

# MATLAB EXPO 2018

## Deploying Deep Learning Networks to Embedded GPUs and CPUs

Daniel Aronsson



# Agenda

- What is CUDA code?
- What is GPU Coder?
- Why use GPU Coder?
- How to use GPU Coder?
- How fast is GPU Coder?
- Key takeaways

# Algorithm Design to Embedded Deployment Workflow



```

1 %_demo1: Embedded Deployment Test
2 % This script demonstrates the deployment of a MATLAB algorithm to an embedded GPU.
3 % It shows the compilation of the algorithm into a C++ executable and its execution on the embedded GPU.
4 % The algorithm is a simple matrix multiplication.
5
6 % Create the MATLAB algorithm
7 % Create a MATLAB script that performs a matrix multiplication
8 % Create a MATLAB function that performs a matrix multiplication
9
10 % Compile the MATLAB algorithm into a C++ executable
11 % Use the MATLAB C++ compiler to compile the MATLAB algorithm into a C++ executable
12 % The MATLAB C++ compiler is located in the MATLAB installation directory
13 % The MATLAB C++ compiler is used to compile the MATLAB algorithm into a C++ executable
14 % The MATLAB C++ compiler is used to compile the MATLAB algorithm into a C++ executable
15
16 % Execute the C++ executable on the embedded GPU
17 % Use the MATLAB C++ compiler to compile the MATLAB algorithm into a C++ executable
18 % The MATLAB C++ compiler is used to compile the MATLAB algorithm into a C++ executable
19 % The MATLAB C++ compiler is used to compile the MATLAB algorithm into a C++ executable
20
21 % End of script

```

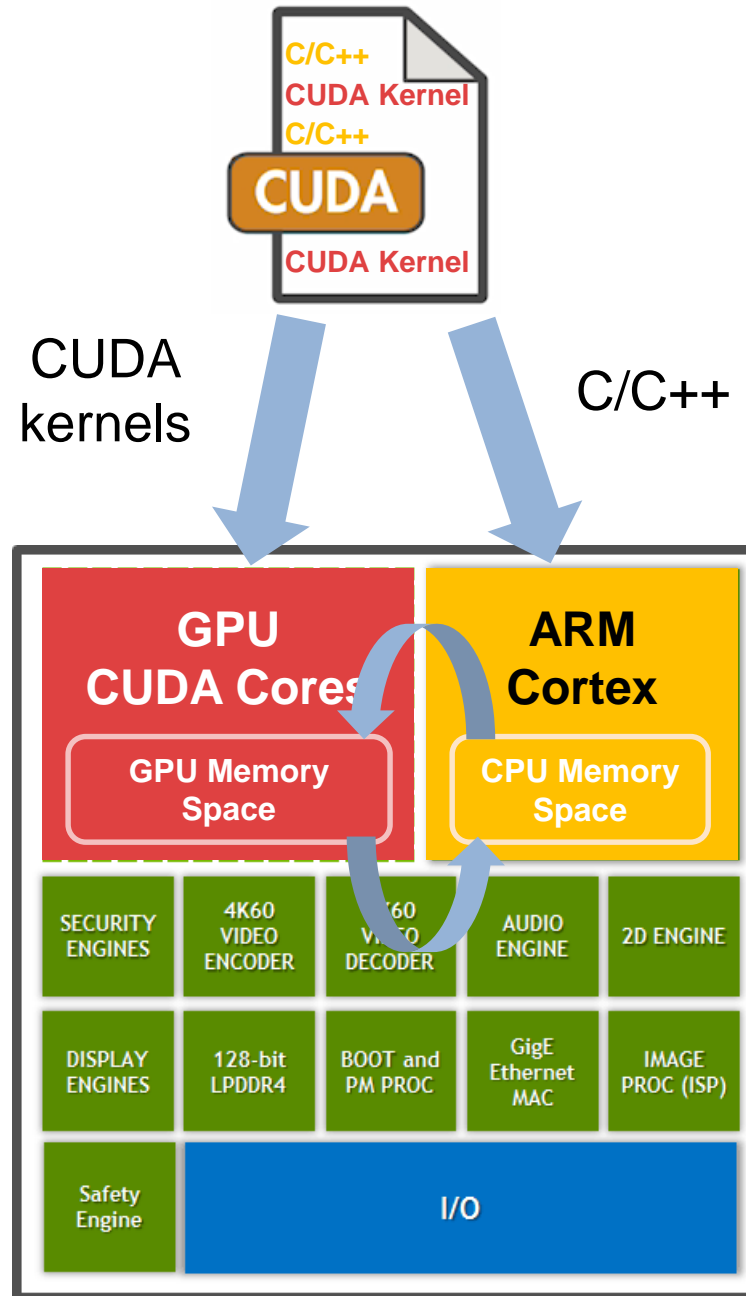
MATLAB algorithm  
(functional reference)



