

# MATLAB EXPO 2018

## Deploying Deep Learning Networks to Embedded GPUs and CPUs

Rainer Mümmler



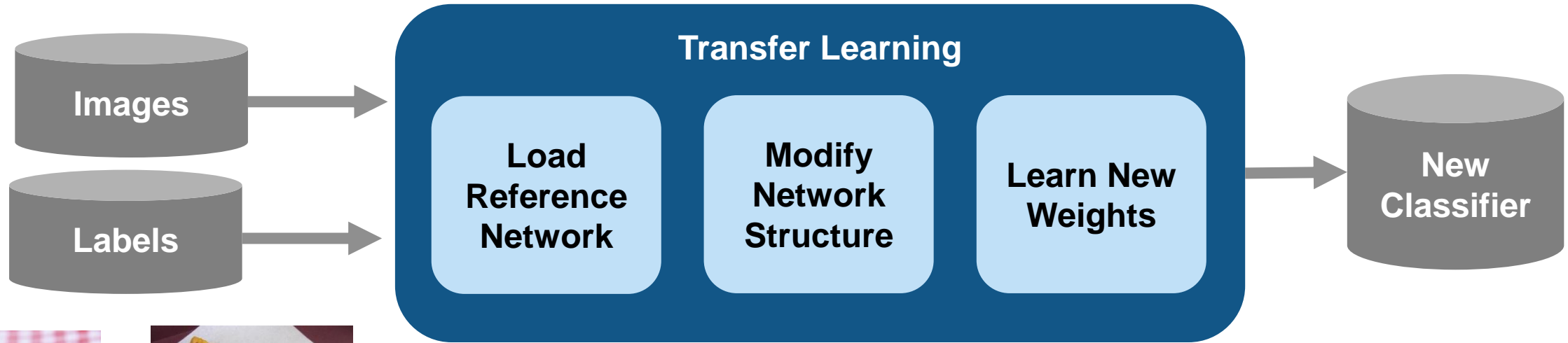
# MATLAB Deep Learning Framework



- **Manage** large image sets
- **Automate** image labeling
- **Easy access** to models
- **Acceleration** with GPU's
- **Scale** to clusters
- **Automate compilation to GPUs and CPUs using GPU Coder**

# Design Deep Learning & Vision Algorithms

## Transfer Learning Workflow



**Labels:** Hot dogs, Pizzas, Ice cream, Chocolate cake, French fries

# Example: Transfer Learning in MATLAB

## Set up training dataset

```
%% set up training dataset
cifarFolder = 'cifar10Train';
categories = {'Cars', 'Trucks', 'BigTrucks', 'Suvs', 'Vans'};
imds = imageDatastore(fullfile(cifarFolder, categories), ...
    'LabelSource', 'foldernames');

imds = splitEachLabel(imds, 500, 'randomize'); % we only need 500 images per class
imds.ReadFcn = @readFunctionTrain;
```

## Load Reference Network

```
%% load reference network
net = alexnet;
layers = net.Layers;
```

## Modify Network Structure

```
%% modify network
layers = layers(1:end-3);

layers(end+1) = fullyConnectedLayer(64, 'Name', 'special_2');
layers(end+1) = reluLayer;
layers(end+1) = fullyConnectedLayer(5, 'Name', 'fc8_2 ');
layers(end+1) = softmaxLayer;
layers(end+1) = classificationLayer();
```

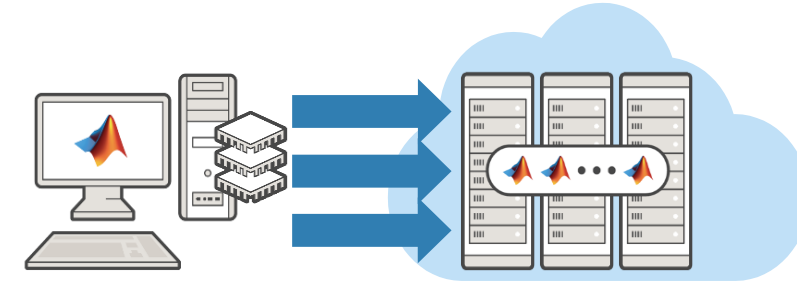
## Learn New Weights

```
%% train!
options = trainingOptions('sgdm', ...
    'LearnRateSchedule', 'none', ...
    'InitialLearnRate', .0001, ...
    'MaxEpochs', 20, ...
    'MiniBatchSize', 128);

myConvnet = trainNetwork(imds, layers, options);
```

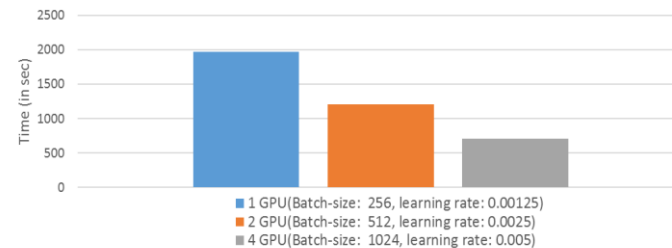
# Scaling Up Model Training Performance

```
'ExecutionEnvironment', 'parallel' );
```



**Training on the AWS (EC2)**

```
'ExecutionEnvironment', 'multi-gpu' );
```

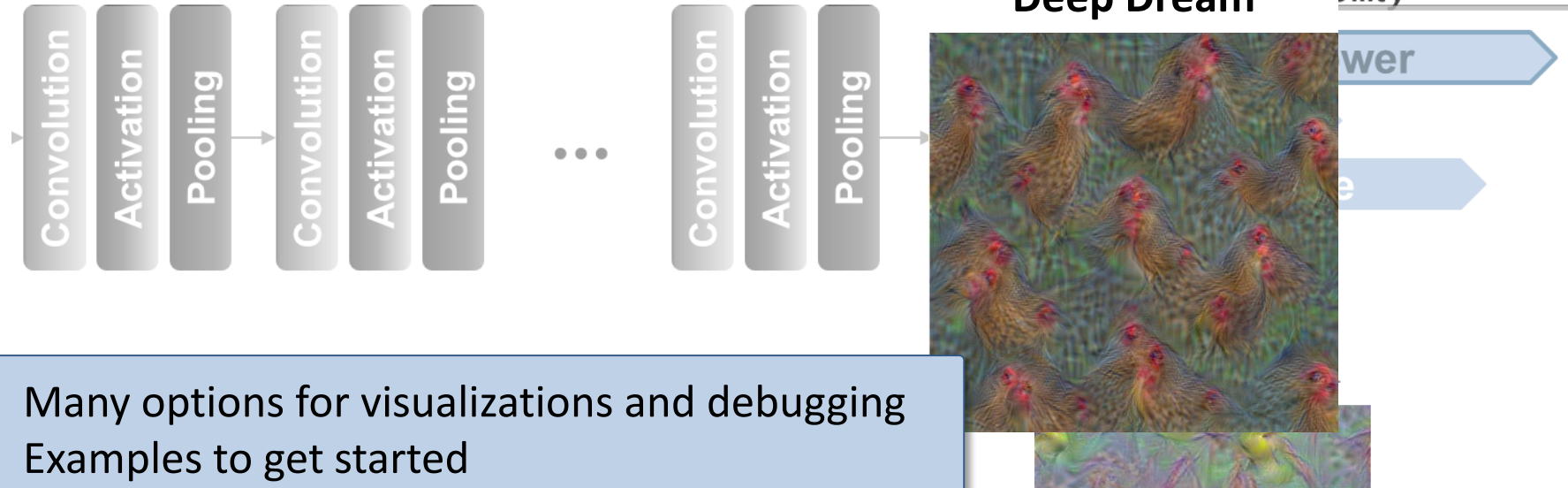
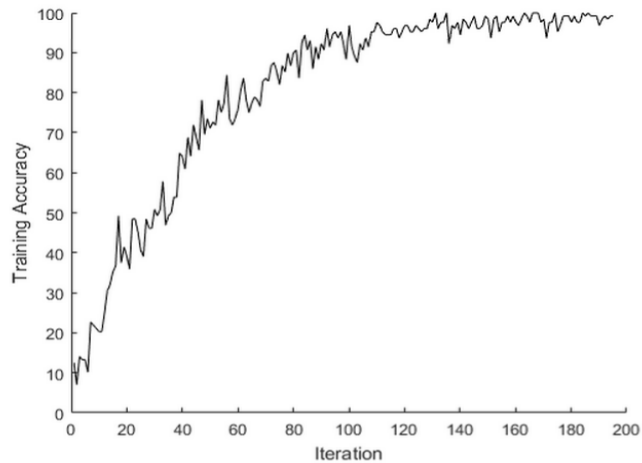


**Multiple GPU support**

```
opts = trainingOptions('sgdm', ...  
    'MaxEpochs', 100, ...  
    'MiniBatchSize', 250, ...  
    'InitialLearnRate', 0.00005, ...  
    'ExecutionEnvironment', 'auto' );
```

