# MATLAB EXPO 2018

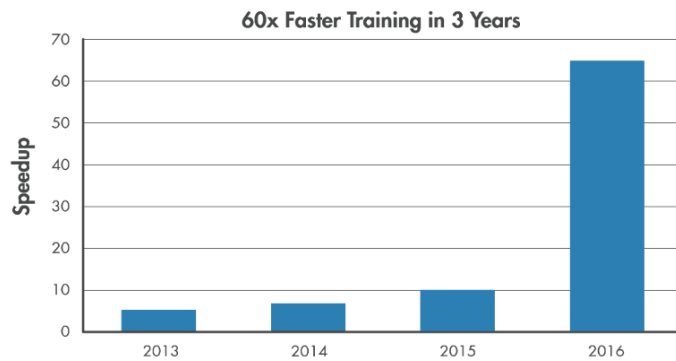## Deploying Deep Learning Networks to Embedded GPUs and CPUs

Pierre Nowodzienski
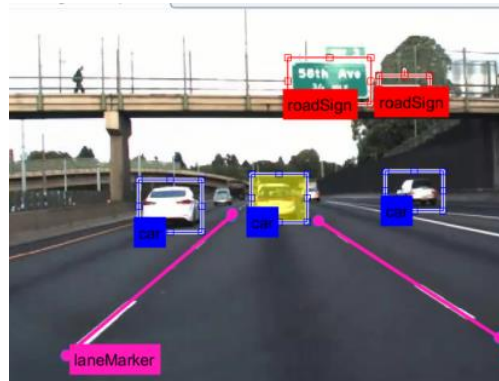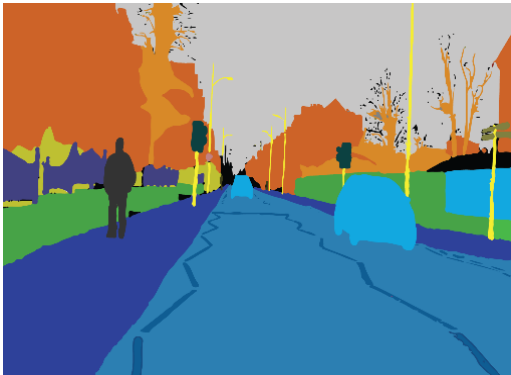
# Deep Learning enablers

## Increased GPU acceleration



60x Faster Training in 3 Years

## Labeled public datasets



## World-class models to be leveraged

**AlexNet**
PRETRAINED MODEL

**VGG-16**
PRETRAINED MODEL

**ResNet**
PRETRAINED MODEL

**Caffe**
M O D E L S

**GoogLeNet**
PRETRAINED MODEL

**TensorFlow/Keras**
M O D E L S



Human Accuracy

machine learning          deep learning

# Deep Learning Applications:
# Image classification, speech recognition, autonomous driving, etc…
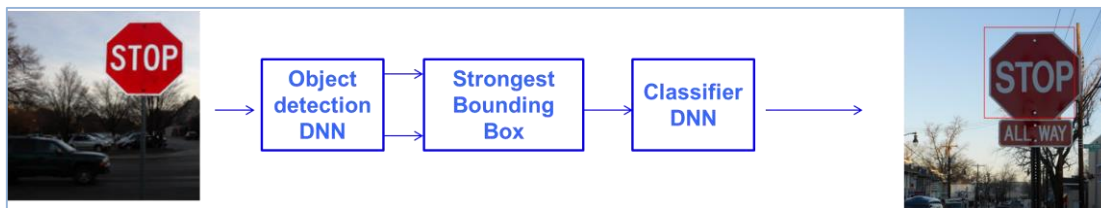

