

# MATLAB EXPO 2018

## Deploying Deep Learning Networks to Embedded GPUs and CPUs

Rishu Gupta, PhD

*Senior Application Engineer, Computer Vision*

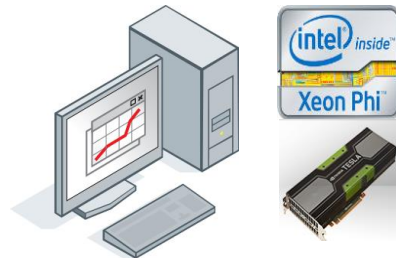


# MATLAB Deep Learning Framework



- **Manage** large image sets
- **Automate** image labeling
- **Easy access** to models
- **Acceleration** with GPU's
- **Scale** to clusters

# Multi-Platform Deep Learning Deployment



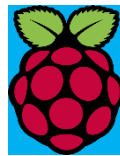
Desktop



Data-center



Nvidia  
TX1, TX2, TK1



Raspberry pi



Mobile



Beagle bone



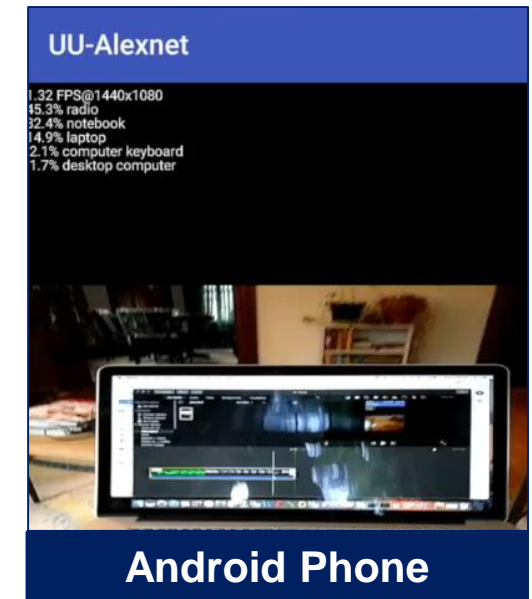
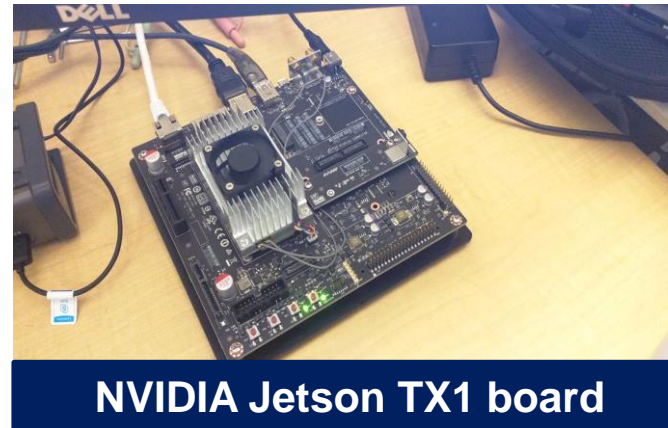
Embedded

# Multi-Platform Deep Learning Deployment

- Need code that takes advantage of:
  - NVIDIA<sup>®</sup> CUDA libraries, including cuDNN and TensorRT
  - Intel<sup>®</sup> Math Kernel Library for Deep Neural Networks (MKL-DNN) for Intel processors
  - ARM<sup>®</sup> Compute libraries for ARM processors

# Multi-Platform Deep Learning Deployment

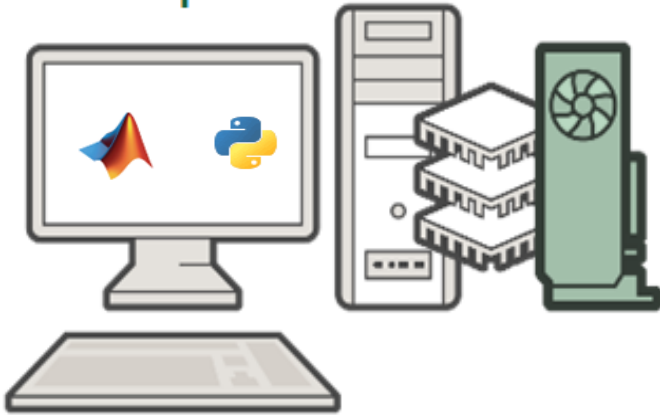
- Need code that takes advantage of:
  - NVIDIA® CUDA libraries, including cuDNN and TensorRT
  - Intel® Math Kernel Library for Deep Neural Networks (MKL-DNN) for Intel processors
  - ARM® Compute libraries for ARM processors



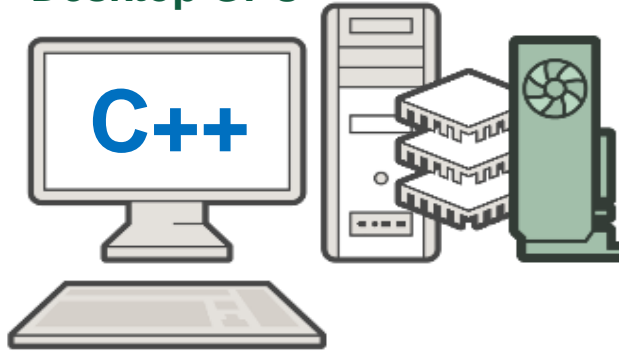
# Algorithm Design to Embedded Deployment Workflow

## *Conventional Approach*

Desktop GPU



Desktop GPU



**Challenges**

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

1

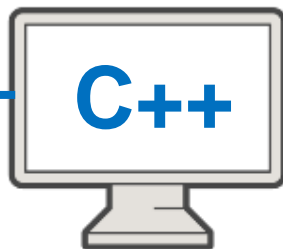
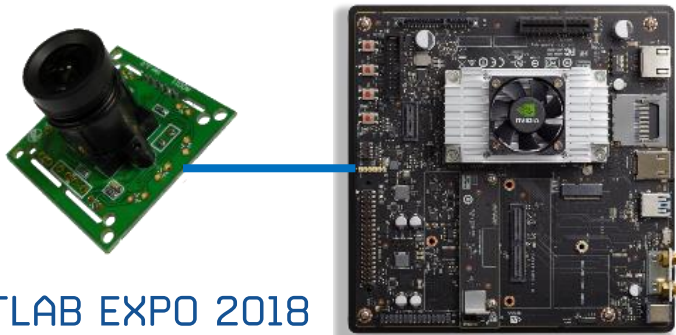
High-level language  
Deep learning framework  
Large, complex software stack

2

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



Embedded GPU



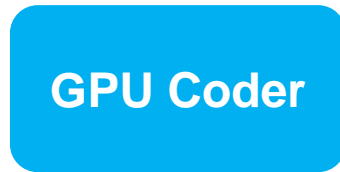
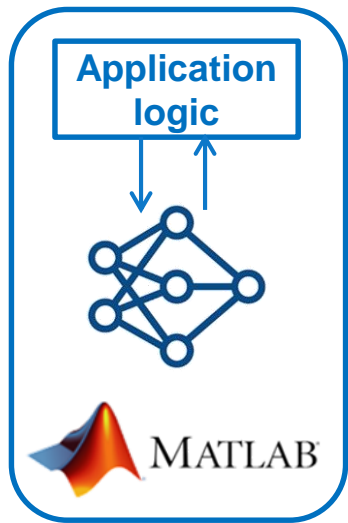
3