# Visualizing and Debugging Intermediate Results

## Training Accuracy Visualization

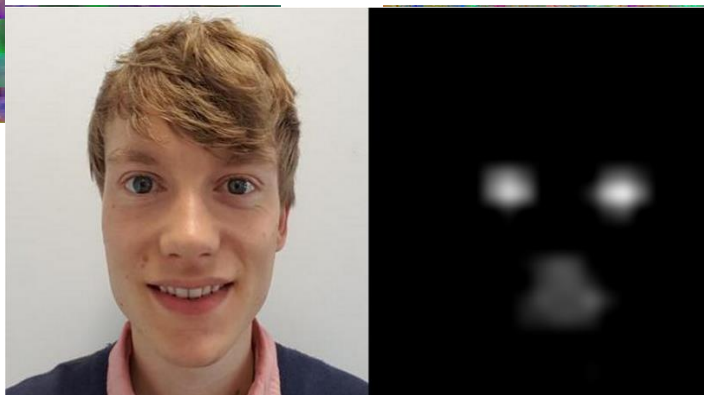


- Many options for visualizations and debugging
- Examples to get started

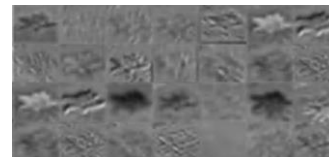
## Filters



## Layer Activations



## Activations

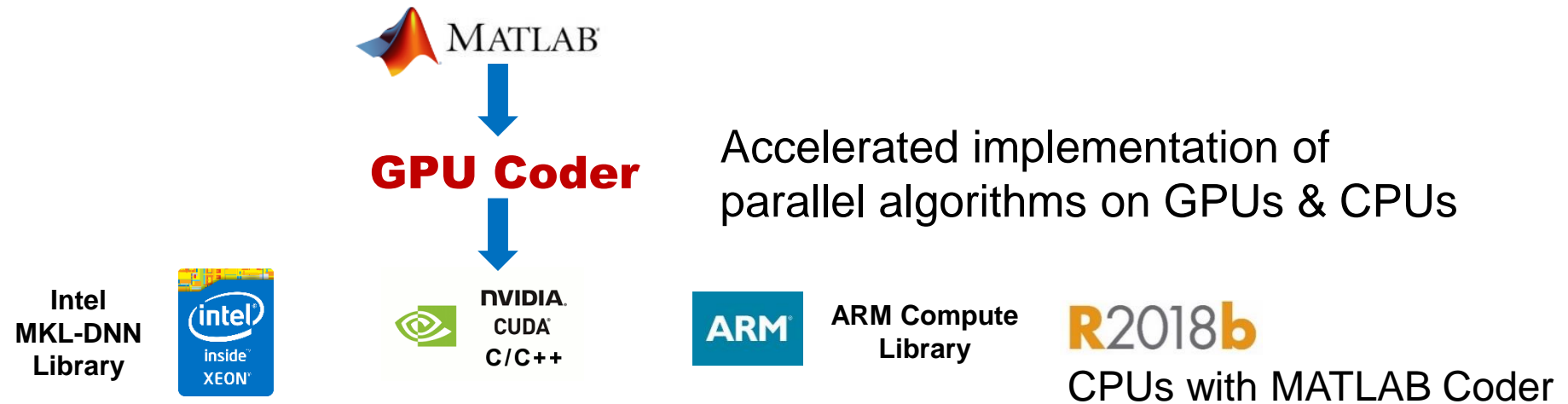


## Feature Visualization



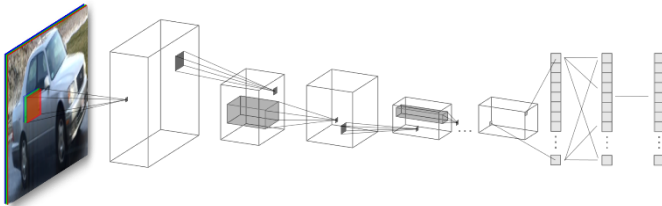


# GPU Coder for Deployment



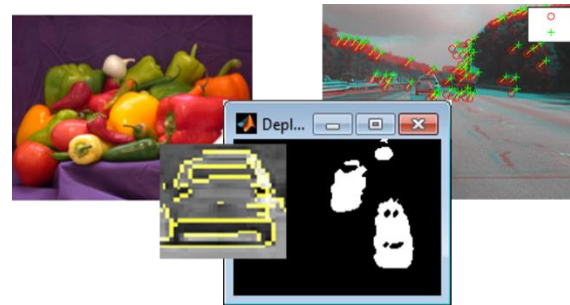
## Deep Neural Networks

Deep Learning, machine learning



## Image Processing and Computer Vision

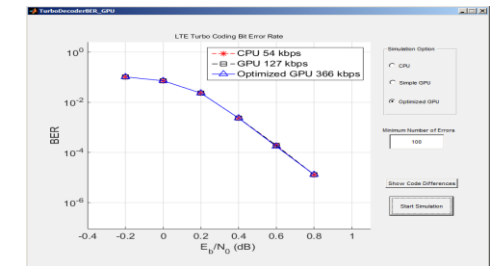
Image filtering, feature detection/extraction