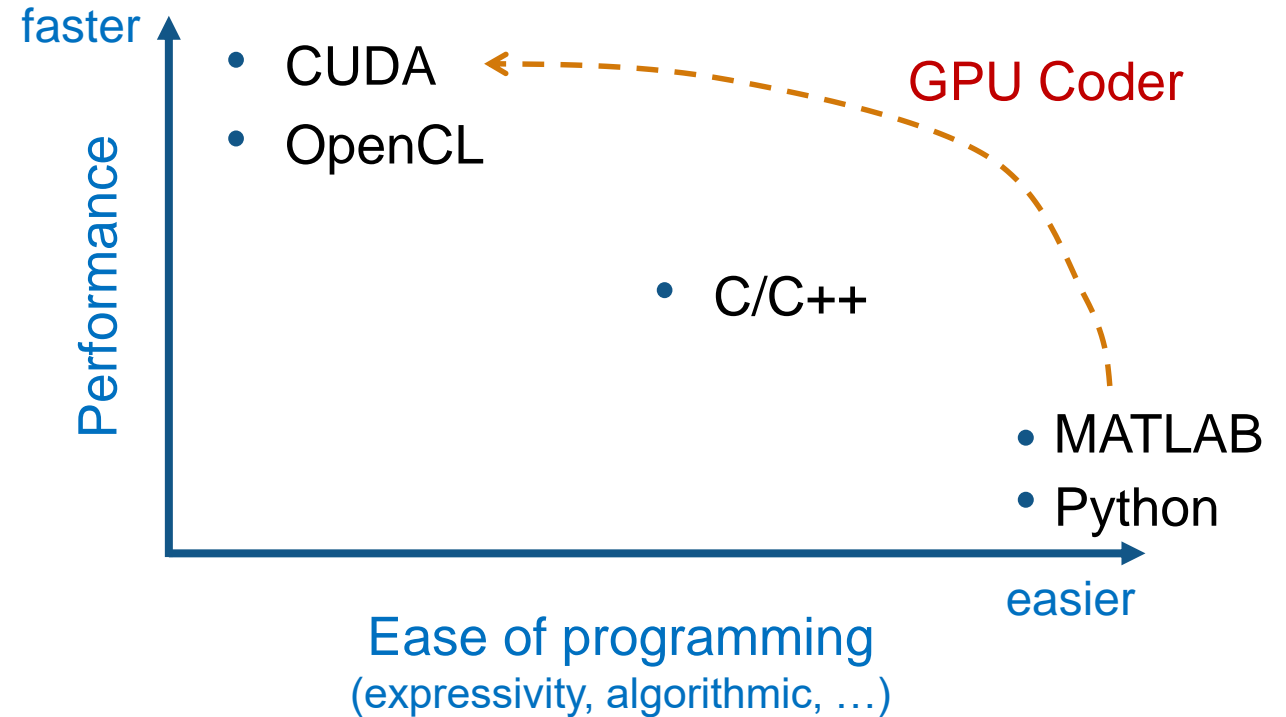
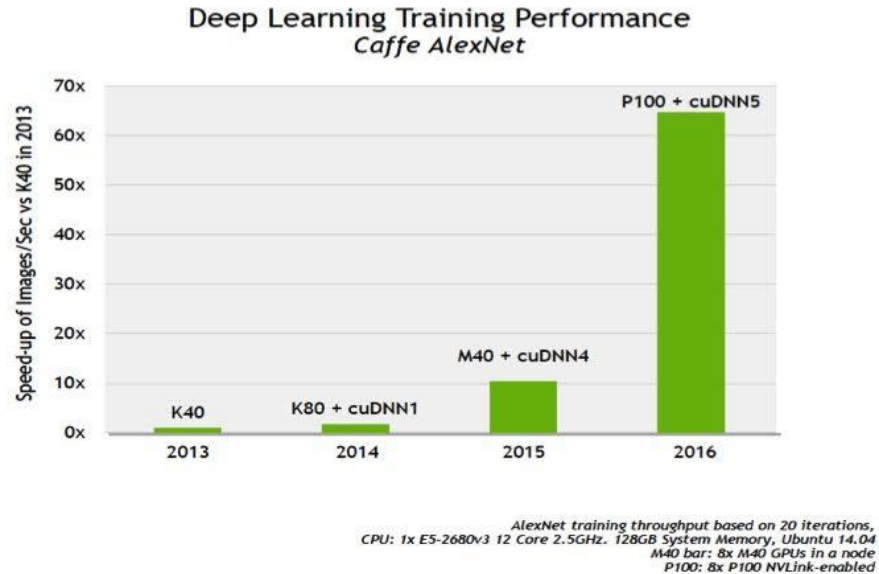*Detection of cars and road in autonomous driving systems*