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

MATLAB EXPO 2018

# Solution- GPU Coder for Deep Learning Deployment

Target Libraries



Intel  
MKL-DNN  
Library



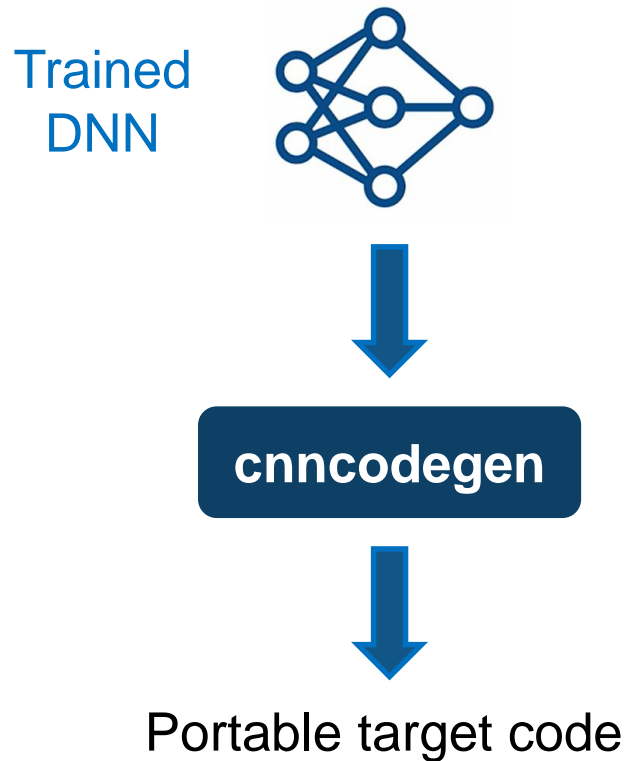
NVIDIA  
TensorRT &  
cuDNN  
Libraries



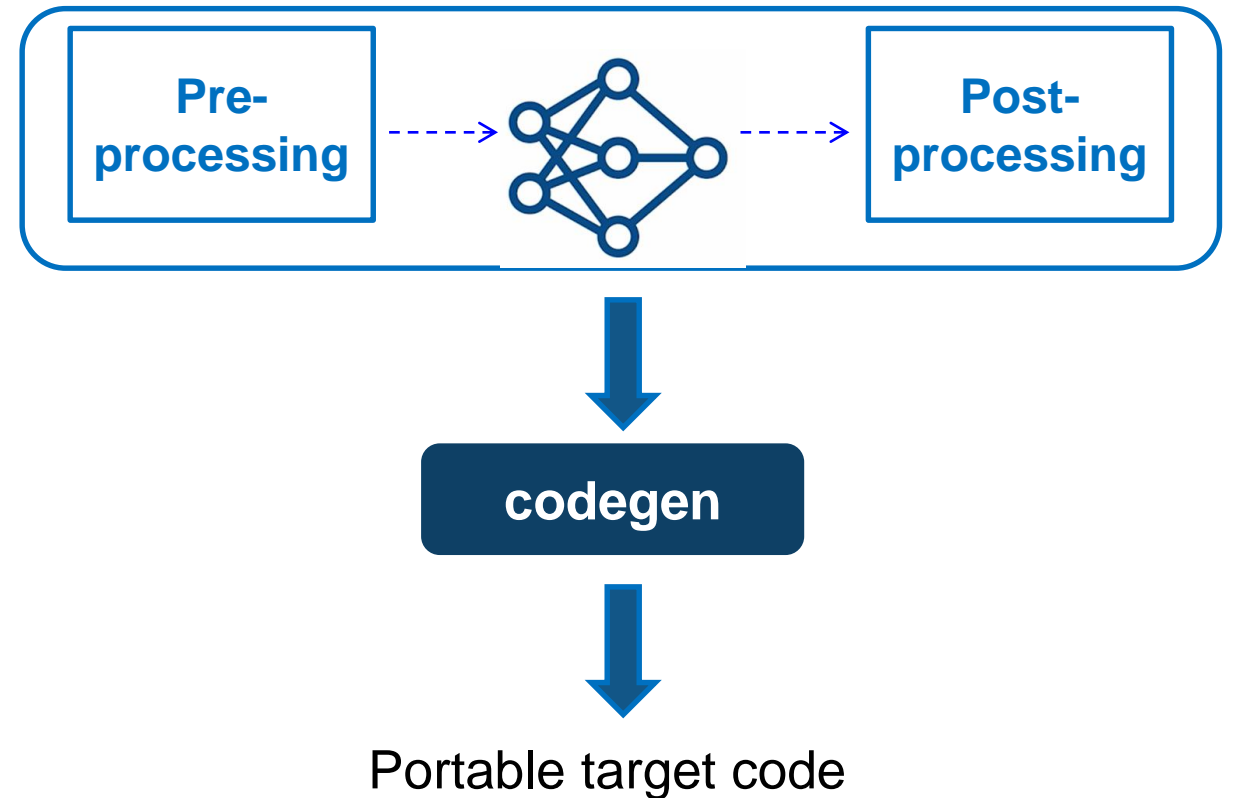
ARM  
Compute  
Library

# Deep Learning Deployment Workflows

## INFERENCE ENGINE DEPLOYMENT



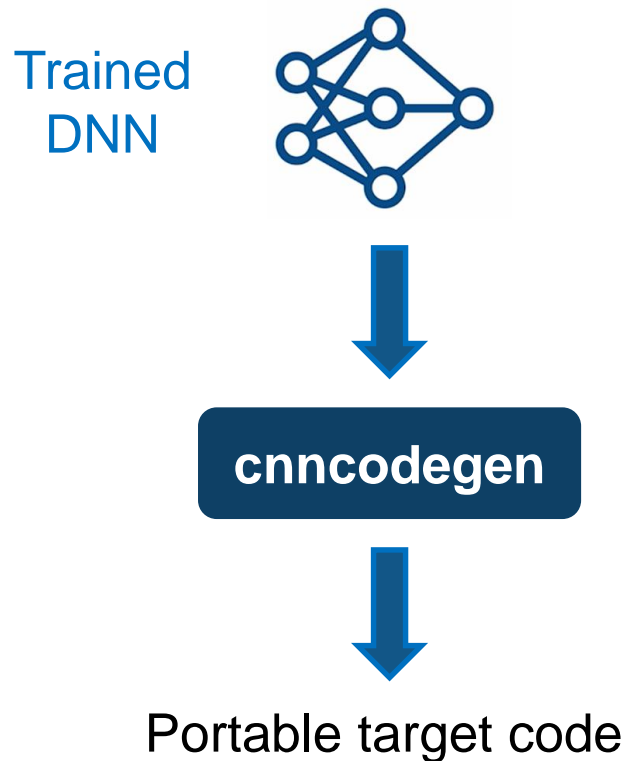
## INTEGRATED APPLICATION DEPLOYMENT





# Workflow for Inference Engine Deployment

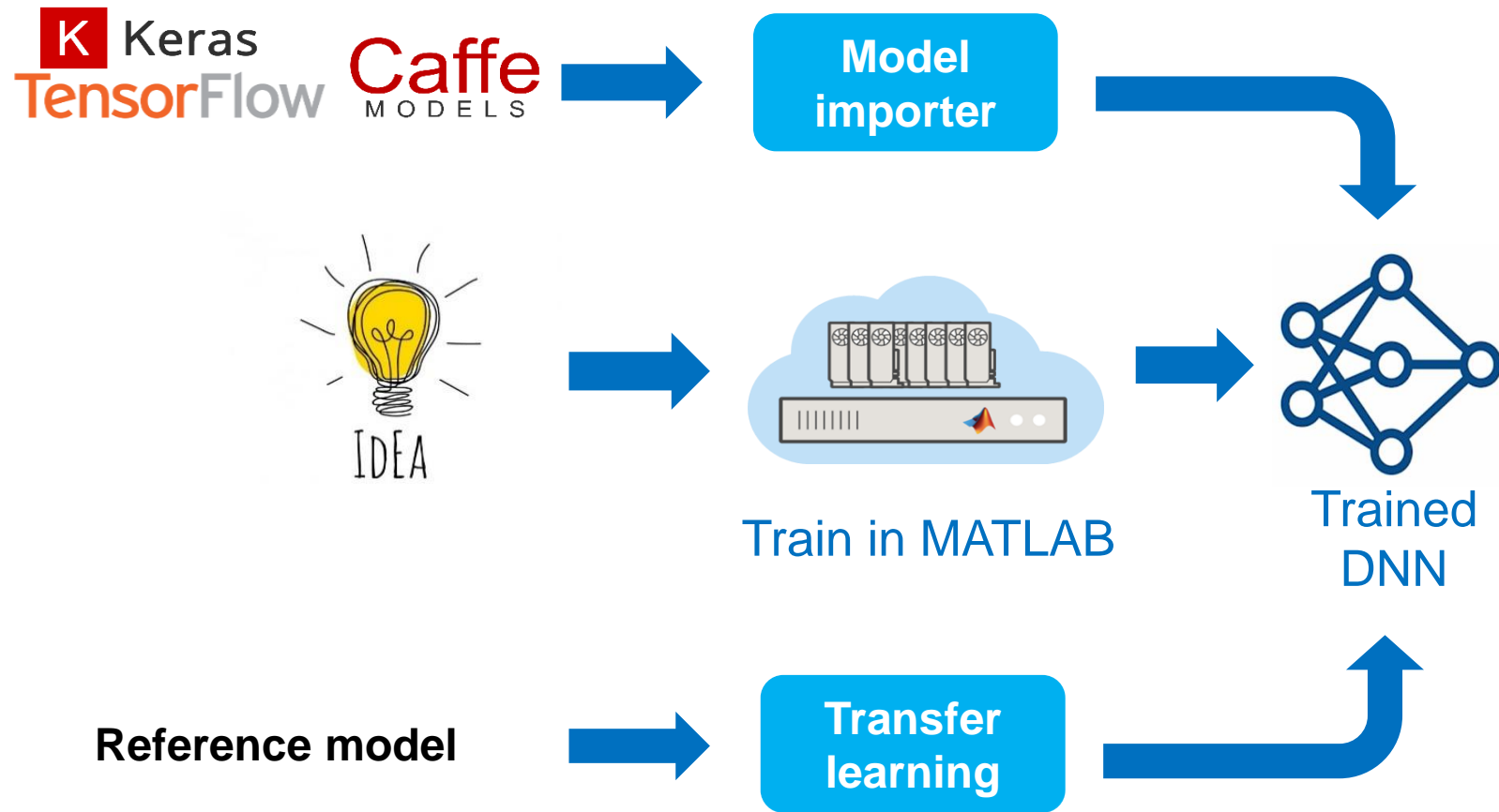
## INFERENCE ENGINE DEPLOYMENT



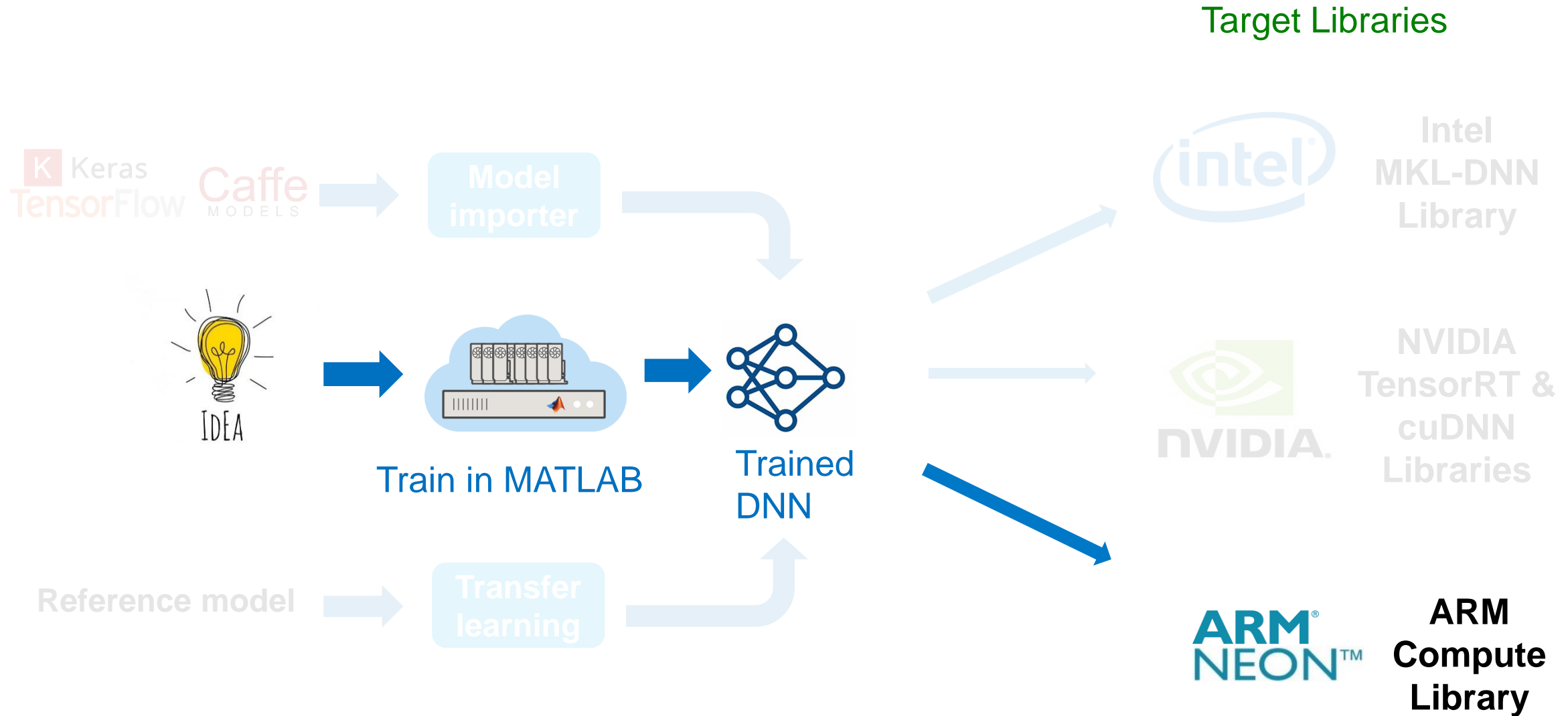
### Steps for inference engine deployment

1. Generate the code for trained model  
`>> cnncodegen(net, 'targetlib', 'cudnn')`
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 ./ .....`

# How to get a Trained DNN into MATLAB?



# Deep Learning Inference Deployment



# Building DNN from Scratch

Load Training Data



Build Layer Architecture



Set Training Options



Train Network

```
%% Create a datastore
imds = imageDatastore('Data',...
    'IncludeSubfolders',true,'LabelSource','foldernames');
num_classes = numel(unique(imds.Labels));

%% Build layer architecture
layers = [imageInputLayer([64 32 3])
    convolution2dLayer(5,20)
    reluLayer()
    maxPooling2dLayer(2,'Stride',2)
    fullyConnectedLayer(512)
    fullyConnectedLayer(2)
    softmaxLayer()
    classificationLayer()];

%% Set Training Options
trainOpts = trainingOptions('sgdm',...
    'MiniBatchSize', miniBatchSize,...
    'Plots', 'training-progress');

%% Train Network
net = trainNetwork(imds, layers, trainOpts);
```