**60x faster** than CPUs  
for stereo disparity

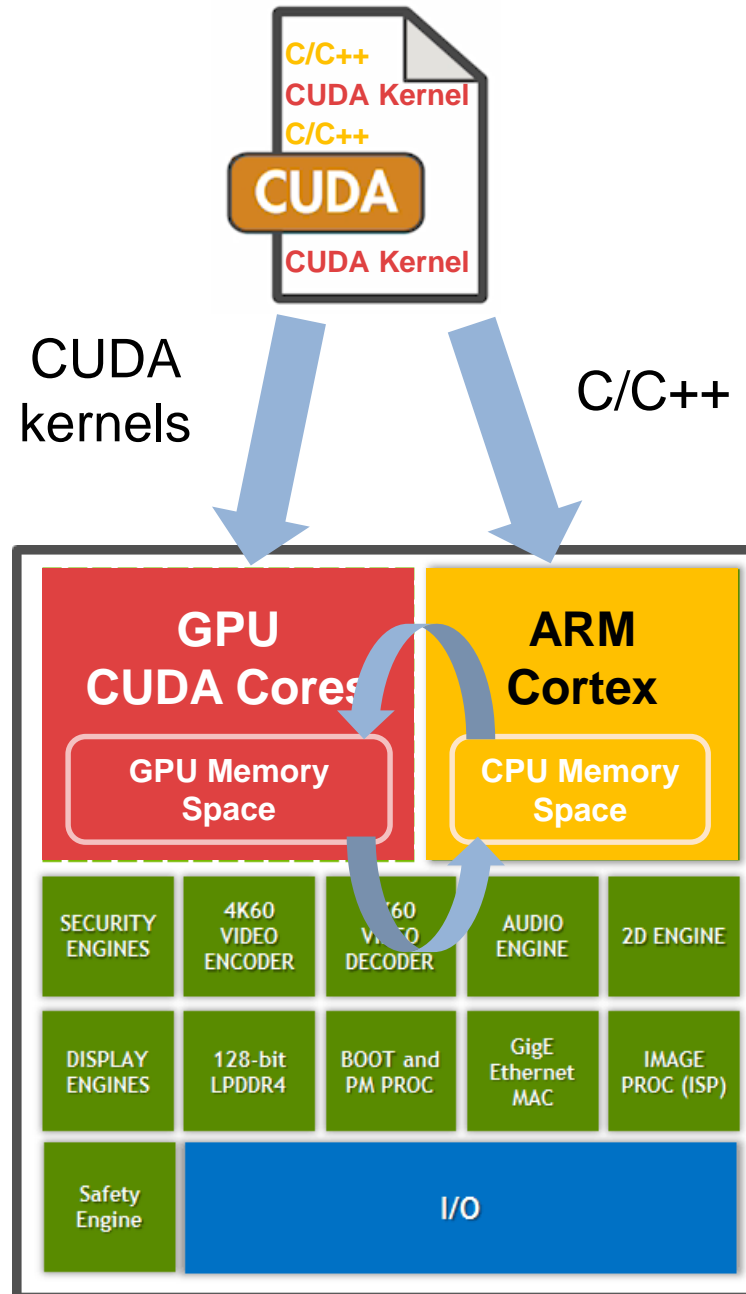
## Signal Processing and Communications

FFT, filtering, cross correlation,



**20x faster** than  
CPUs for FFTs

# GPUs and CUDA

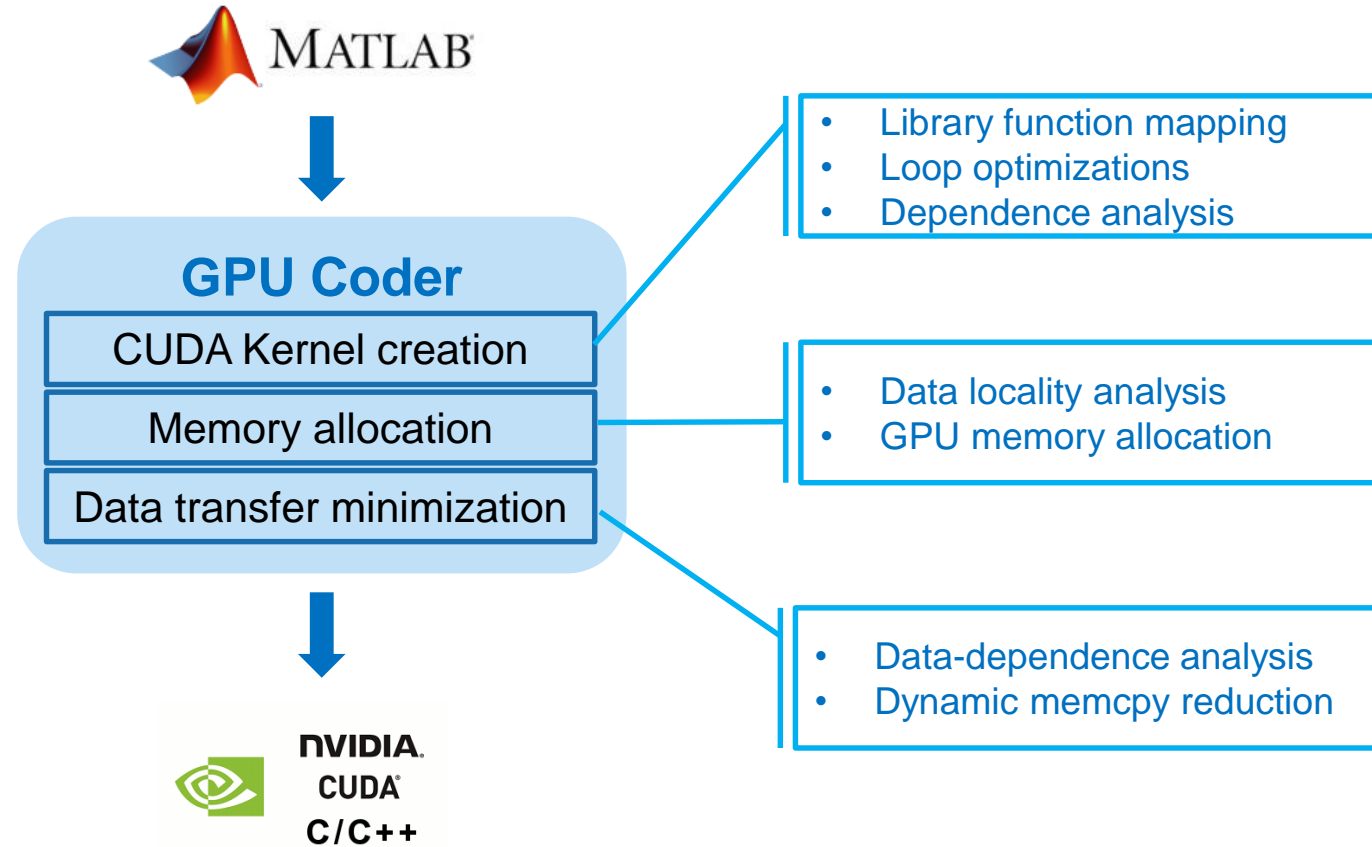




# Challenges of Programming in CUDA for GPUs

- 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



# GPU Coder Generates CUDA from MATLAB: saxpy

## Scalarized MATLAB

```
for i = 1:length(x)
    z(i) = a .* x(i) + y(i);
end
```



GPU Coder

## Vectorized MATLAB

```
z = a .* x + y;
```



Loops and matrix operations are directly compiled into kernels

## CUDA

```
cudaMalloc(&gpu_z, 8388608UL);
cudaMalloc(&gpu_x, 4194304UL);
cudaMalloc(&gpu_y, 4194304UL);
cudaMemcpy((void *)gpu_y, (void *)y, 4194304UL, cudaMemcpyHostToDevice);
cudaMemcpy((void *)gpu_x, (void *)x, 4194304UL, cudaMemcpyHostToDevice);
saxpy_kernel1<<<dim3(2048U, 1U, 1U), dim3(512U, 1U, 1U)>>>(&gpu_y, &gpu_x, a,
    &gpu_z);
cudaMemcpy((void *)z, (void *)gpu_z, 8388608UL, cudaMemcpyDeviceToHost);
cudaFree(gpu_y);
cudaFree(gpu_x);
cudaFree(gpu_z);
```

## CUDA kernel for GPU parallelization

```
static __global__ __launch_bounds__(512, 1) void saxpy_kernel1(const real32_T *y,
    const real32_T *x, real32_T a, real_T *z)
{
    int32_T i;

    i = (int32_T)((((gridDim.x * gridDim.y * blockIdx.z + gridDim.x * blockIdx.y)
        + blockIdx.x) * (blockDim.x * blockDim.y * blockDim.z) +
        threadIdx.z * blockDim.x * blockDim.y) + threadIdx.y *
        blockDim.x) + threadIdx.x);

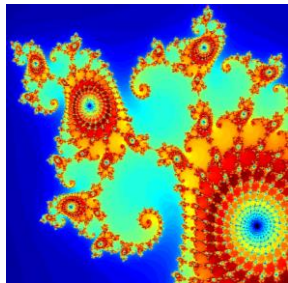
    if (!(i >= 1048576)) {
        z[i] = (real_T)(a * x[i] + y[i]);
    }
}
```

# Generated CUDA Optimized for Memory Performance

Kernel data allocation is automatically optimized

```

z = z0;
for n = 0:maxIterations
    z = z.*z + z0;
    inside = abs( z ) <= 2;
    count = count + inside;
end
count = log( count );
  
```



Mandelbrot space

GPU Coder

## CUDA kernel for GPU parallelization

```

static __global__ __launch_bounds__(512, 1) void kernel3(creal_T *z0, real_T
*count, creal_T *z)
{
    real_T z_im;
    real_T y[1000000];
    int32_T threadIdx;
    threadIdx = (int32_T)(blockDim.x * blockIdx.x + threadIdx.x);
    if (!(threadIdx >= 1000000)) {
        z_im = z[threadIdx].re * z[threadIdx].im + z[threadIdx].im * z[threadIdx].re;
        z[threadIdx].re = (z[threadIdx].re * z[threadIdx].re - z[threadIdx].im *
            z[threadIdx].im) + z0[threadIdx].re;
        z[threadIdx].im = z_im + z0[threadIdx].im;
        y[threadIdx] = hypot(z[threadIdx].re, z[threadIdx].im);
        count[threadIdx] += (real_T)(y[threadIdx] <= 2.0);
    }
}
  
```

## CUDA

...

```

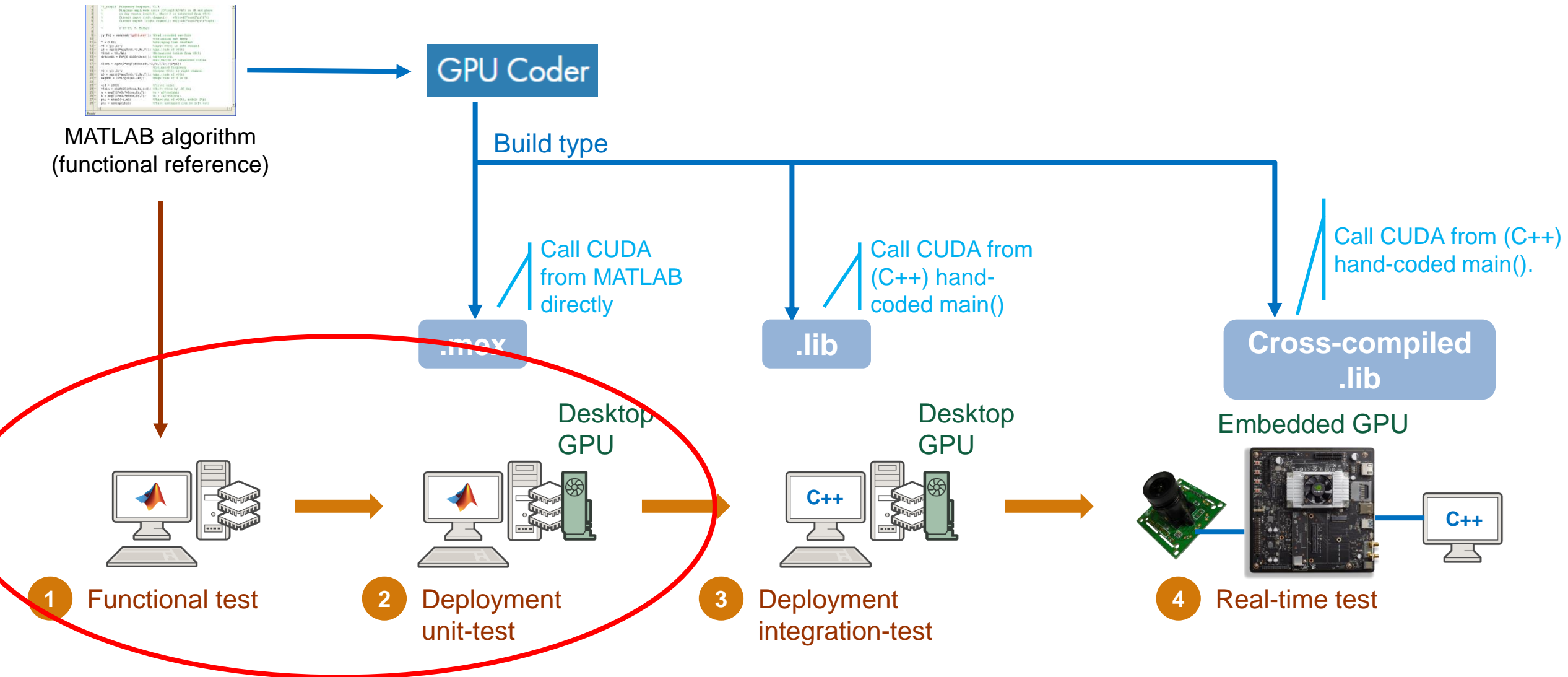
cudaMalloc(&gpu_xGrid, 8000000U);
cudaMalloc(&gpu_yGrid, 8000000U);

/* mandelbrot computation */
cudaMemcpy(gpu_yGrid, yGrid, 8000000U, cudaMemcpyHostToDevice);
cudaMemcpy(gpu_xGrid, xGrid, 8000000U, cudaMemcpyHostToDevice);
kernel1<<<dim3(1954U, 1U, 1U), dim3(512U, 1U, 1U)>>>(gpu_yGrid, gpu_xGrid,
    gpu_z, gpu_count, gpu_z0);
for (n = 0; n < (int32_T)(maxIterations + 1.0); n++) {
    kernel3<<<dim3(1954U, 1U, 1U), dim3(512U, 1U, 1U)>>>(gpu_z0, gpu_count,
        gpu_z);
}

kernel2<<<dim3(1954U, 1U, 1U), dim3(512U, 1U, 1U)>>>(gpu_count);
cudaMemcpy(count, gpu_count, 8000000U, cudaMemcpyDeviceToHost);
cudaFree(gpu_yGrid);
  
```

...