*Rain Detection and Removal[1]*

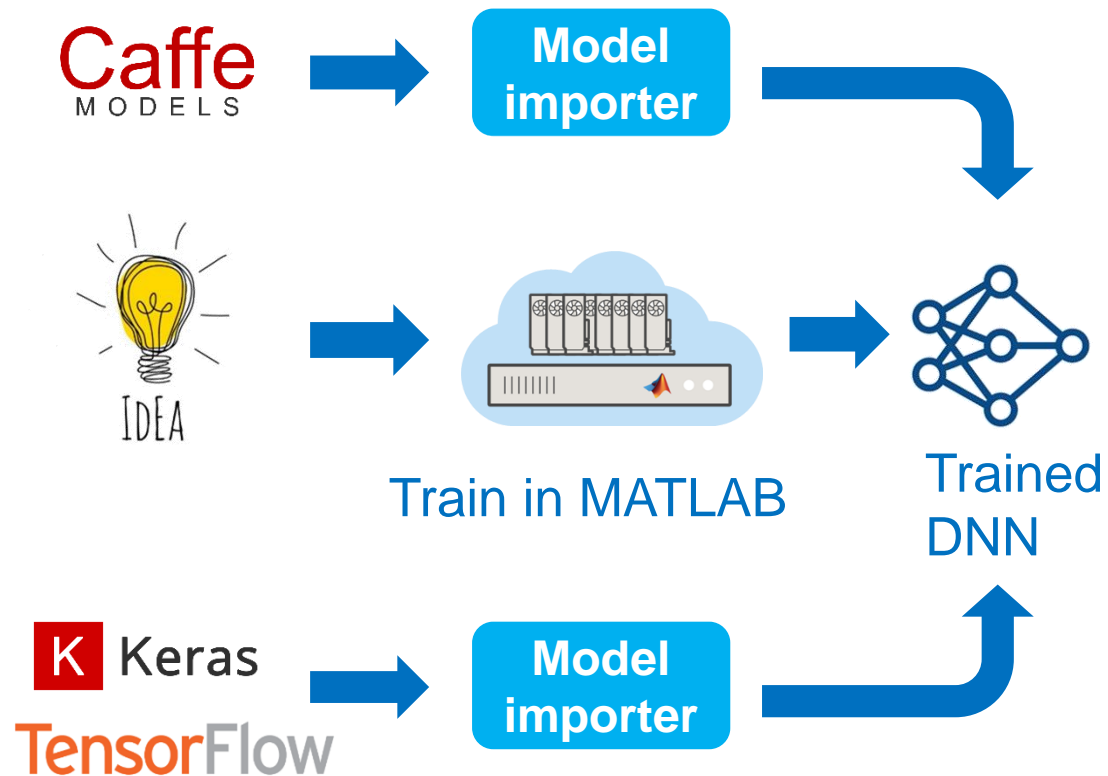1. *Deep Joint Rain Detection and Removal from a Single Image*



**Object detection DNN** → **Strongest Bounding Box** → **Classifier DNN**

Traffic Sign Recognition

# GPUs and CUDA programming



Deep Learning Training Performance
*Caffe AlexNet*

AlexNet training throughput based on 20 iterations,
CPU: 1x E5-2680v3 12 Core 2.5GHz. 128GB System Memory, Ubuntu 14.04
M40 bar: 8x M40 GPUs in a node
P100: 8x P100 NVLink-enabled

faster

Performance

- CUDA
- OpenCL

GPU Coder

- C/C++

- MATLAB
- Python

easier

Ease of programming
(expressivity, algorithmic, …)