# Pedestrian Detection DNN Deployment on ARM Processor

```
layers = [imageInputLayer([64 32 3])  
convolution2dLayer(5,20)  
reluLayer()  
maxPooling2dLayer(2,'Stride',2)  
CrossChannelNormalizationLayer(5,'K',1);  
convolution2dLayer(5,20)  
reluLayer()  
maxPooling2dLayer(2,'Stride',2)  
fullyConnectedLayer(512)  
fullyConnectedLayer(2)  
softmaxLayer()  
classificationLayer()];
```

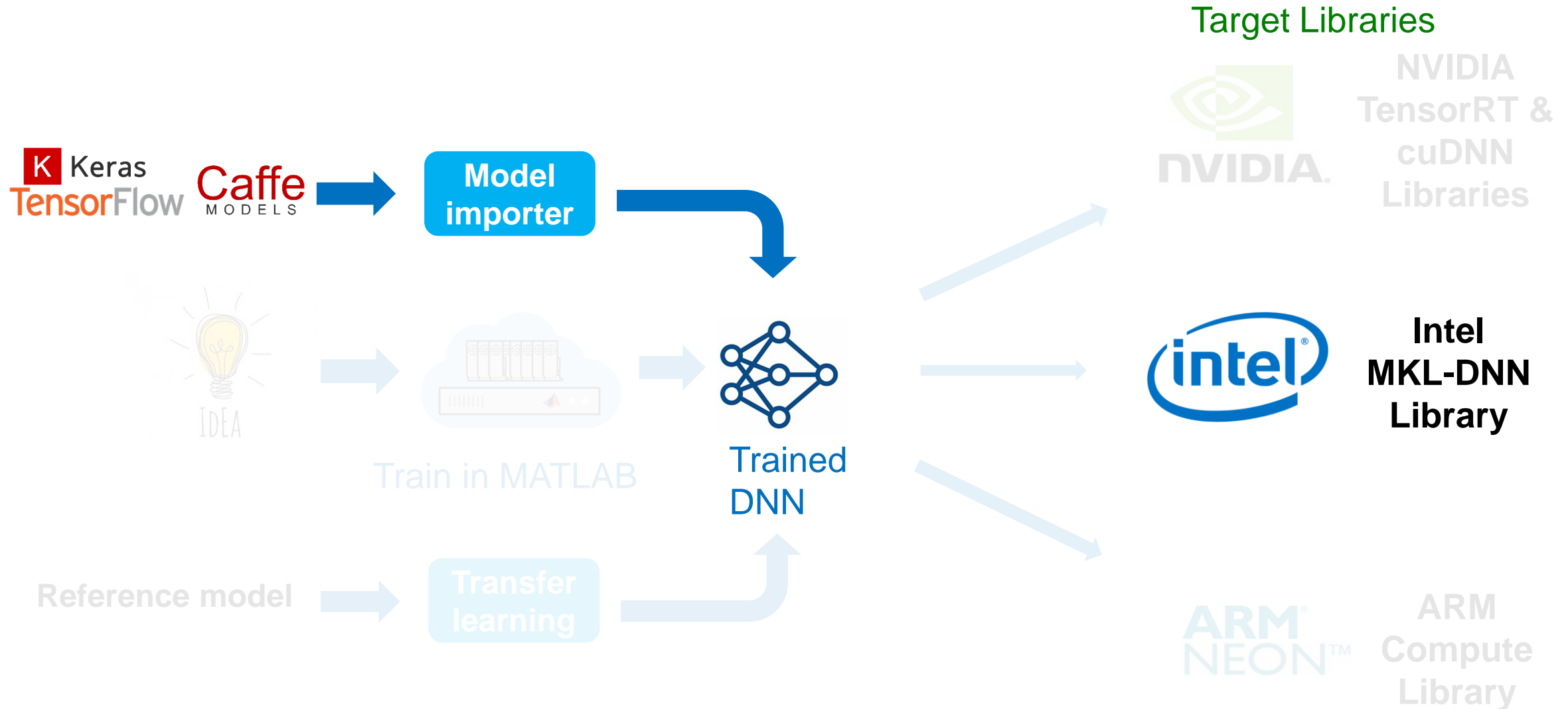


# Pedestrian Detection DNN Deployment on ARM Processor

- ARM Neon instruction set architecture
  - Example: ARM Cortex A
- ARM Compute Library
  - Low-level Software functions
  - Computer vision, machine learning etc...
- Pedestrian detection on Raspberry pi



# Deep Learning Inference Deployment



# Importing DNN from Open Source Framework

## Caffe Model Importer (including Caffe Model Zoo)



Caffe  
MODELS

- `importCaffeLayers`

- `importCaffeNetwork`

```
network = importCaffeNetwork(protofile, 'yolo.caffemodel');
```