# Algorithm Design to Embedded Deployment Workflow



# Demo: Alexnet Deployment with 'mex' Code Generation

vravicha-deb8-64:2 (vravicha) - TigerVNC

Applications Menu | MATLAB R2017b | Editor - /tmp/webcam... | vravicha@vravicha-d... | vravicha@vravicha-d... | 18:00 | Vandana Ravichandran

**MATLAB R2017b**

HOME | PLOTS | APPS

Search Documentation | Log In

New Script | New Live Script | New | Open | Find Files | Compare | Import Data | Save Workspace | New Variable | Open Variable | Clear Workspace | Analyze Code | Run and Time | Clear Commands | Simulink | Layout | Preferences | Set Path | Parallel | Add-Ons | Help | Community | Request Support | Learn MATLAB

FILE | VARIABLE | CODE | SIMULINK | ENVIRONMENT | RESOURCES

Current Folder: /tmp/webcamtt

- test\_alexnet\_codegen.m
- synsetWords.txt
- peppers\_out.png
- peppers.png
- old\_workspace.mat
- getAlexnet.m
- cleanup.m
- alexnet\_webcam.m
- alexnet\_predict.prj
- alexnet\_predict.m
- alexnet\_live.m

alexnet\_predict.m (Function)

Workspace

Name	Value
------	-------