**GPUs are "hardware on steroids",     …     but, programming them is hard**

# Deep learning workflow in MATLAB

**Deep Neural Network**
**Design + Training**



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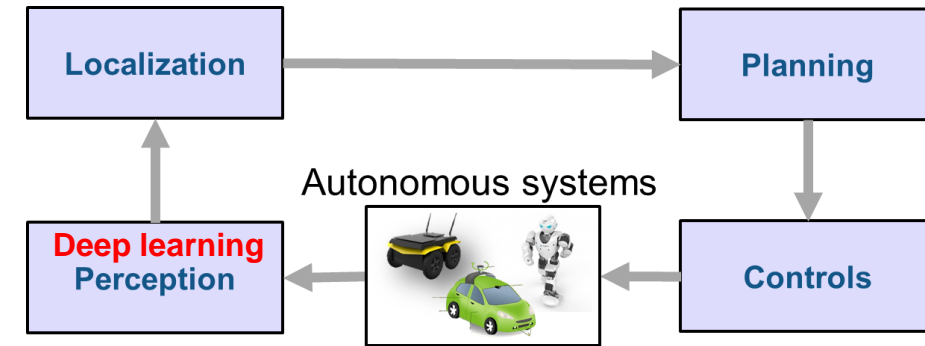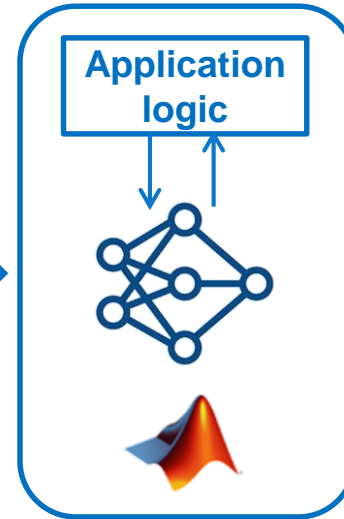  - **Easy access** to models

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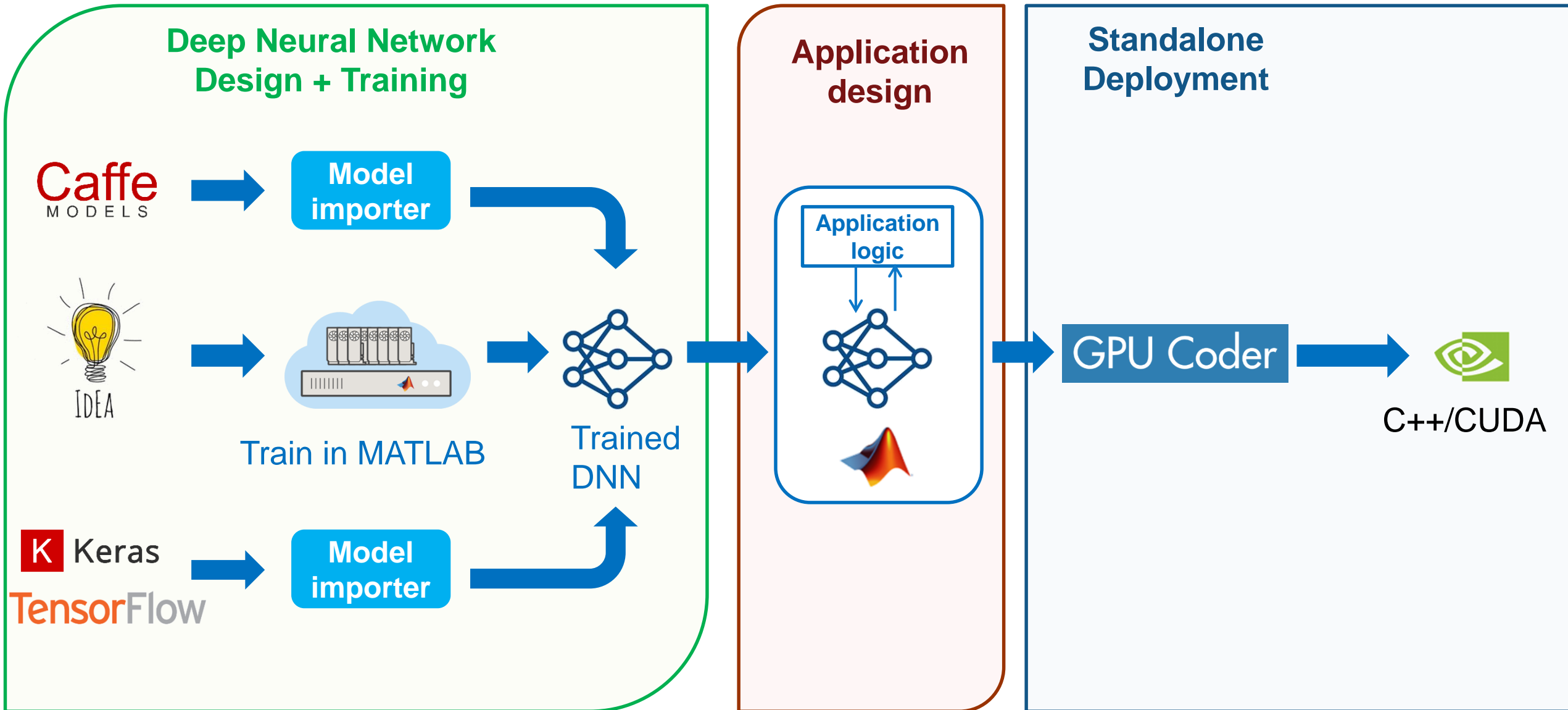  - **Scale** to clusters