Add-Ons



## TensorFlow-Keras Model Importer

- `importKerasLayers`

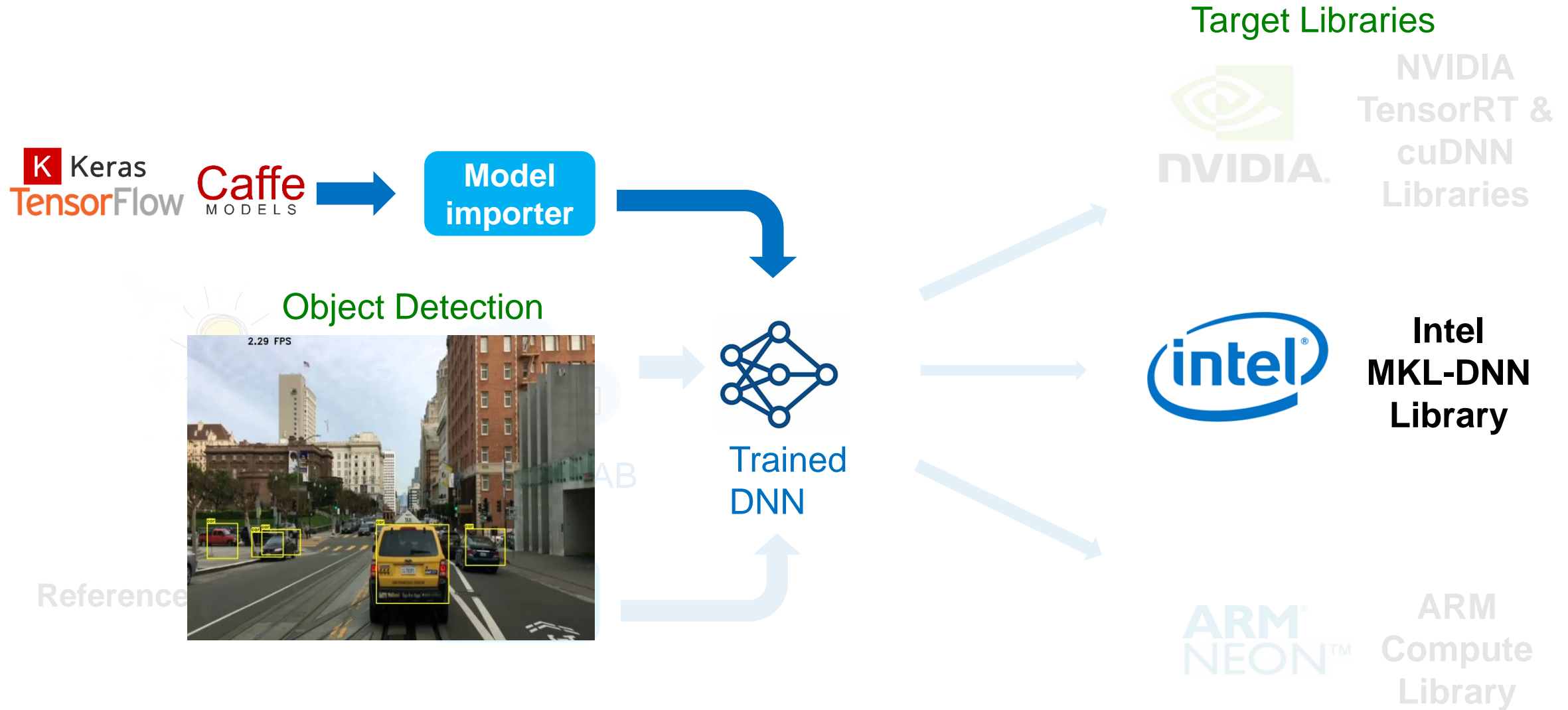
- `importKerasNetwork`

# KERAS IMPORTER

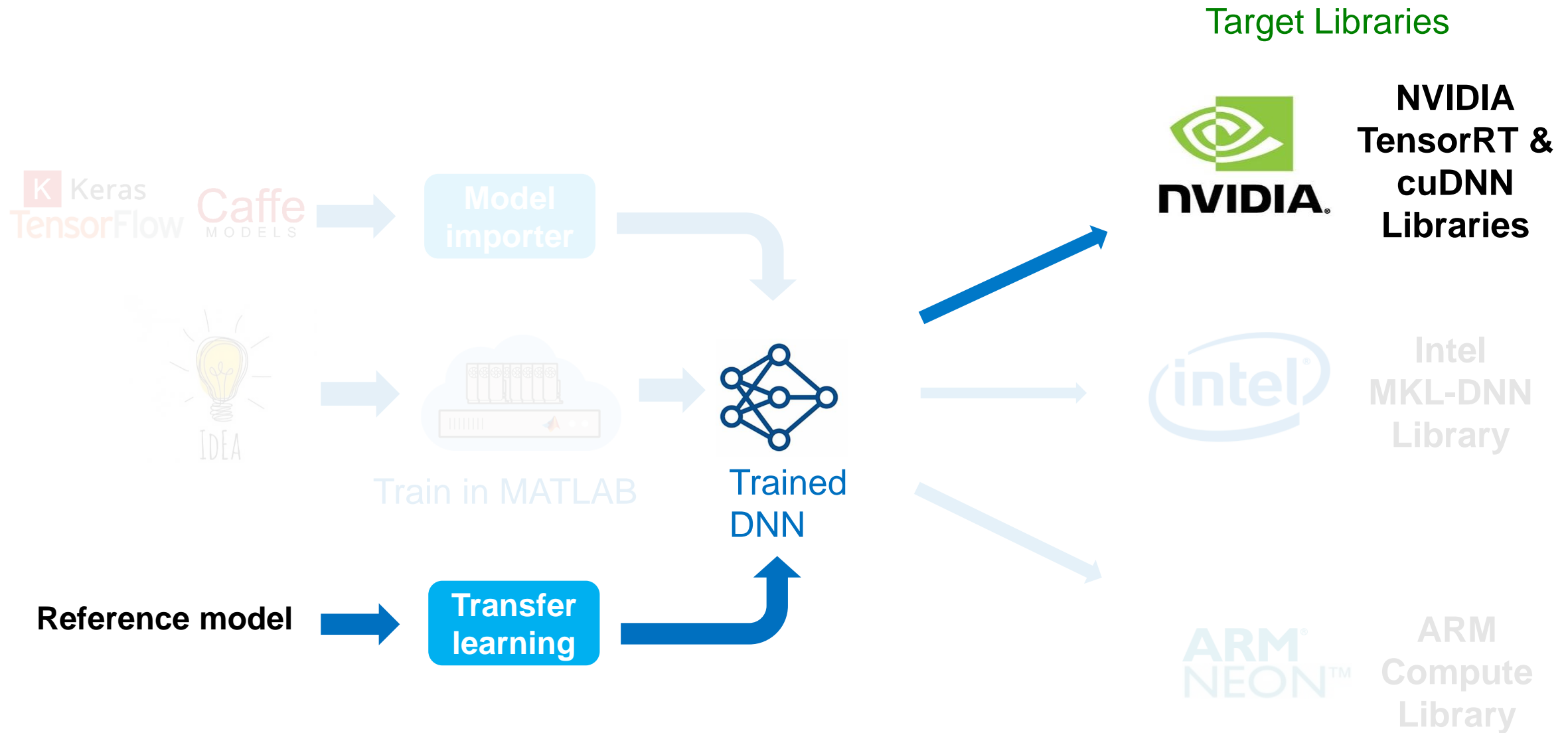
Importer for TensorFlow-Keras Models



# Deep Learning Inference Deployment



# Deep Learning Inference Deployment

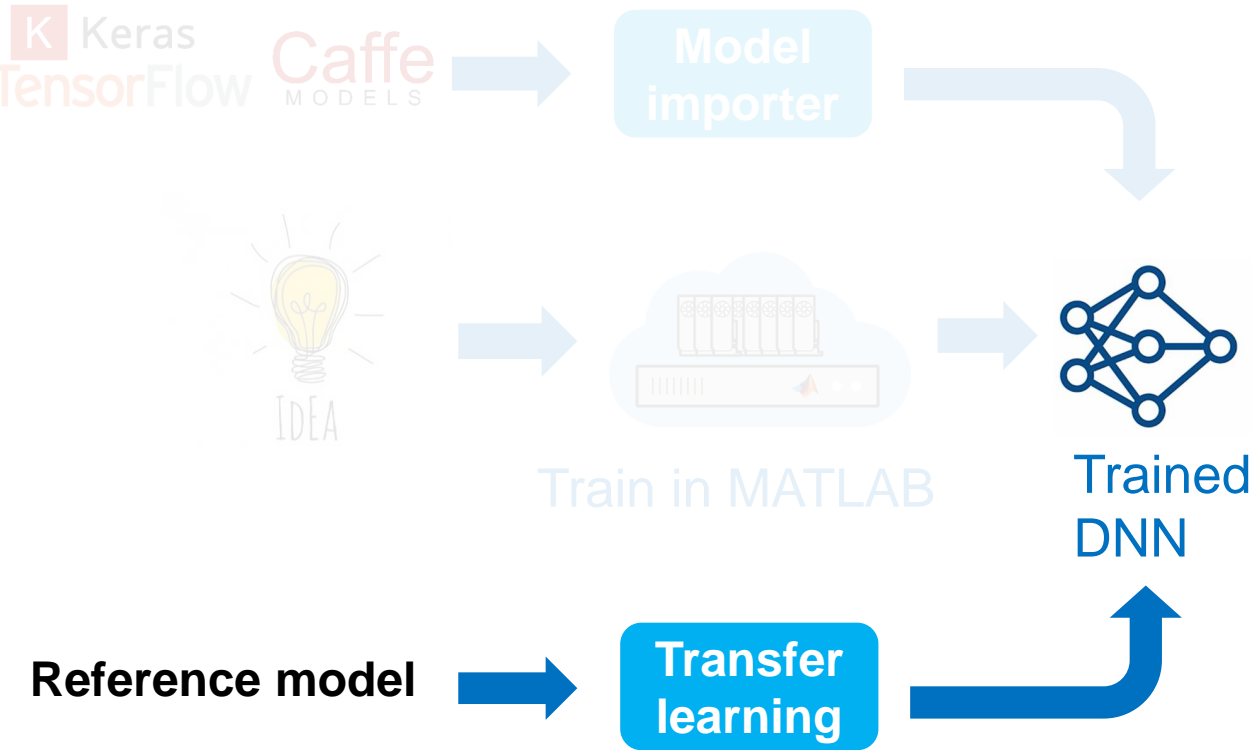


# Deep Learning Inference Deployment

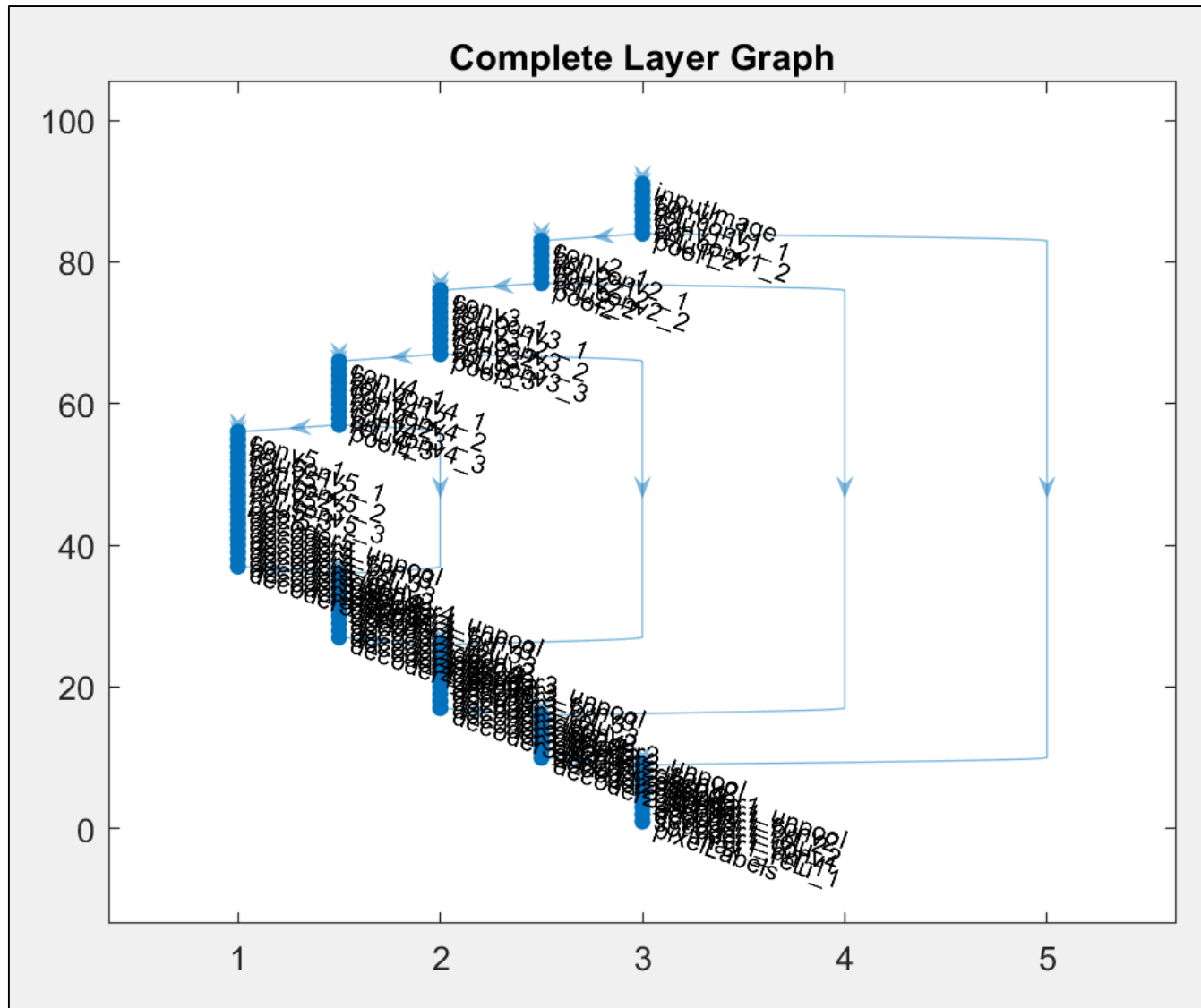
Target Libraries



**NVIDIA  
TensorRT &  
cuDNN  
Libraries**



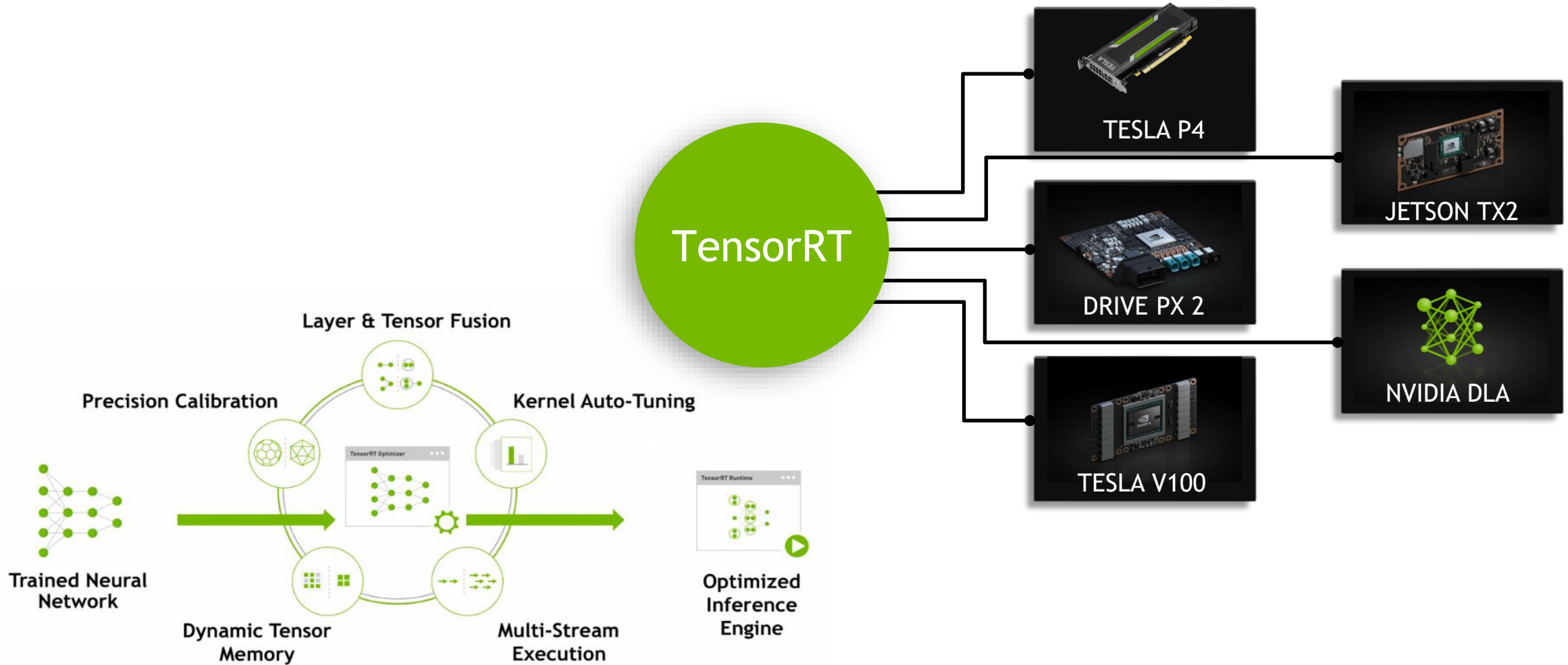
# Layered Architecture for Segnet- Semantic Segmentation



