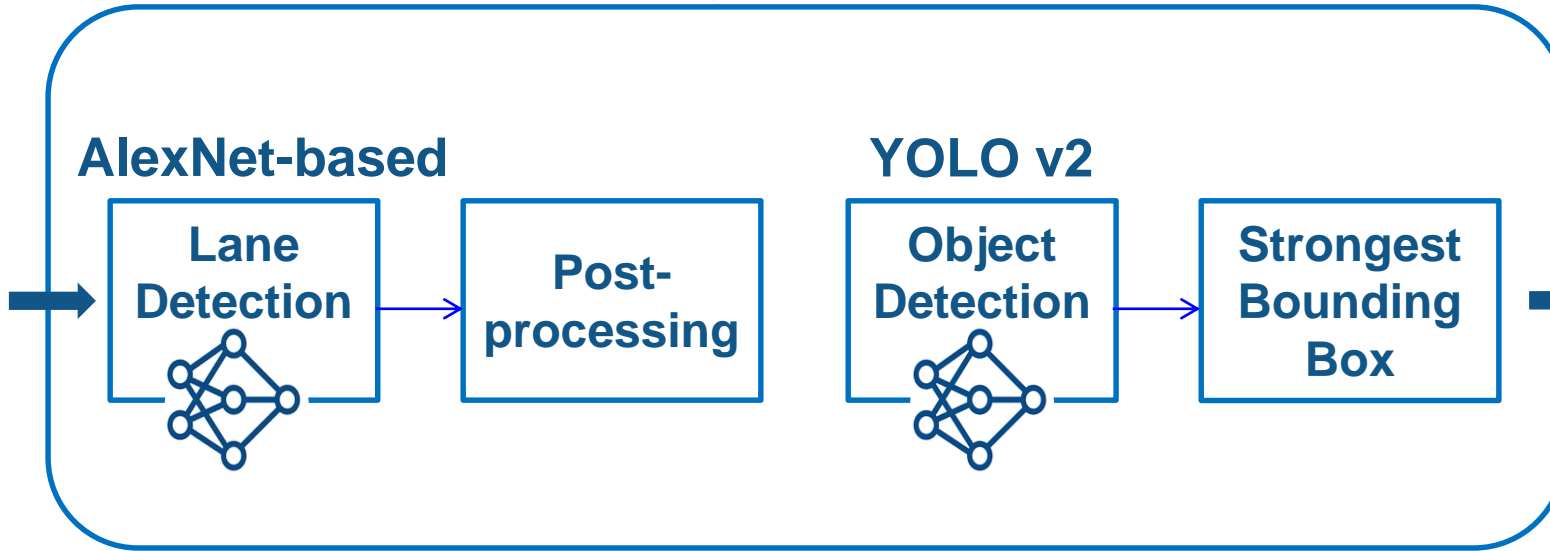# Deep Learning Workflows: Integrated Application Deployment



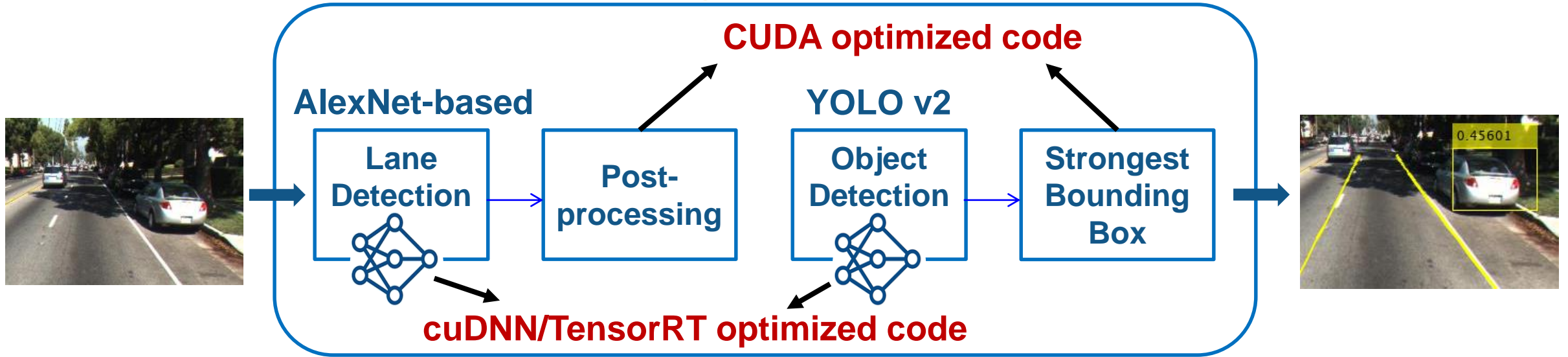Portable target code

# Lane and Object Detection using YOLO v2



**AlexNet-based**

Lane Detection → Post-processing

**YOLO v2**

Object Detection → Strongest Bounding Box

Workflow:

1) Test in MATLAB on CPU

2) Generate code and test on desktop GPU

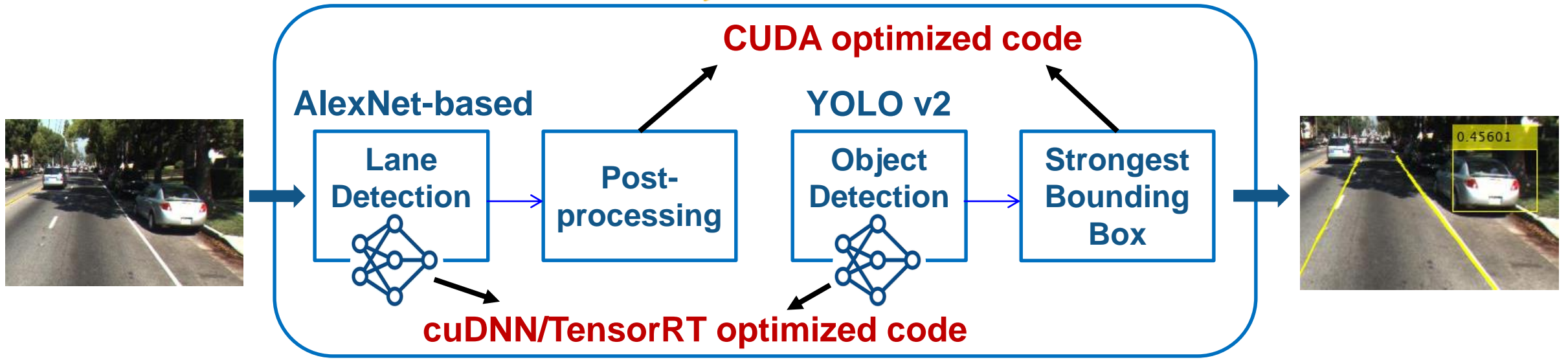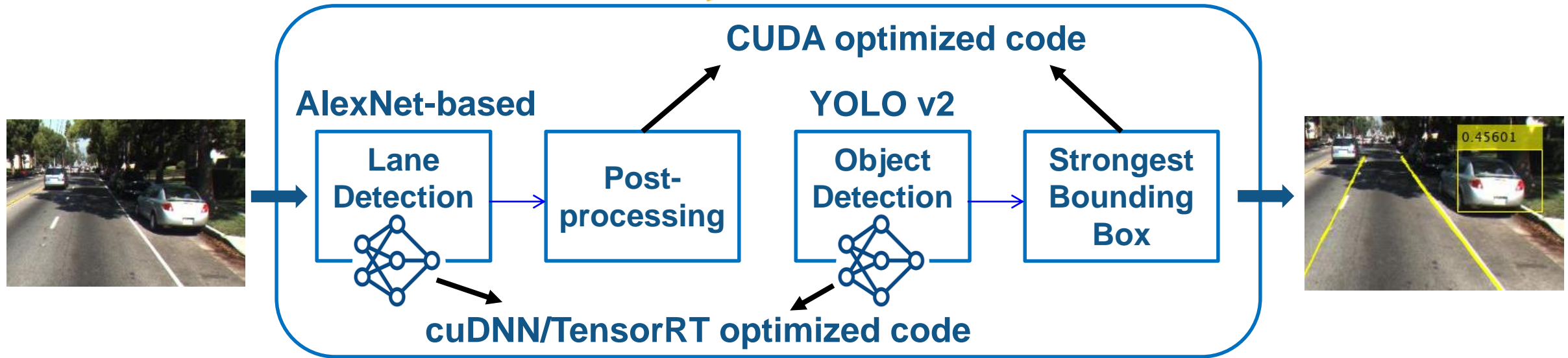3) Generate code and test on Jetson AGX Xavier GPU