# Deep learning workflow in MATLAB

# Deep learning workflow in MATLAB

# GPUs and CUDA

# Challenges of Programming in CUDA for GPUs

- Learning to program in CUDA
  - Need to rewrite algorithms for parallel processing paradigm

- Creating CUDA kernels
  - Need to analyze algorithms to create CUDA kernels that maximize parallel processing

- Allocating memory
  - Need to deal with memory allocation on both CPU and GPU memory spaces

- Minimizing data transfers
  - Need to minimize while ensuring required data transfers are done at the appropriate parts of your algorithm

# GPU Coder Helps You Deploy to GPUs Faster

**MATLAB**

**GPU Coder**

- Library function mapping
- Loop optimizations
- Dependence analysis

CUDA Kernel creation

- Data locality analysis
- GPU memory allocation

Memory allocation

Data transfer minimization

- Data-dependence analysis
- Dynamic memcpy reduction

**NVIDIA. CUDA® C/C++**

# GPU Coder speeds up MATLAB for Image Processing and Computer Vision



Fog removal

**5x speedup**

Frangi filter

**3x speedup**

Distance transform

**8x speedup**

Stereo disparity

**50x speedup**

Ray tracing

**18x speedup**

SURF feature extraction

**700x speedup**

# GPU Coder speeds up MATLAB at least 2x for inference

MATLAB 18a on TitanXP GPU - Linux



GPU Coder
MATLAB

*Single image prediction using Intel® Xeon® CPU - 3.6 GHz, NVIDIA libraries: CUDA8 - cuDNN 7*
*TensorFlow 1.6.0, MXNet 1.1.0, MATLAB 18a*

# With GPU Coder, MATLAB is faster than other frameworks

Single Image Prediction (TitanXP GPU, Linux)

**AlexNet**  **ResNet-50**  **VGG-16**

TensorFlow
MXNet
GPU Coder

*Single image prediction using Intel® Xeon® CPU - 3.6 GHz, NVIDIA libraries: CUDA8 - cuDNN 7*
*TensorFlow 1.6.0, MXNet 1.1.0, MATLAB 18a*

# Embedded GPU Benchmarking: Jetson TX2

# Algorithm Design to Embedded Deployment Workflow



MATLAB algorithm
(functional reference)

GPU Coder

Build type

Call CUDA
from MATLAB
directly

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

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

.mex

.lib

Cross-compiled
.lib

Desktop
GPU

Desktop
GPU

Embedded GPU

C++

C++

**1** Functional test   **2** Deployment unit-test   **3** Deployment integration-test   **4** Real-time test

# Demo: Alexnet Deployment with 'mex' Code Generation

# Algorithm Design to Embedded Deployment on Tegra GPU



MATLAB algorithm
(functional reference)

GPU Coder

Build type

Call CUDA from MATLAB directly

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

Call CUDA from (C++) hand-coded main(). Cross-compiled on host with Linaro toolchain