**DAG Network**  
**Total number of layers: 91**

# NVIDIA TensorRT

## PROGRAMMABLE INFERENCE ACCELERATOR



# Performance Summary (VGG-16) on TitanXP

VGG Demo

Fps = 327.79


rubber eraser 11.77%

ballpoint 4.91%

vase 4.49%

perfume 2.20%

teddy 2.11%

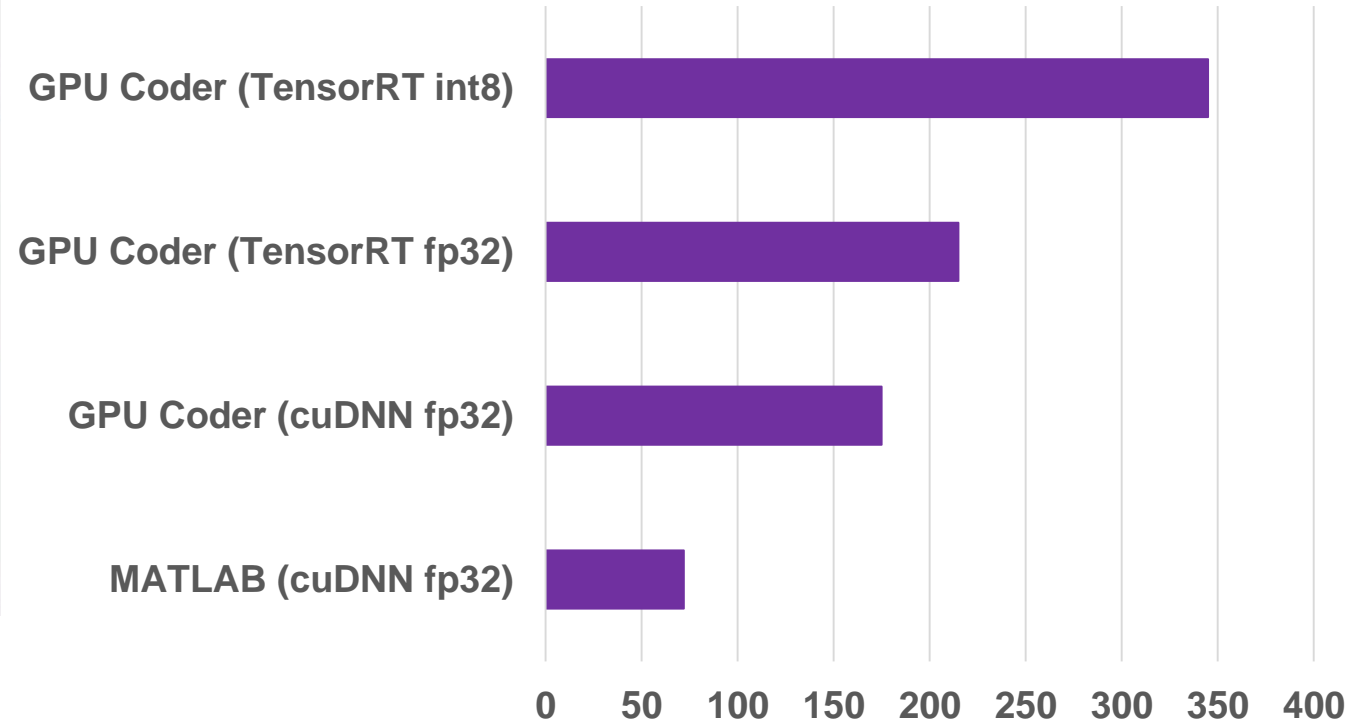


MATLAB on TitanXP: 72 Fps

GPU Coder (cuDNN) on TitanXP: 175 Fps

GPU Coder (TensorRT) on TitanXP: 210 Fps

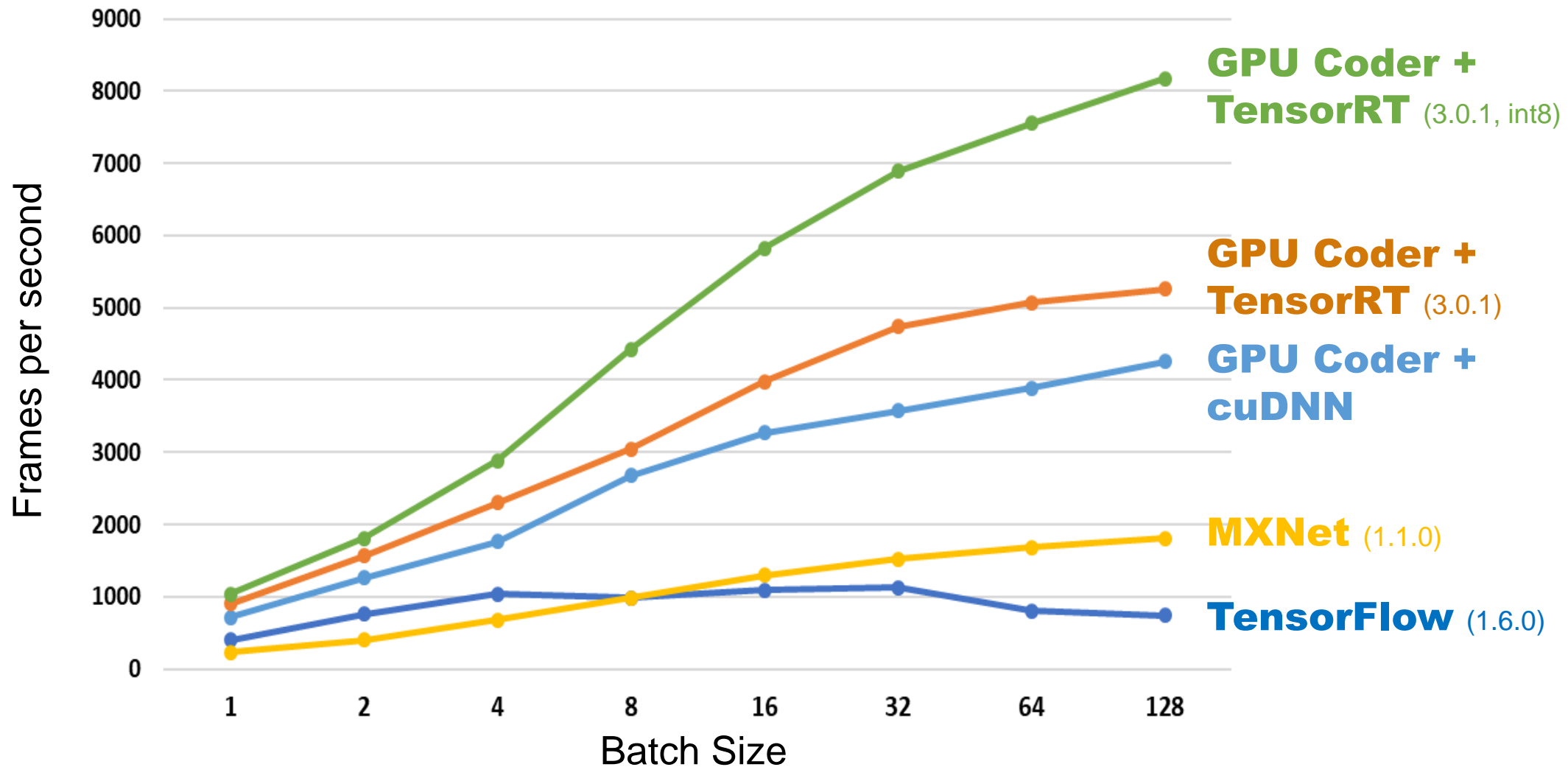
GPU Coder (TensorRT-int8) on TitanXP: 345 Fps



# How Good is Generated Code Performance?

- Performance of CNN inference (Alexnet) on Titan XP GPU
- Performance of CNN inference (Alexnet) on Jetson (Tegra) TX2