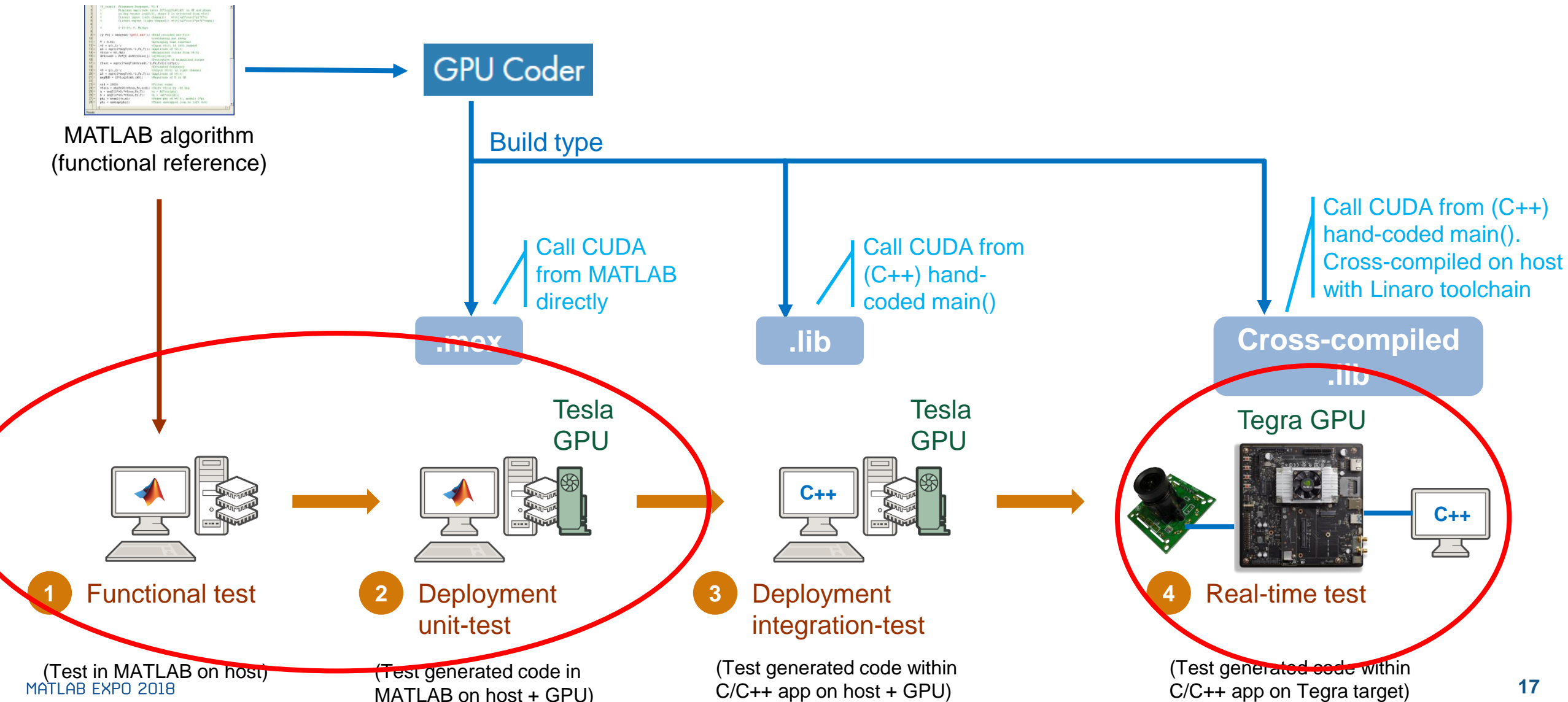
Command Window

New to MATLAB? See resources for [Getting Started](#).

```
fx >>
```

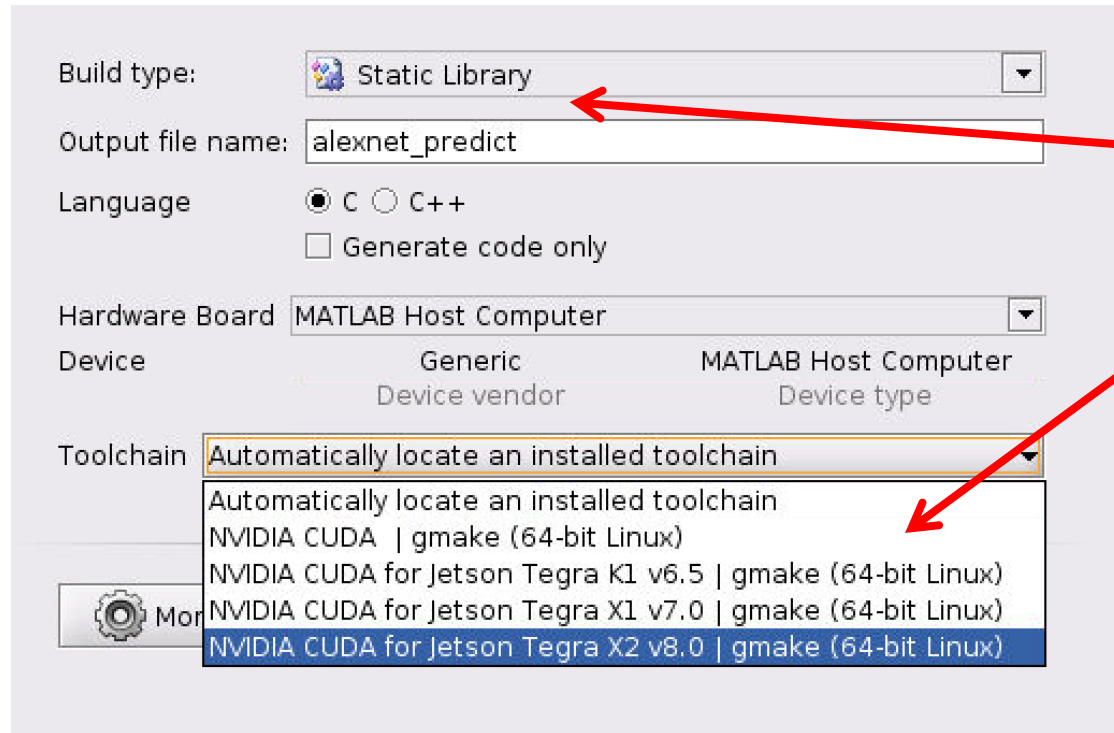
Windows Taskbar: 3:30 AM 10/5/2017

# Algorithm Design to Embedded Deployment on Tegra GPU



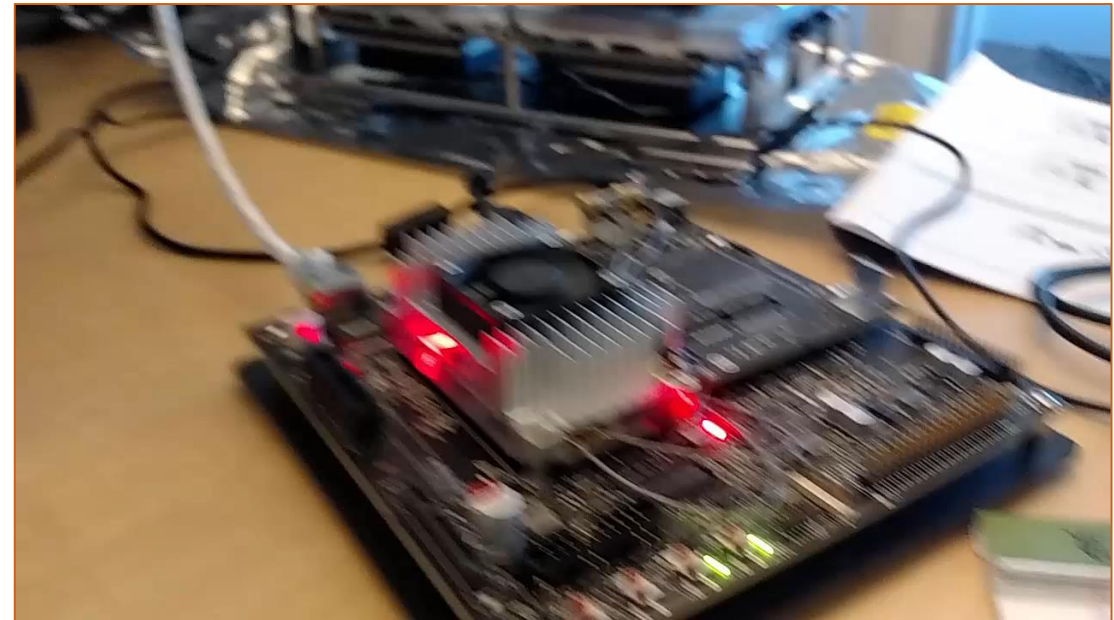


# Alexnet Deployment to Tegra: Cross-Compiled with 'lib'



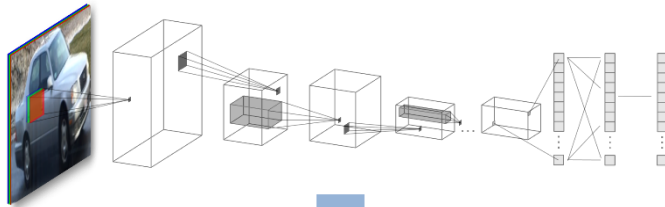
## Two small changes

1. Change build-type to 'lib'
2. Select cross-compile toolchain



# End-to-End Application: Lane Detection

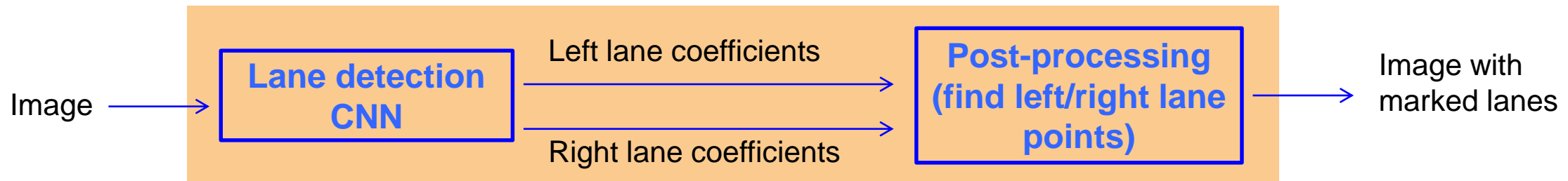
Alexnet



**Transfer Learning**



Output of CNN is lane parabola coefficients according to:  $y = ax^2 + bx + c$

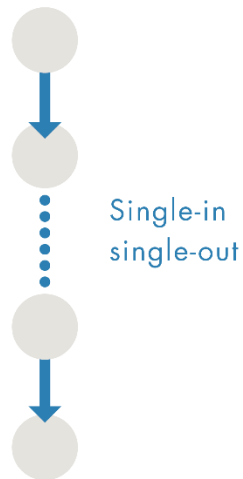


**GPU coder generates code for whole application**



# Deep Learning Network Support (with Neural Network Toolbox)

## SeriesNetwork



GPU Coder: **R2017b**

Networks: MNist  
Alexnet  
YOLO  
VGG  
Lane detection  
Pedestrian detection

## DAGNetwork

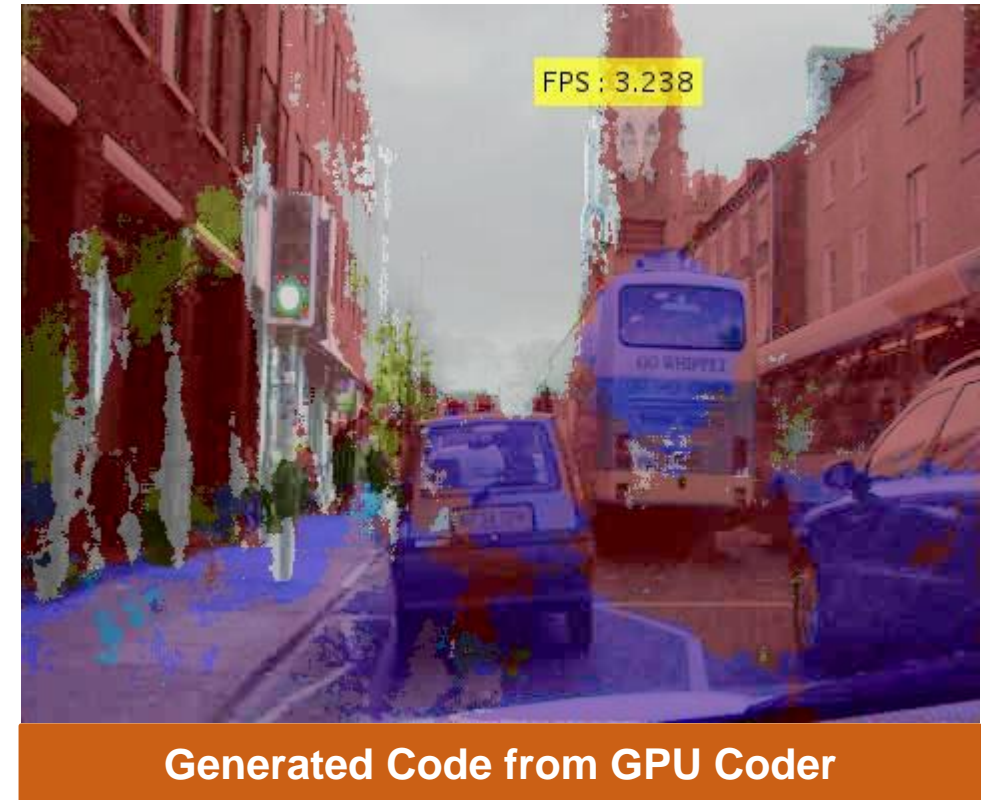
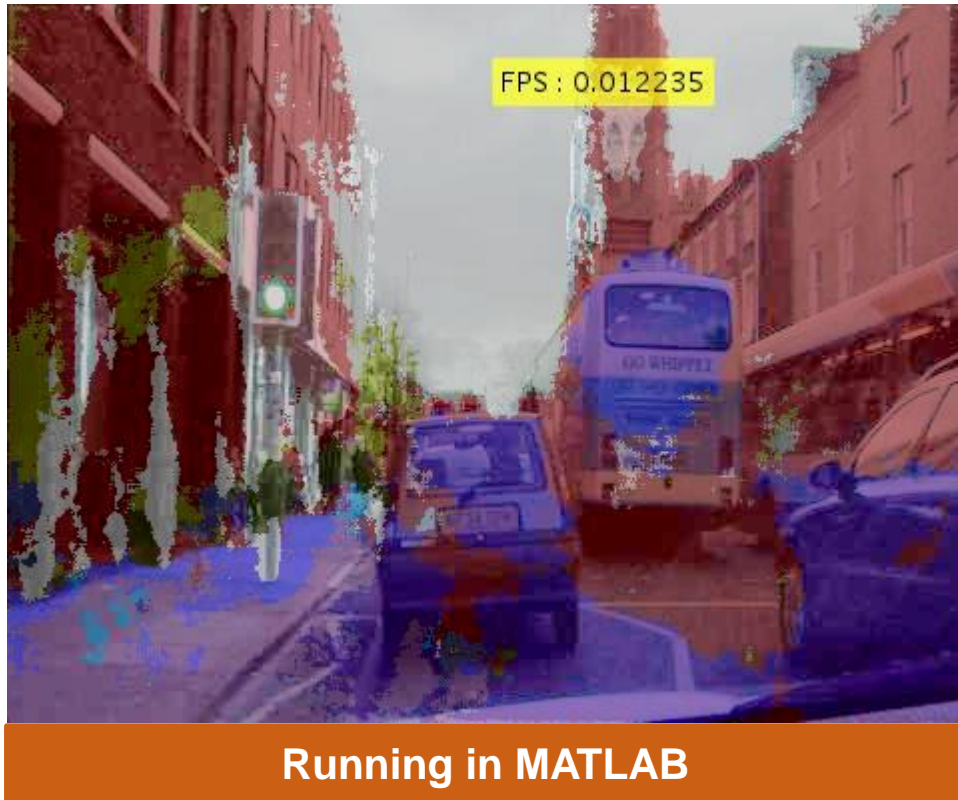