.mex

.lib

Cross-compiled .lib

Tesla GPU

Tesla GPU

Tegra GPU

C++

C++

**1** Functional test

**2** Deployment unit-test

**3** Deployment integration-test

**4** Real-time test
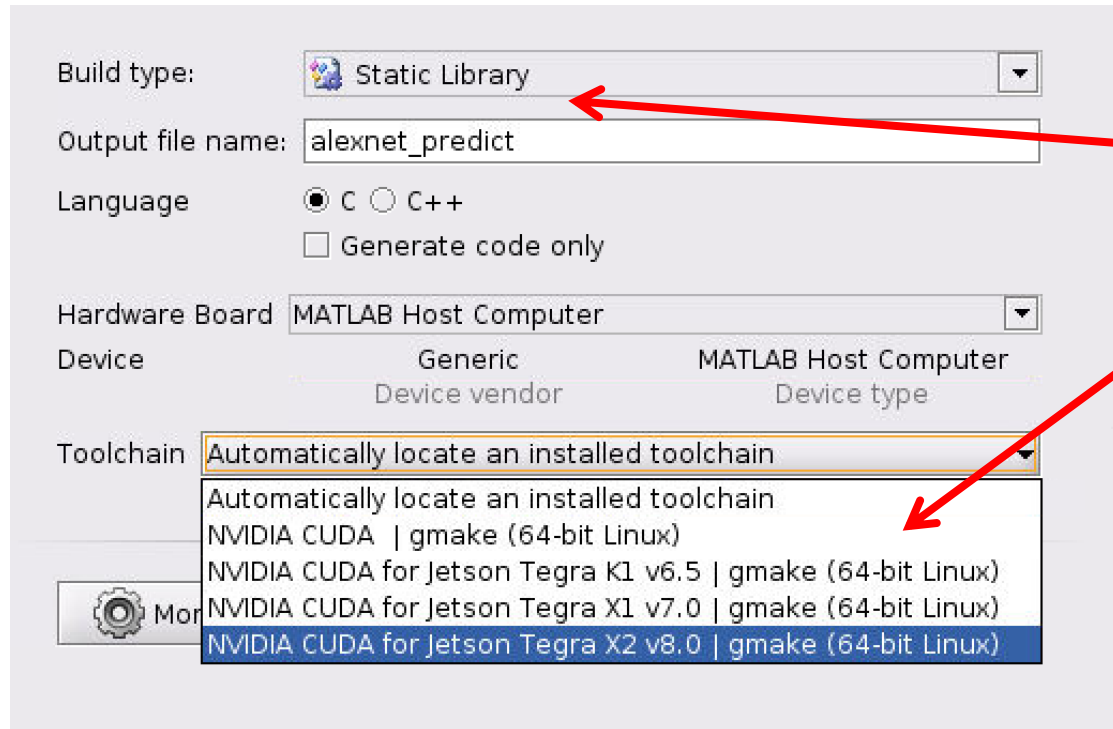
(Test in MATLAB on host)

(Test generated code in MATLAB on host + GPU)