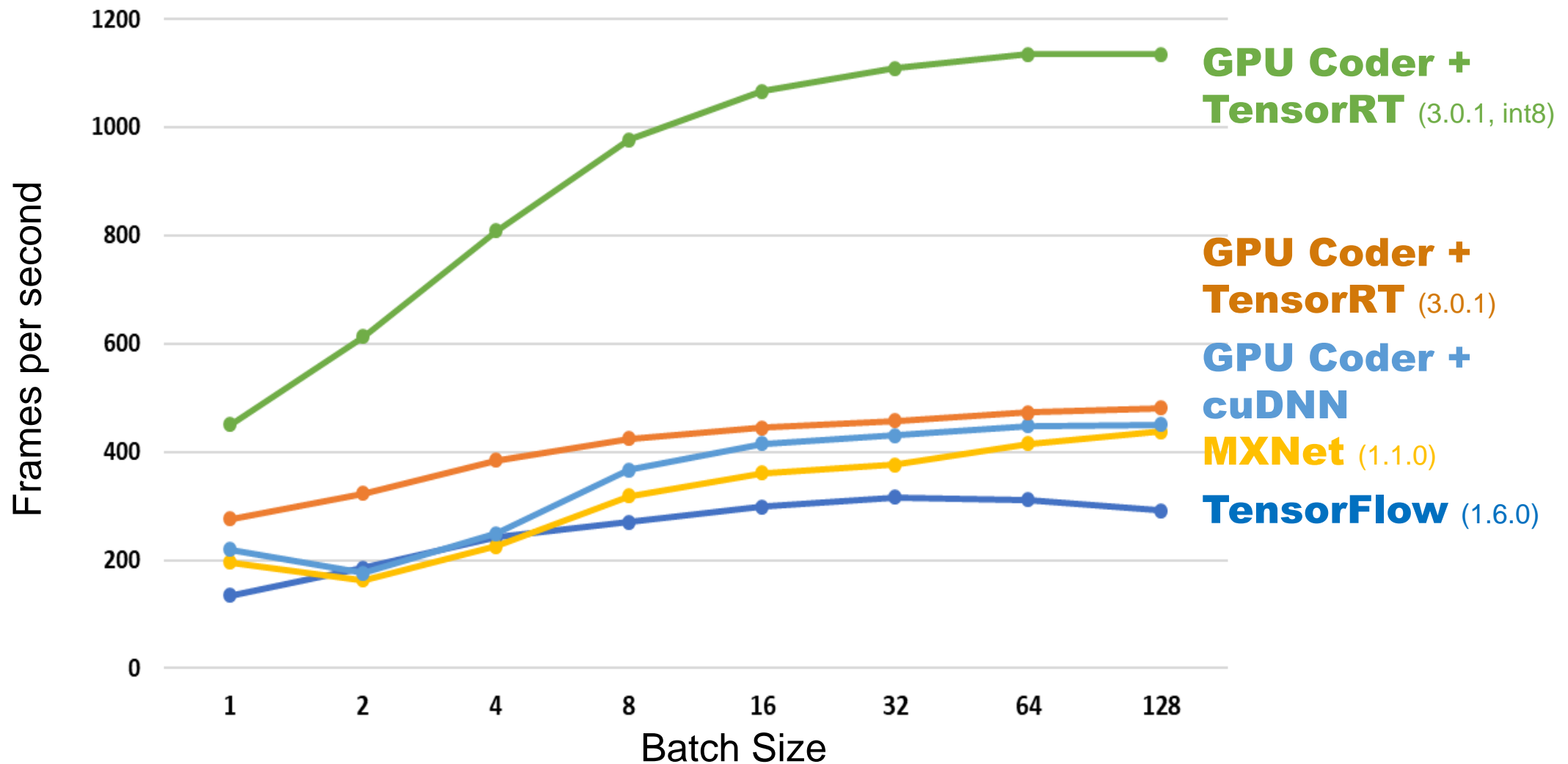
# Alexnet Inference on NVIDIA Titan Xp



CPU	Intel(R) Xeon(R) CPU E5-1650 v4 @ 3.60GHz
GPU	Pascal Titan Xp
cuDNN	v7

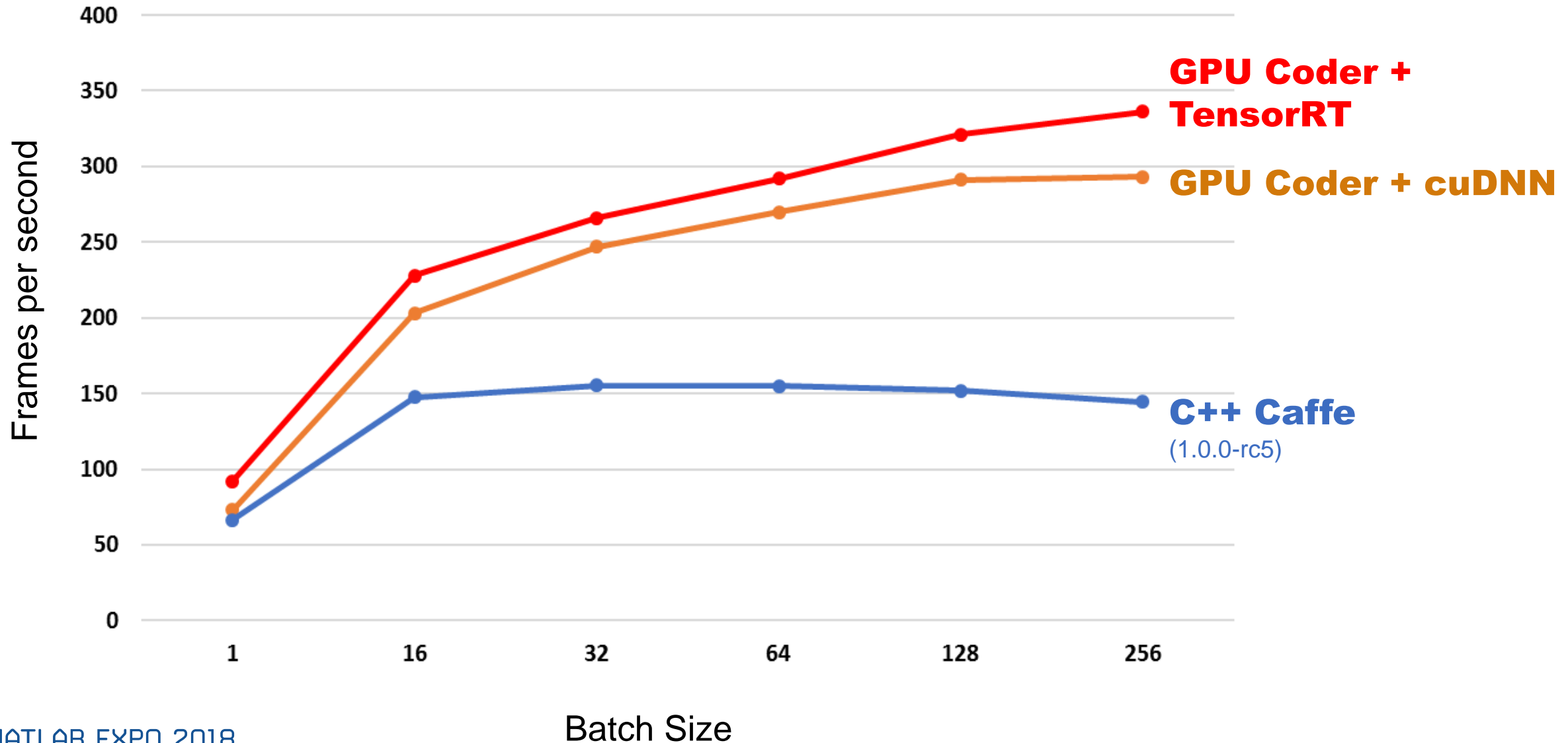


# VGG-16 Inference on NVIDIA Titan Xp



CPU	Intel(R) Xeon(R) CPU E5-1650 v4 @ 3.60GHz
GPU	Pascal Titan Xp
cuDNN	v7

# Alexnet Inference on Jetson TX2: Frame-Rate Performance



# Brief Summary

## **DNN libraries are great for inference, ...**

- GPU coder generates code that takes advantage of:



NVIDIA® CUDA libraries, including cuDNN, and TensorRT



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



ARM® Compute libraries for mobile platforms

# Brief Summary

## DNN libraries are great for inference, ...

- GPU coder generates code that takes advantage of:



NVIDIA® CUDA libraries, including cuDNN, and TensorRT

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



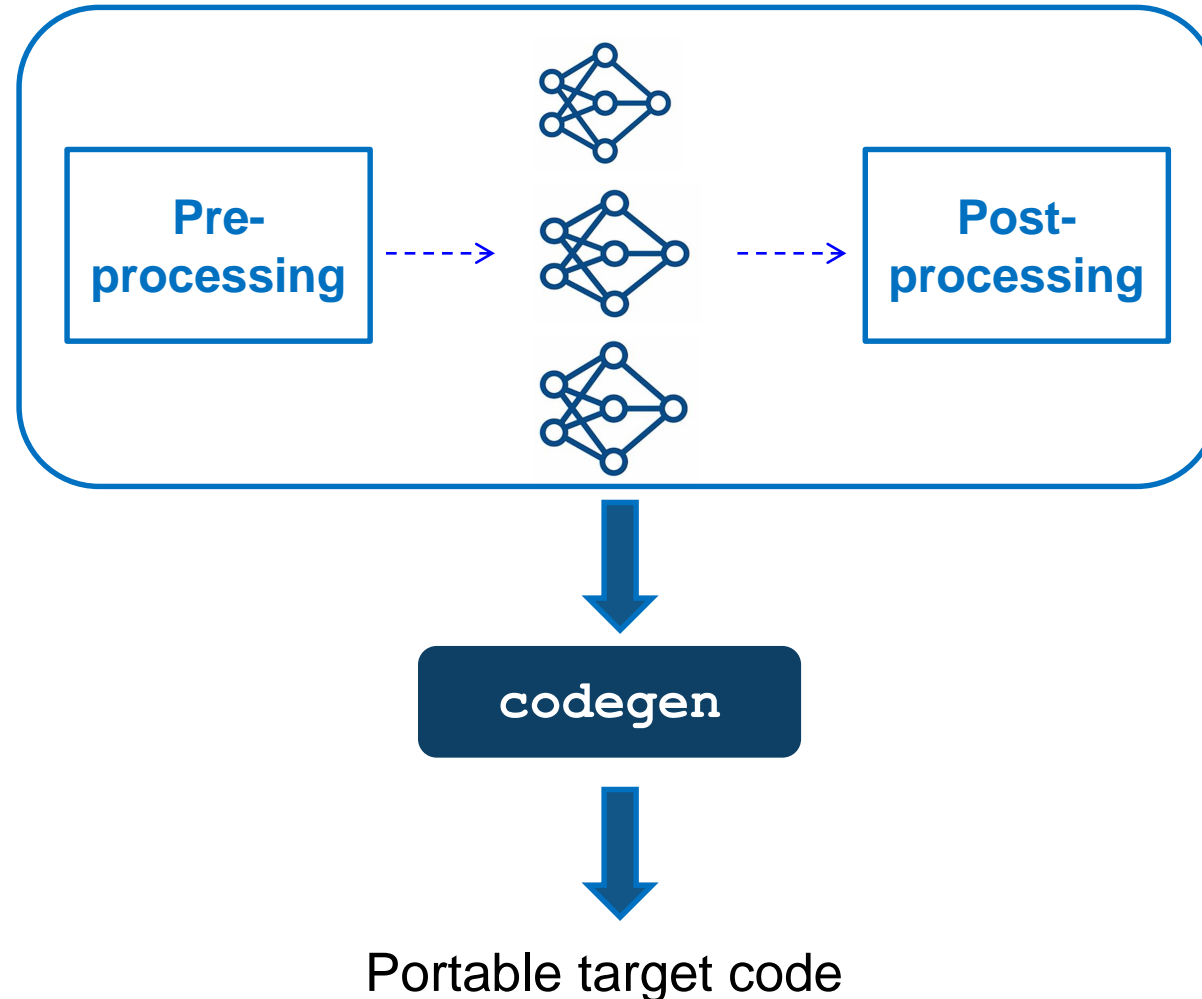
Intel® Math Kernel Library for Deep Neural Networks (MKL-DNN)



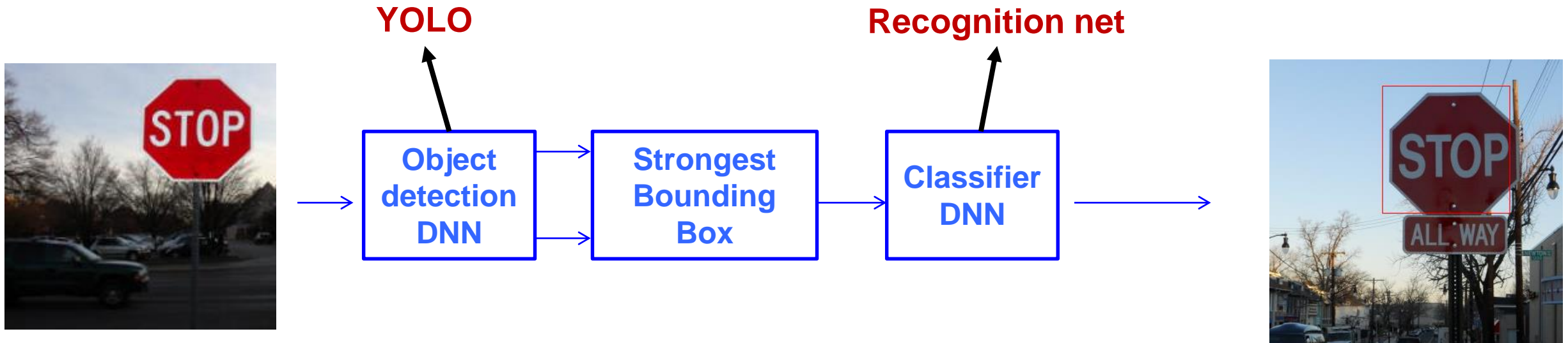
ARM® Compute libraries for mobile platforms