GPU Coder: **R2018a**

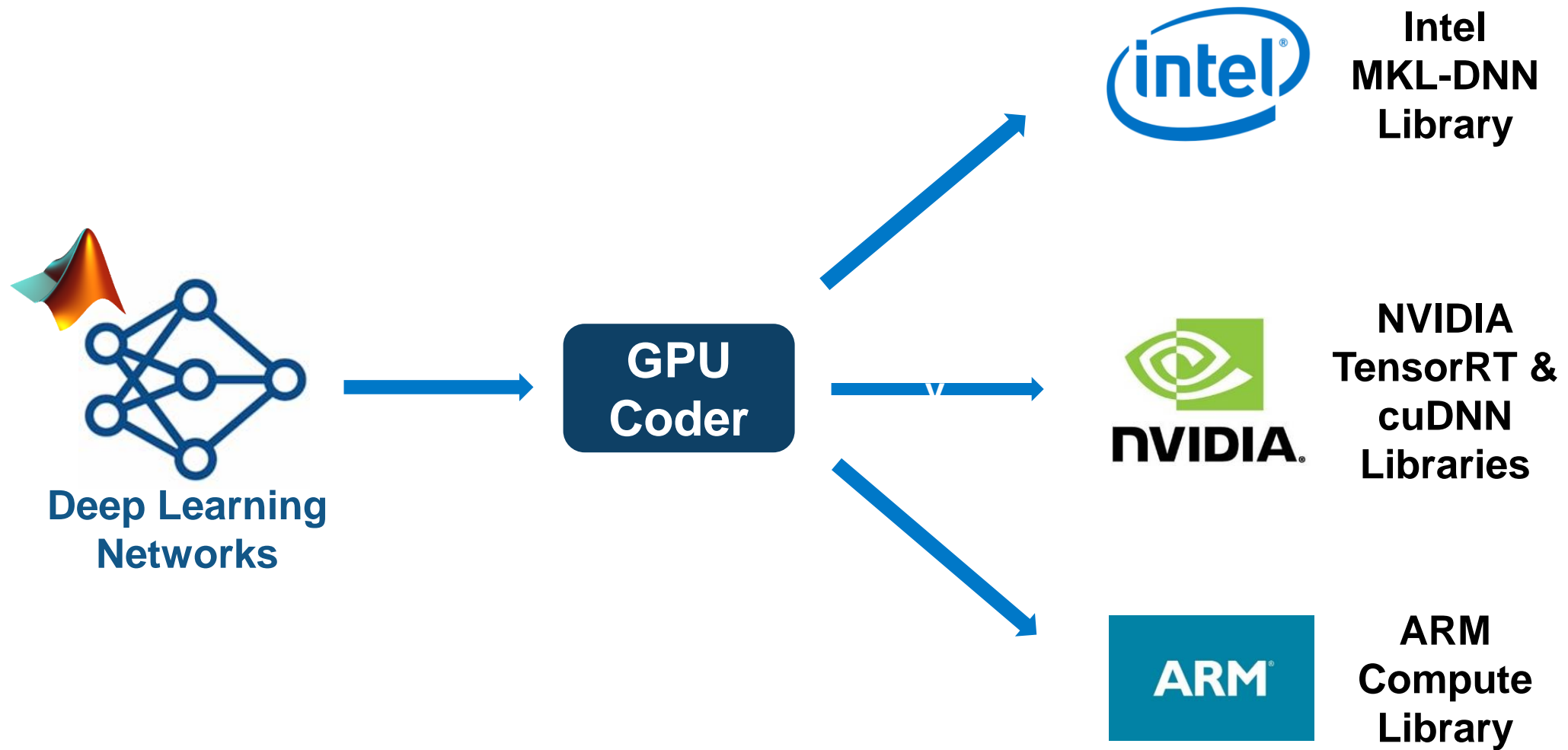
Networks: GoogLeNet  
ResNet  
SegNet  
DeconvNet

} Object  
detection  
} Semantic  
segmentation

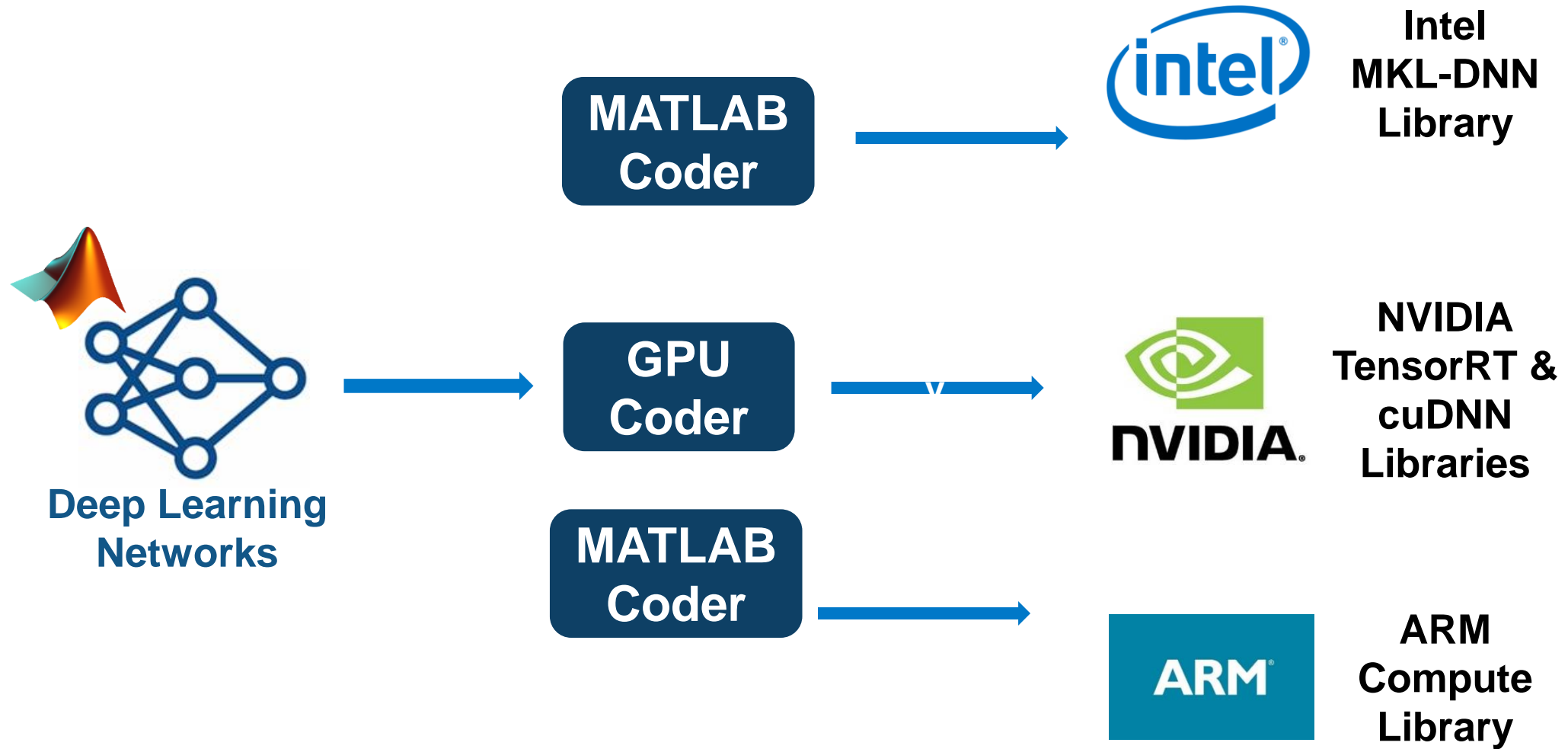
# Semantic Segmentation



# Deploying to CPUs

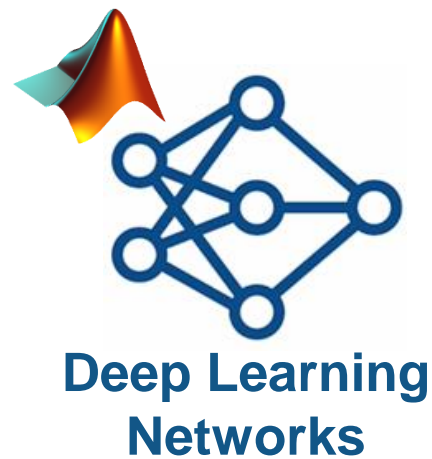


# Deploying to CPUs





# Deploying to CPUs



**MATLAB  
Coder**



**GPU  
Coder**



**NVIDIA  
TensorRT &  
cuDNN  
Libraries**

**MATLAB  
Coder**





# How Good is Generated Code Performance

- Performance of image processing and computer vision
- Performance MATLAB / GPU Coder / TensorFlow / MXNet / PyTorch / etc.

# GPU Coder for Image Processing and Computer Vision



Fog removal



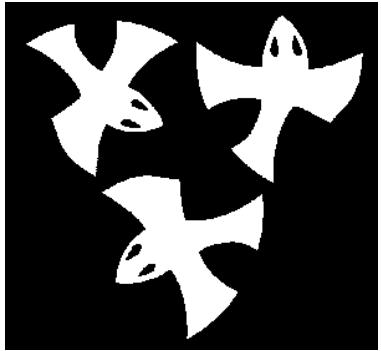
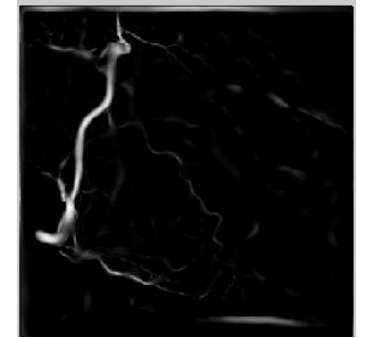
5x speedup