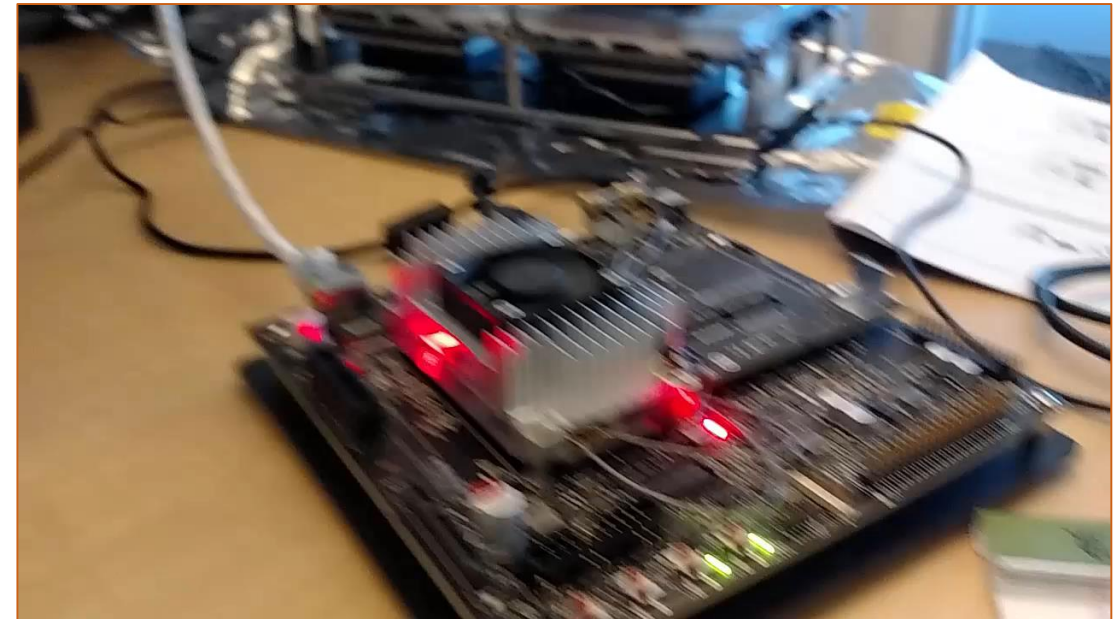(Test generated code within C/C++ app on host + GPU)

(Test generated code within C/C++ app on Tegra target)

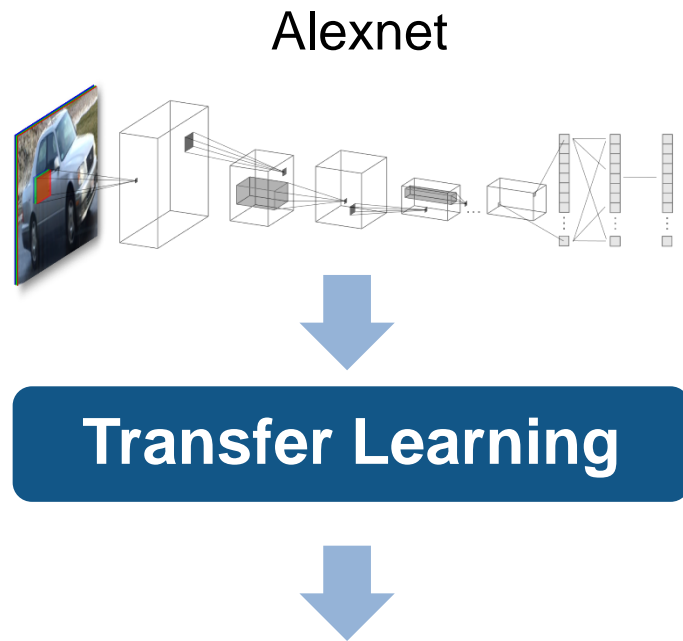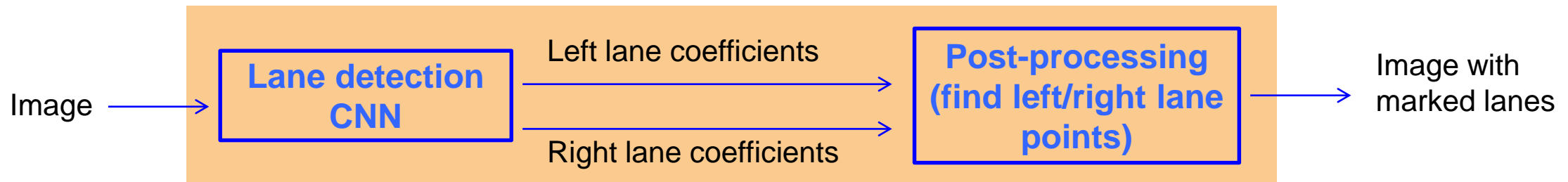# Alexnet Deployment to Tegra: Cross-Compiled with 'lib'



Two small changes
1. Change build-type to 'lib'

2. Select cross-compile toolchain

# Deep learning workflow in MATLAB