# Deep learning Workflows- Integrated Application Deployment

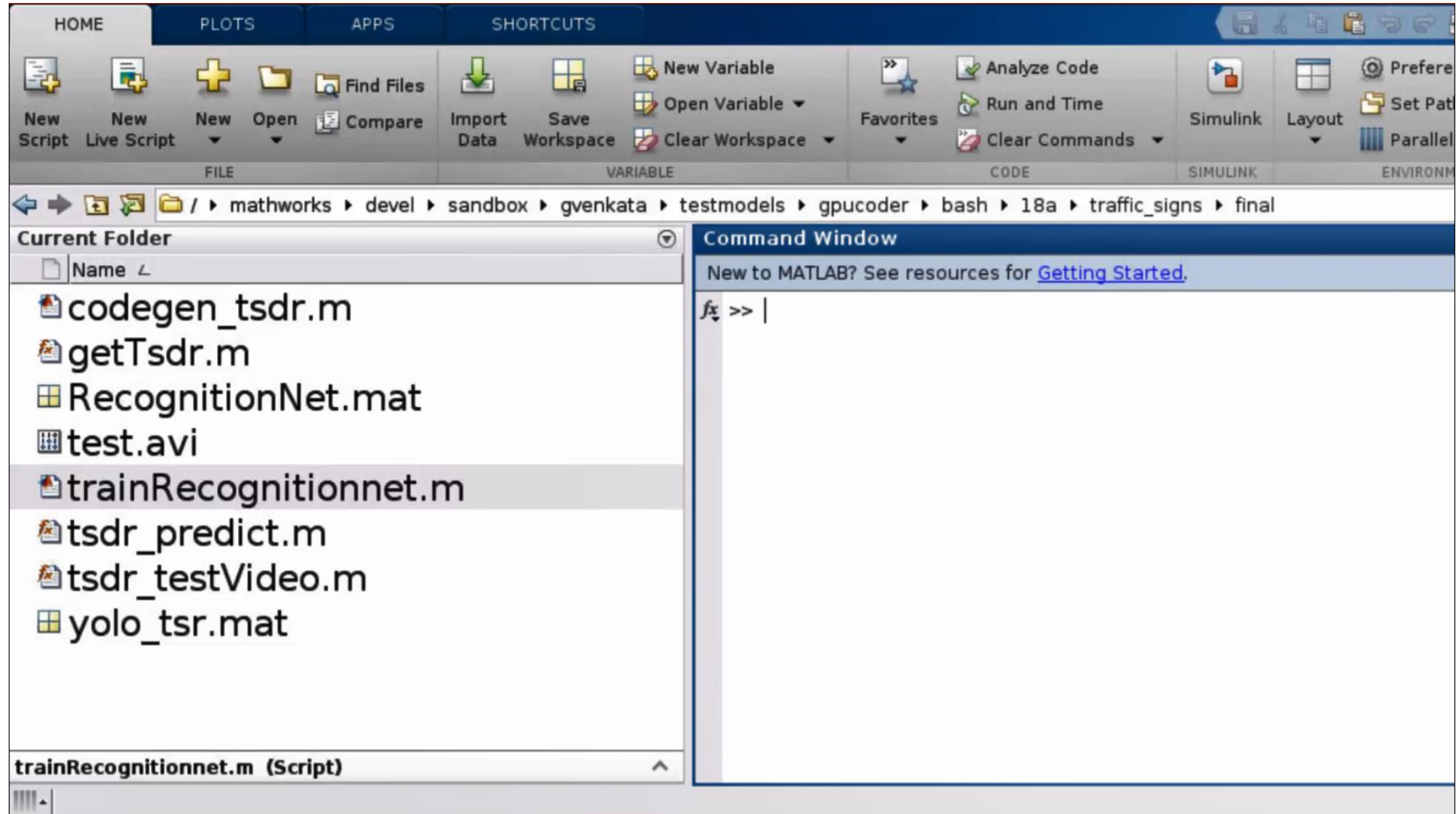
## INTEGRATED APPLICATION DEPLOYMENT



# Traffic sign detection and recognition



# Traffic sign detection and recognition

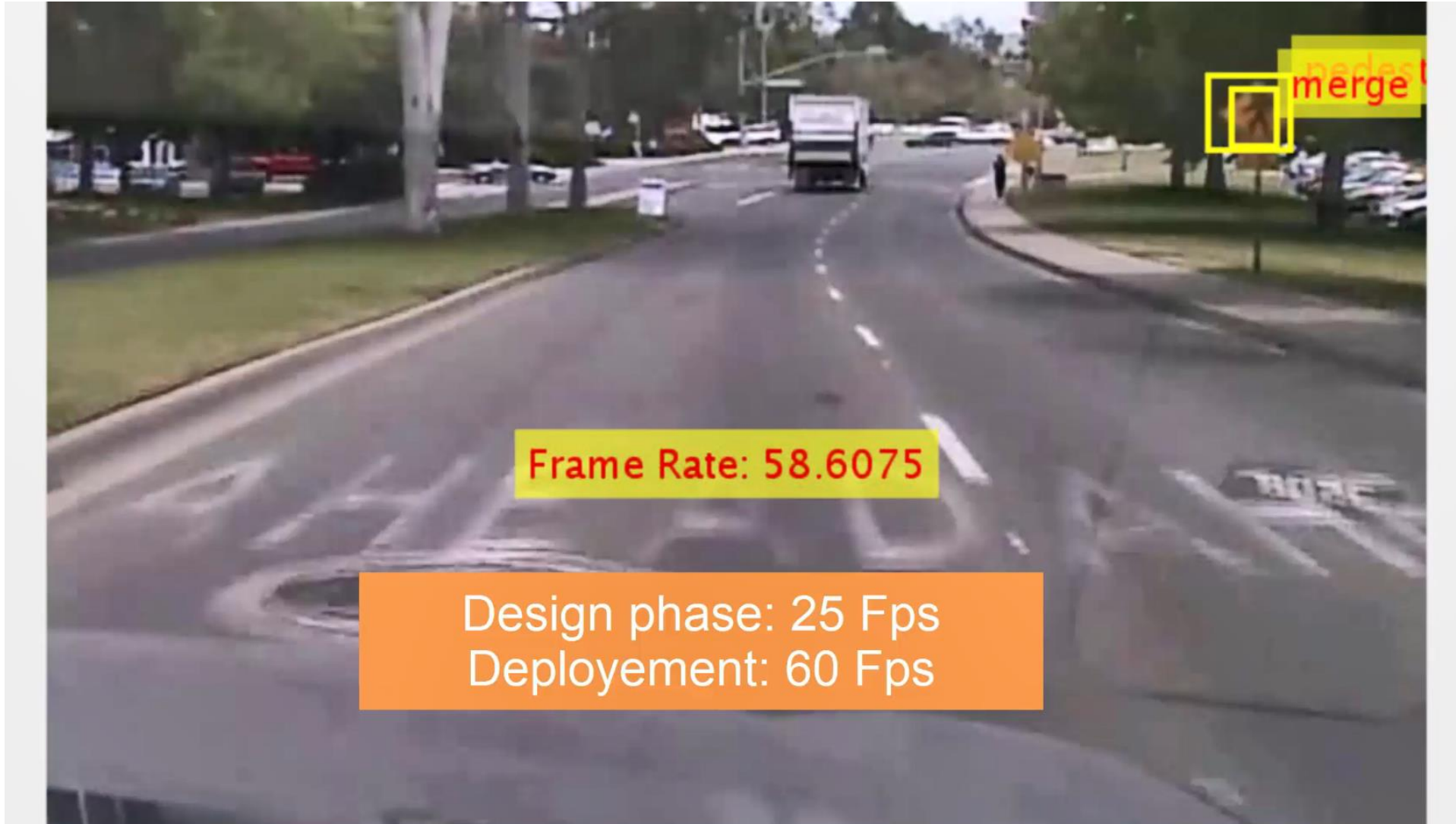


The screenshot displays the MATLAB IDE interface. The top menu bar includes tabs for HOME, PLOTS, APPS, and SHORTCUTS. Below the menu bar, there are several toolbars with icons for file operations (New Script, New Live Script, New, Open, Compare, Import Data, Save Workspace, Clear Workspace), variable management (New Variable, Open Variable, Clear Workspace), code execution (Analyze Code, Run and Time, Clear Commands), and environment settings (Simulink, Layout, Preferences, Set Path, Parallel). The current folder path is shown as `/ > mathworks > devel > sandbox > gvenkata > testmodels > gpucoder > bash > 18a > traffic_signs > final`. The file explorer on the left shows a list of files in the current folder:

- codegen\_tsdr.m
- getTsdr.m
- RecognitionNet.mat
- test.avi
- trainRecognitionnet.m
- tsdr\_predict.m
- tsdr\_testVideo.m
- yolo\_tsr.mat

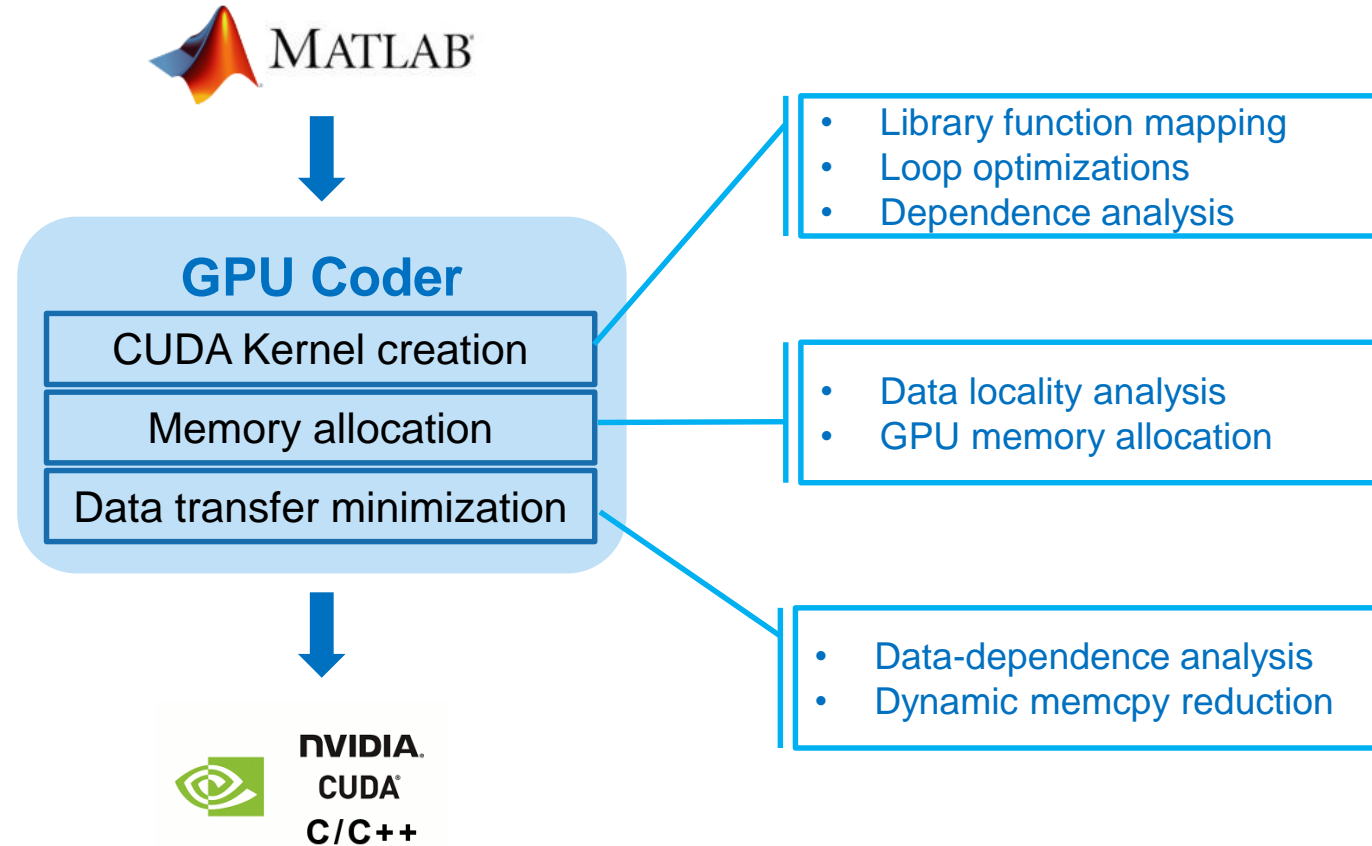
The Command Window on the right shows the prompt `>> |` and a message: "New to MATLAB? See resources for [Getting Started](#)."