# (1) Test in MATLAB on CPU
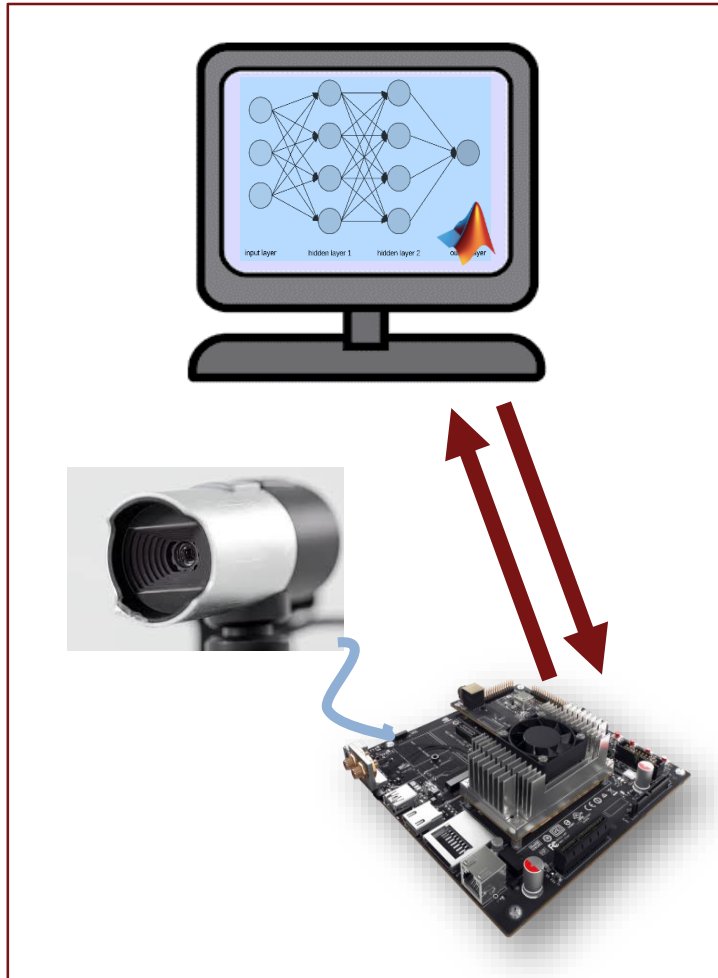
# (2) Generate Code and Test on Desktop GPU



**CUDA optimized code**

AlexNet-based

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

YOLO v2

**cuDNN/TensorRT optimized code**

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



**CUDA optimized code**

AlexNet-based

YOLO v2

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

**cuDNN/TensorRT optimized code**

# 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

Inference Speed - ResNet-50 (Img/Sec)

289

120 — TensorFlow XLA

123 — PyTorch JIT

GPU Coder

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

# Generated Code Calls Optimized Libraries

**cuFFT, cuBLAS, cuSolver, Thrust Libraries**

**Pre-processing** → **Post-processing**

**Intel MKL-DNN Library**

**NVIDIA TensorRT & cuDNN Libraries**

**ARM Compute Library**

# Deep Learning Network Optimization

Network

Layer fusion
Optimized computation

Buffer minimization
Optimized memory

# Coding Patterns: Stencil Kernels

▪ Automatically applied for image processing functions (e.g. imfilter, imerode, imdilate, conv2, …)
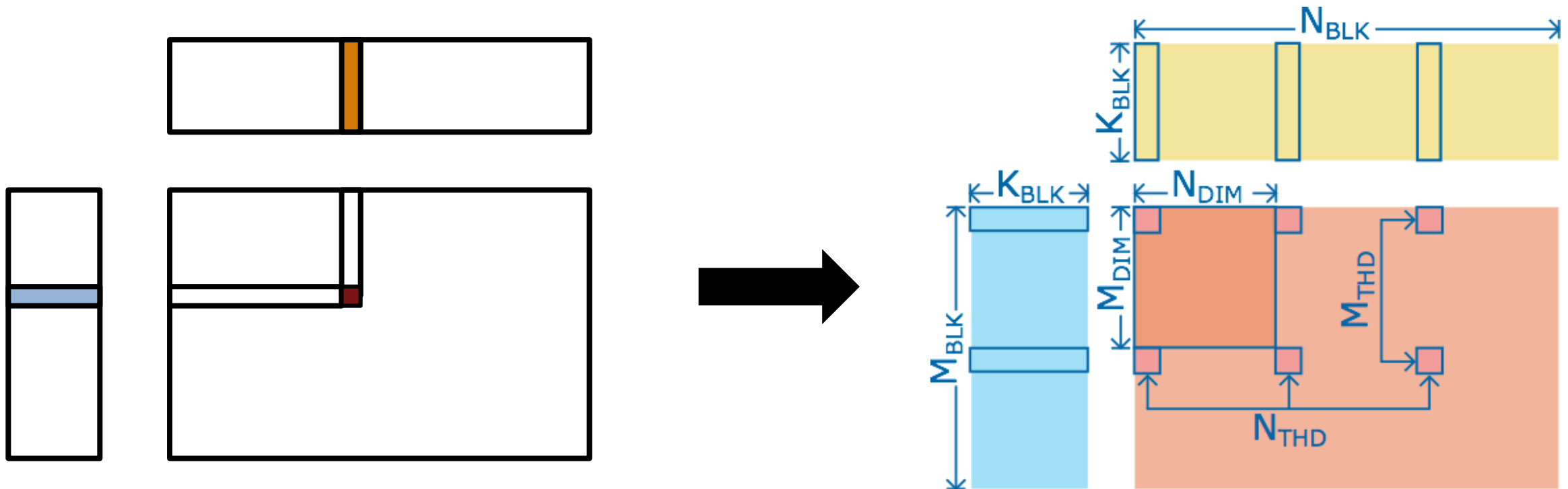
▪ Manually apply using *gpucoder.stencilKernel()*