Frangi filter



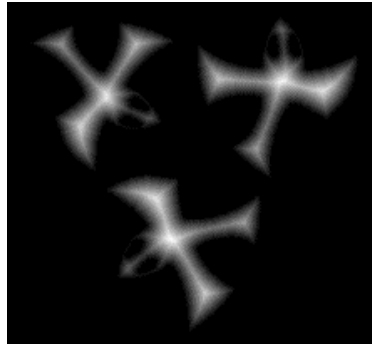
3x speedup



Distance transform



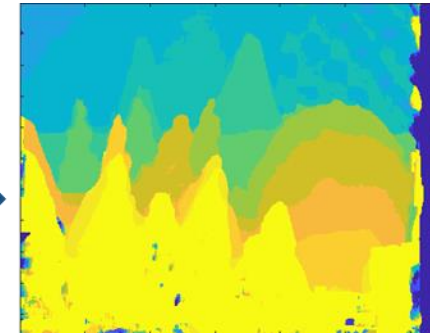
8x speedup



Stereo disparity



50x speedup



Ray tracing



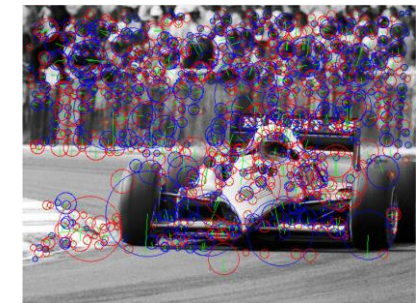
18x speedup



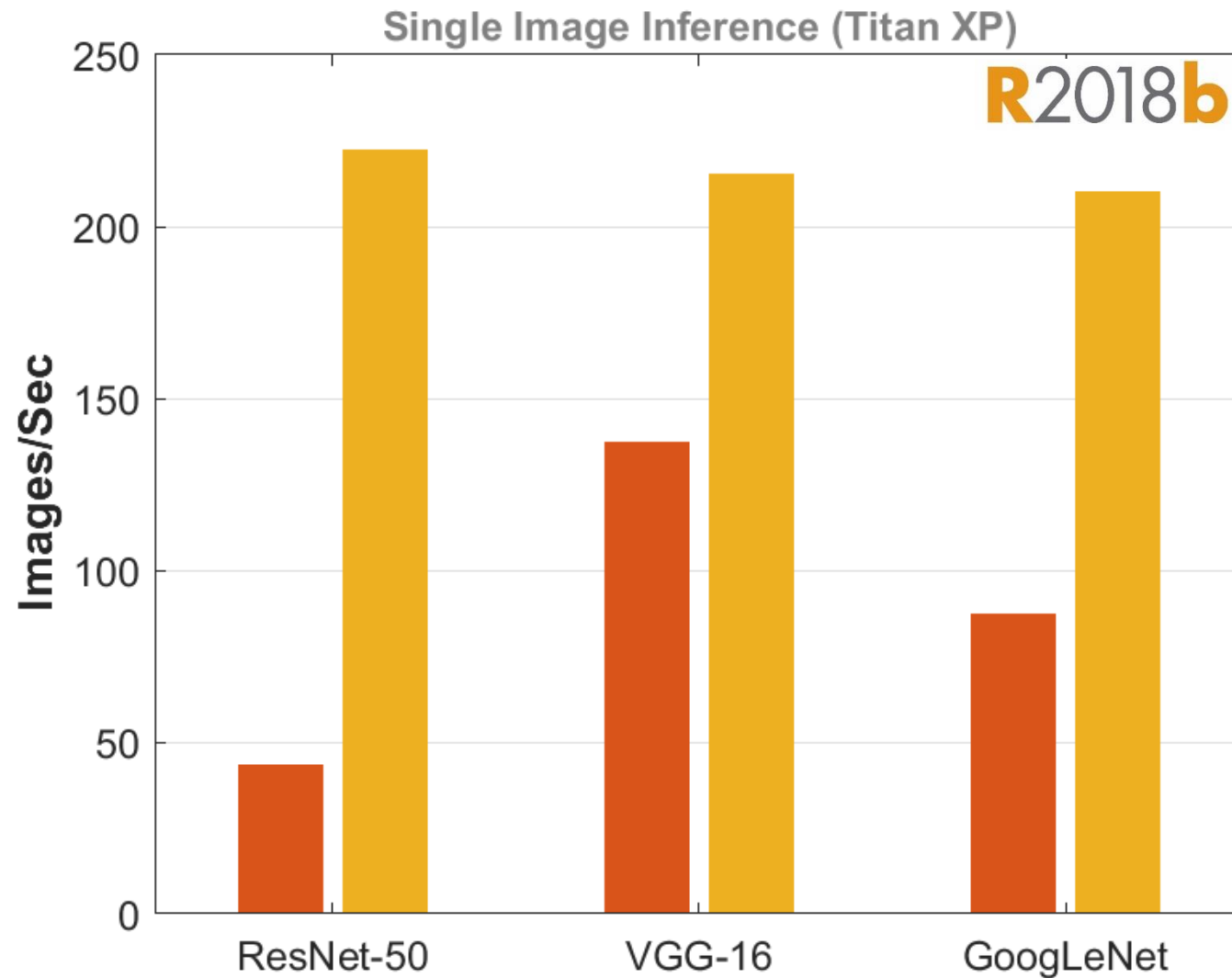
SURF feature extraction



700x speedup

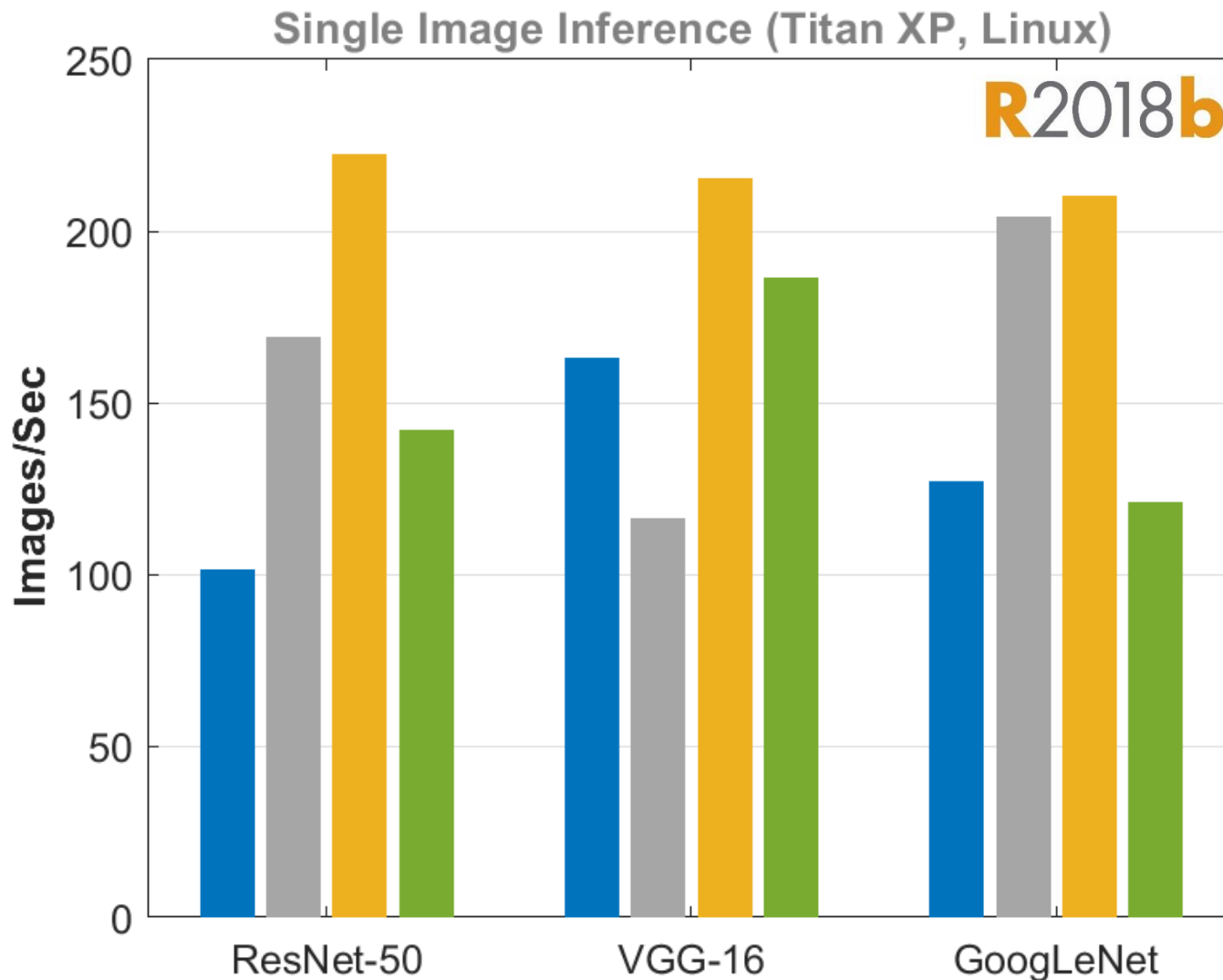


# GPU Coder speeds up MATLAB about 2x



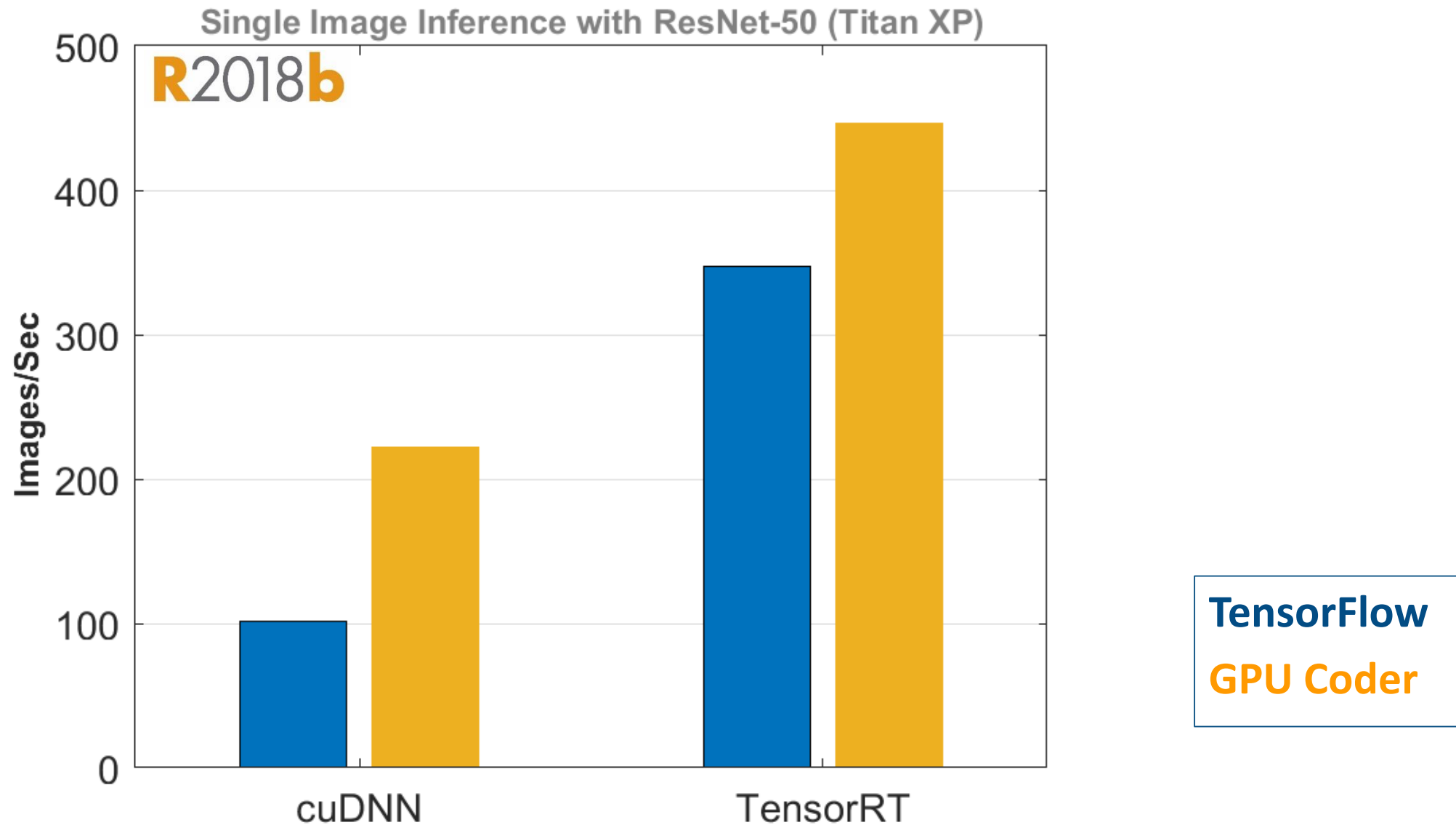
Intel® Xeon® CPU 3.6 GHz - NVIDIA libraries: CUDA9 - cuDNN 7