# Traffic sign detection and recognition



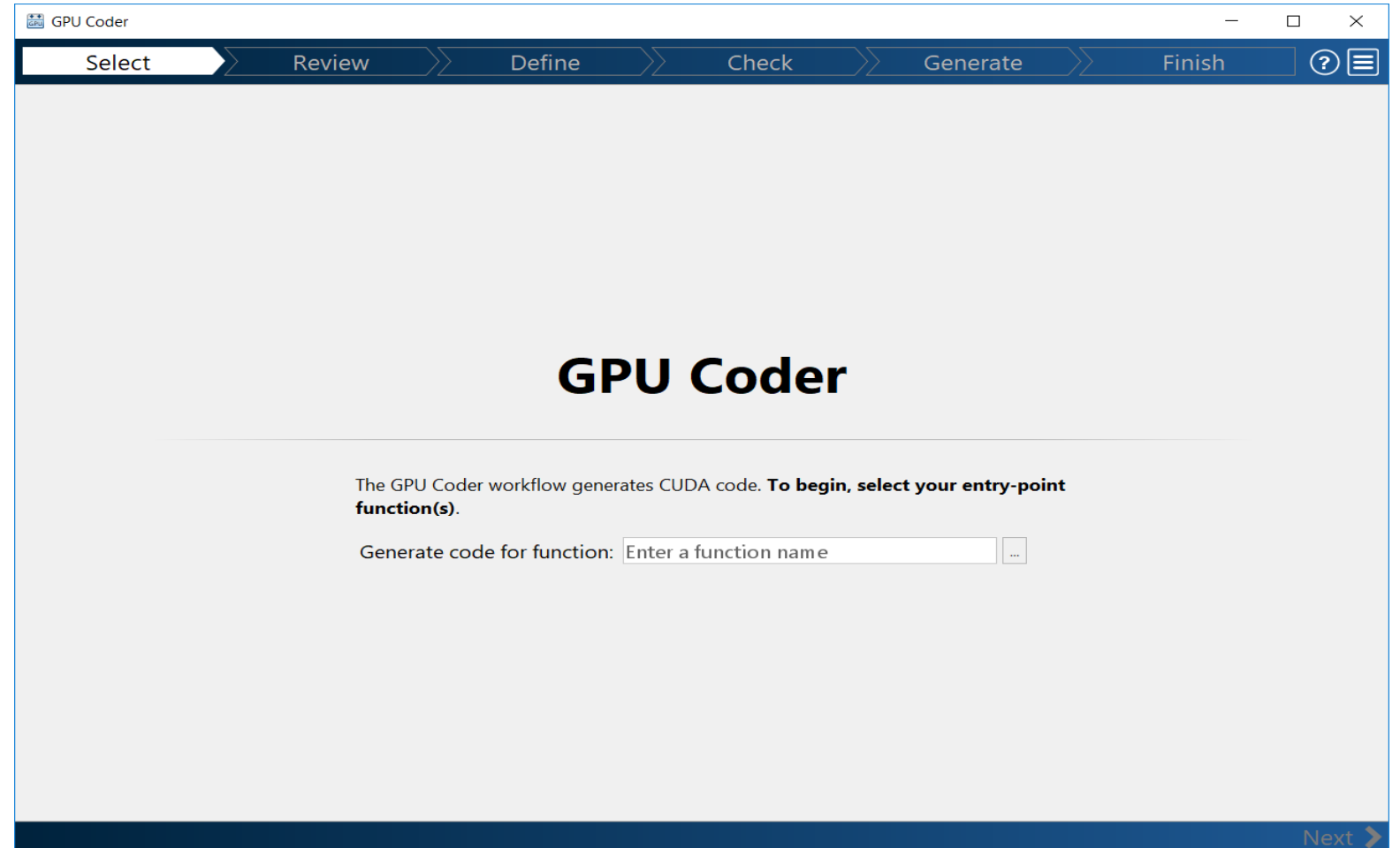


# GPU Coder Helps You Deploy Applications to GPUs Faster



# CUDA Code Generation from GPU Coder app

**Integrated editor  
and simplified  
workflow for code  
generation**



# Summary- GPU Coder

```

%_MEXFILE_ Preprocess MATLAB file
%
% This function precompiles MATLAB code into an MEX file.
% It is used to generate MEX files for MATLAB.
%
% Syntax
%
% mex('filename.m')
% mex('filename.m', 'options')
%
% Example
%
% mex('mycode.m')
%
% See also
%
% mexCompileOptions, mexOptions, mexLinkOptions, mexLinkOptions, mexLinkOptions
%
% Copyright 2012-2013 MathWorks, Inc.

```

MATLAB algorithm  
(functional reference)

GPU Coder

Build type

Call CUDA  
from MATLAB  
directly

.mex

Call CUDA from  
(C++) hand-  
coded main()

.lib

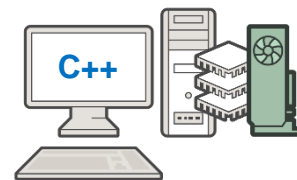
Call CUDA from (C++)  
hand-coded main().

Cross-compiled  
.lib

Desktop  
GPU

Desktop  
GPU

Embedded GPU



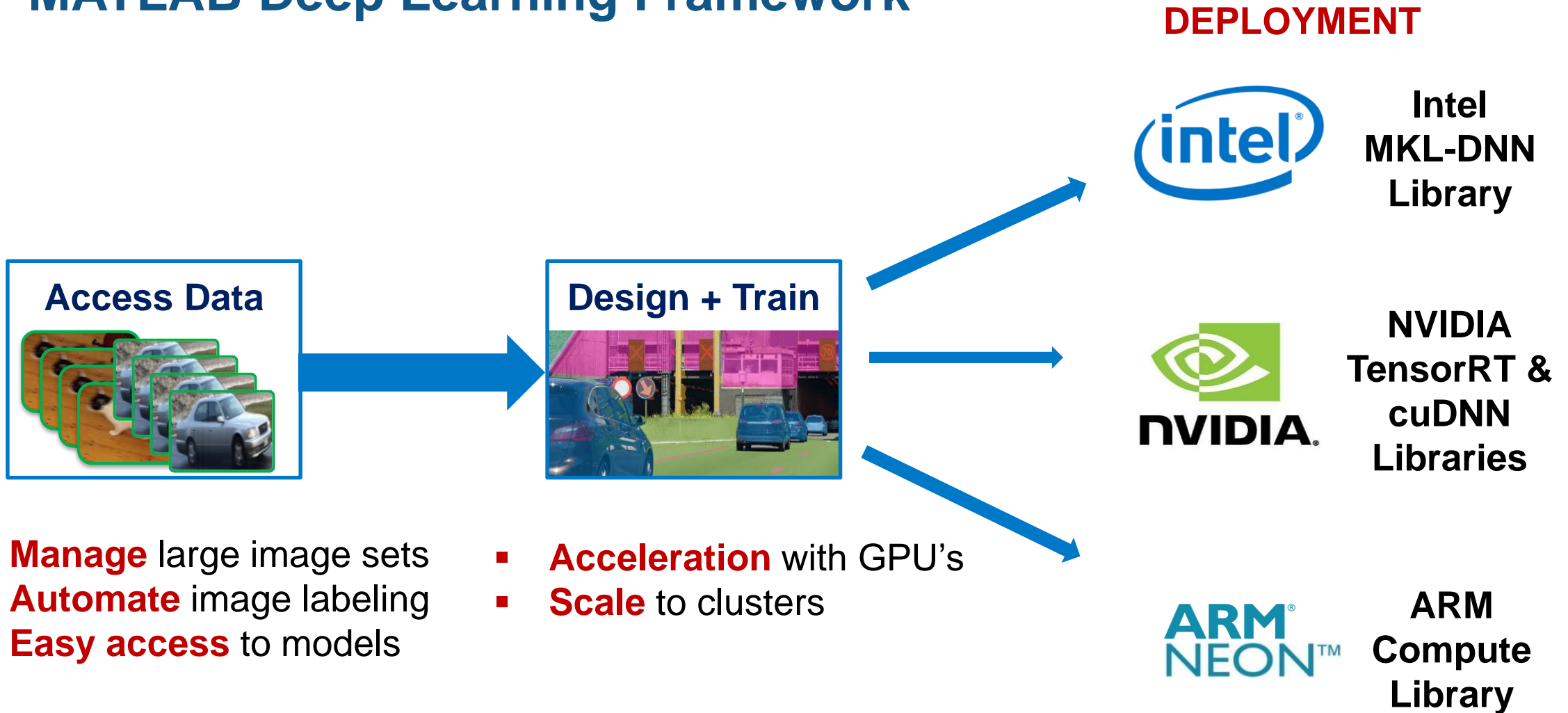
1 Functional test

2 Deployment  
unit-test

3 Deployment  
integration-test

4 Real-time test

# MATLAB Deep Learning Framework



- **Manage** large image sets
  - **Automate** image labeling
  - **Easy access** to models
- **Acceleration** with GPU's
  - **Scale** to clusters

## Speaker Details

Email: [Rishu.Gupta@mathworks.in](mailto:Rishu.Gupta@mathworks.in)

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

## Contact MathWorks India

Products/Training Enquiry Booth

Call: 080-6632-5749

Email: [info@mathworks.in](mailto:info@mathworks.in)

- **Share your experience with MATLAB & Simulink on Social Media**

- Use #MATLABEXPO
- I use #MATLAB because..... Attending #MATLABEXPO
- Examples
  - I use #MATLAB because it helps me be a data scientist! Attending #MATLABEXPO
  - Learning new capabilities in #MATLAB and #Simulink at #MATLABEXPO.

- **Share your session feedback:**

Please fill in your feedback for this session in the feedback form