Input image

Conv. kernel

Output image

cols

rows
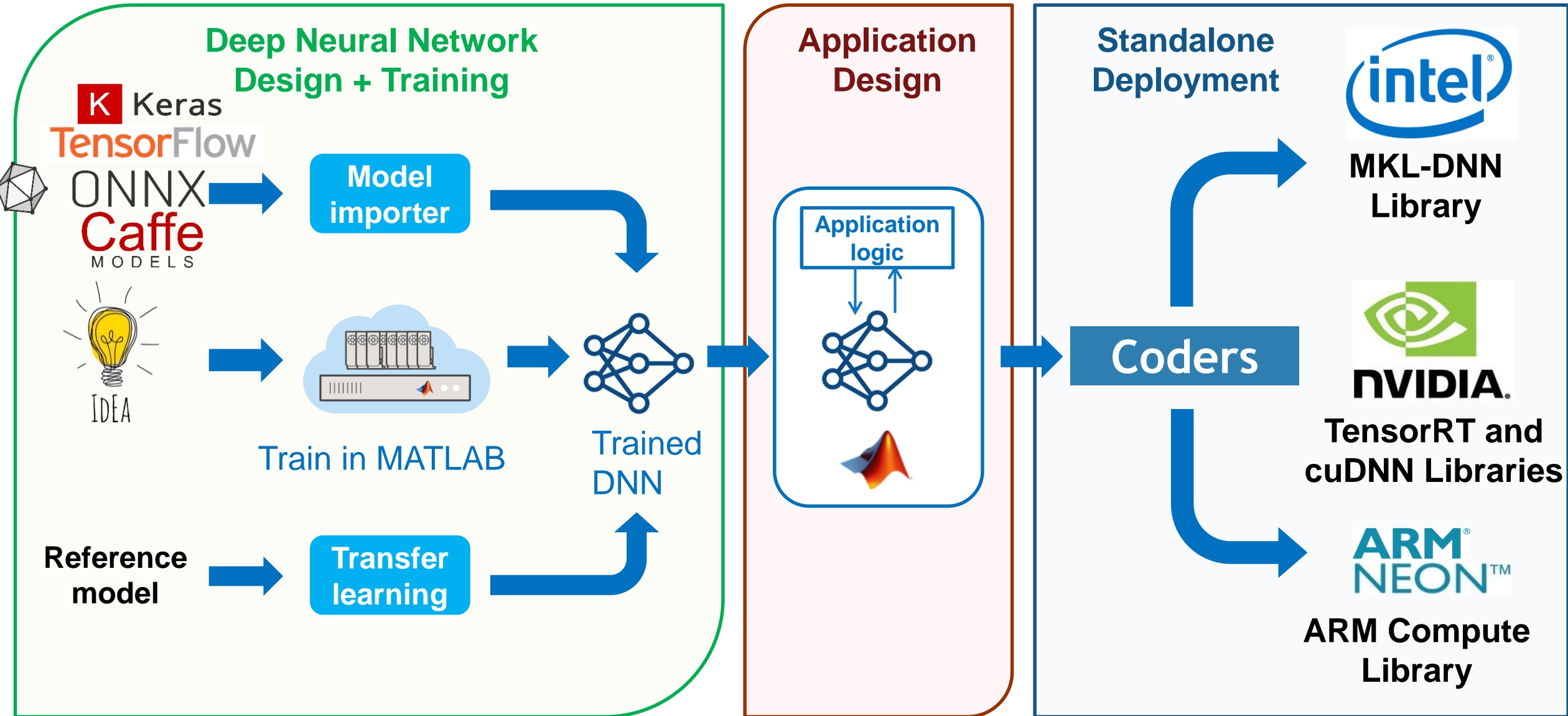
kh

kw

Dotprod

# Coding Patterns: Matrix-Matrix Kernels

- Automatically applied for many MATLAB functions (e.g. matchFeatures SAD, SSD, pdist, …)

- Manually apply using *gpucoder.matrixMatrixKernel()*

# Deep Learning Workflow in MATLAB



**Deep Neural Network Design + Training**

**Application Design**

**Standalone Deployment**

# Deep Learning Workflow in MATLAB

# Call to action

- Visit the Deep Learning Booth!

- Related upcoming talks:
    - AI Techniques for Signals, Time-series, and Text Data
    - Sensor Fusion and Tracking for Autonomous Systems
    - Deploying Deep Neural Networks to Embedded GPUs and CPUs