# With GPU Coder, MATLAB is fast



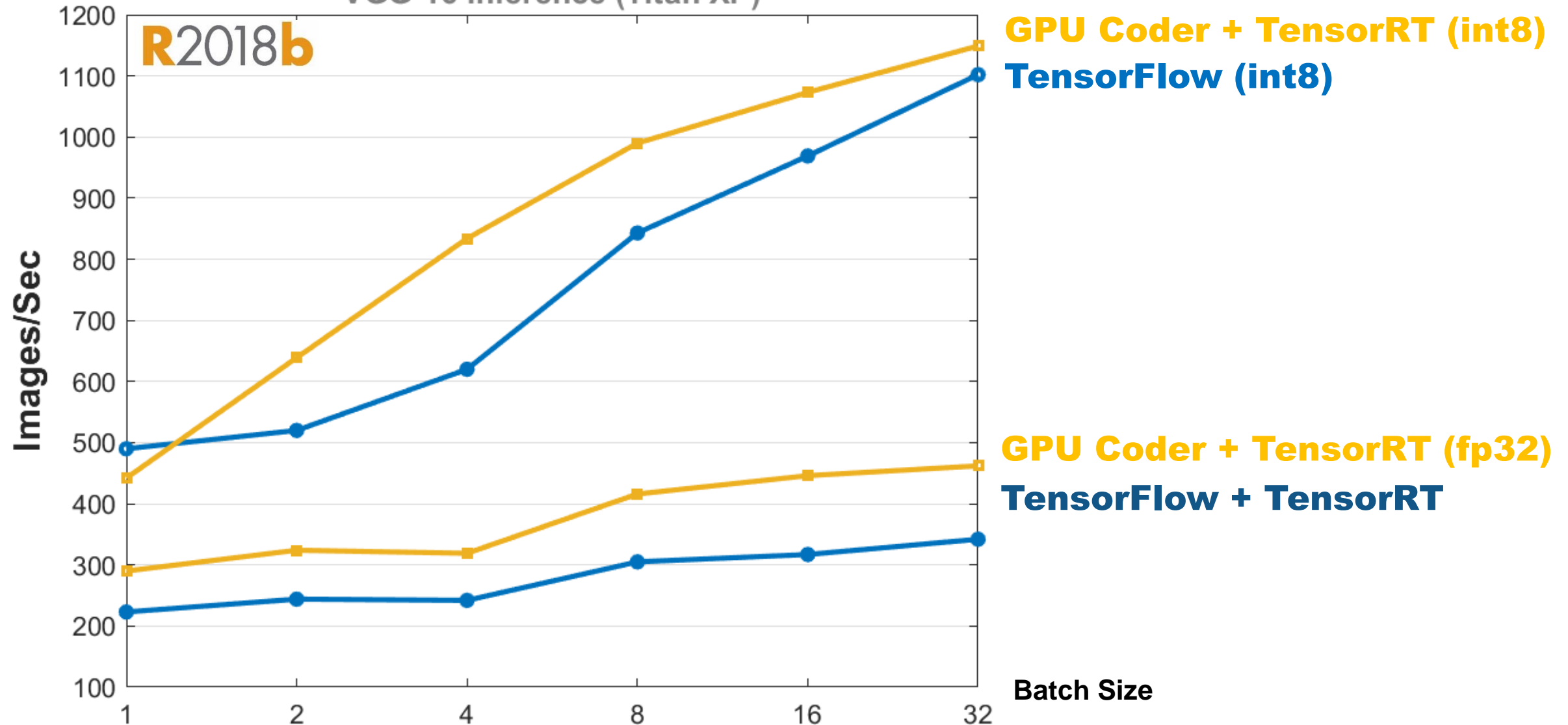
**Faster than TensorFlow,  
MXNet, and PyTorch**

# TensorRT speeds up inference for TensorFlow and GPU Coder

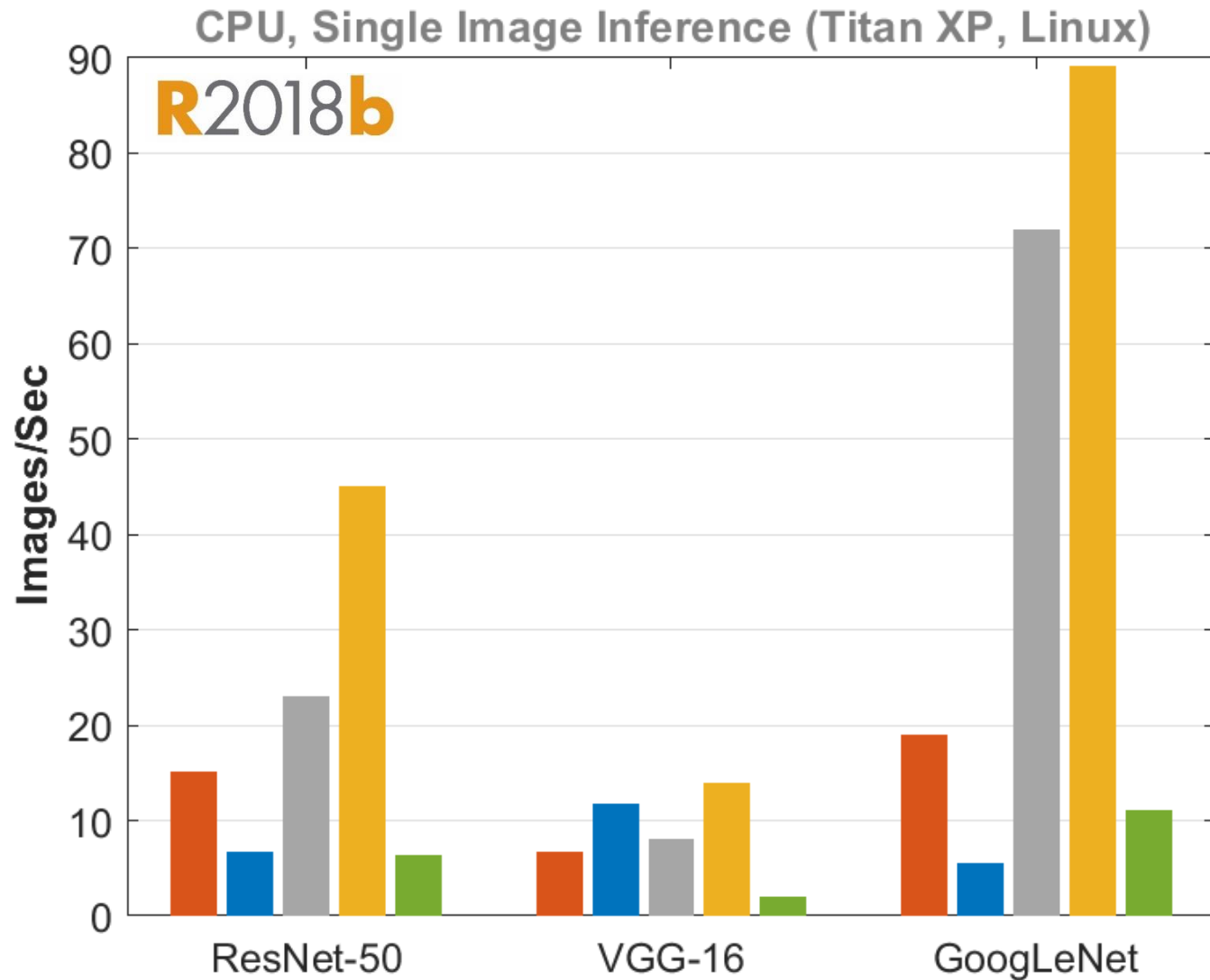


# Even higher Speeds with Integer Arithmetic (int8)

VGG-16 Inference (Titan XP)



# CPU Performance



**GPU Coder** also fastest on CPU

Only **MXNet** faster than **MATLAB**



# Design Your DNNs in MATLAB, Deploy with GPU Coder



- **Manage** large image sets
- **Automate** image labeling
- **Easy access** to models
- **Acceleration** with GPU's
- **Scale** to clusters
- **Automate compilation to GPUs and CPUs using GPU Coder**