Embedded GPU

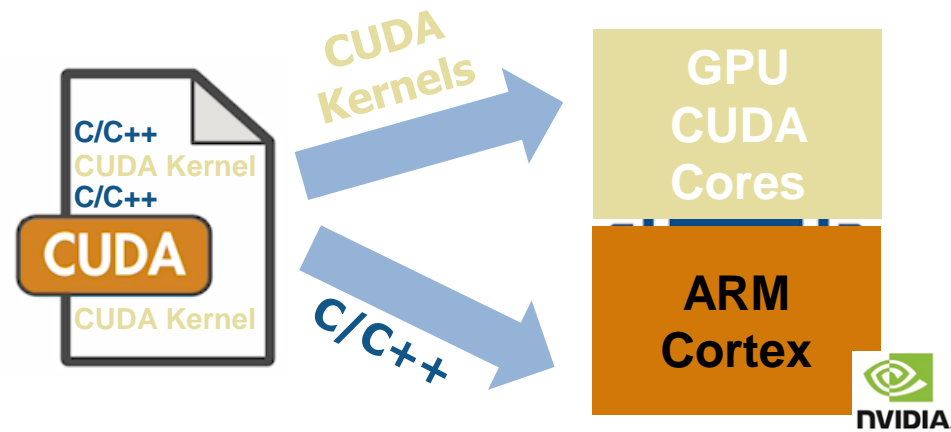


# GPUs and CUDA



# What is CUDA code?

CUDA extends C/C++ code with constructs for parallel computing



# What does CUDA code look like?

```

void foo(const real_T A[100000000], const real_T B[100000000],
real_T C[100000000])
{
    real_T *gpu_B;
    real_T *gpu_A;
    real_T *gpu_C;
    cudaMalloc(&gpu_C, 800000000ULL);
    cudaMalloc(&gpu_A, 800000000ULL);
    cudaMalloc(&gpu_B, 800000000ULL);
    cudaMemcpy((void *)gpu_B, (void *)&B[0], 800000000ULL,
cudaMemcpyHostToDevice);
    cudaMemcpy((void *)gpu_A, (void *)&A[0], 800000000ULL,
cudaMemcpyHostToDevice);
    foo_kernell<<<dim3(313U, 313U, 1U), dim3(32U, 32U,
1U)>>>(gpu_B, gpu_A, gpu_C);
    cudaMemcpy((void *)&C[0], (void *)gpu_C, 800000000ULL,
cudaMemcpyDeviceToHost);
    cudaFree(gpu_B);
    cudaFree(gpu_A);
    cudaFree(gpu_C);
}

```

```

static __global__ __launch_bounds__(1024, 1)
void foo_kernell(const real_T *B,
const real_T *A, real_T *C)
{
    uint32_T threadId;
    int32_T i0;
    threadId = (uint32_T)mwGetGlobalThreadIndex();
    i0 = (int32_T)threadId;
    if (!(i0 >= 100000000)) {
        C[i0] = A[i0] * B[i0];
    }
}

```

**function C = foo(A, B)**  
**C = A\*B;**

# Challenges for the CUDA programmer

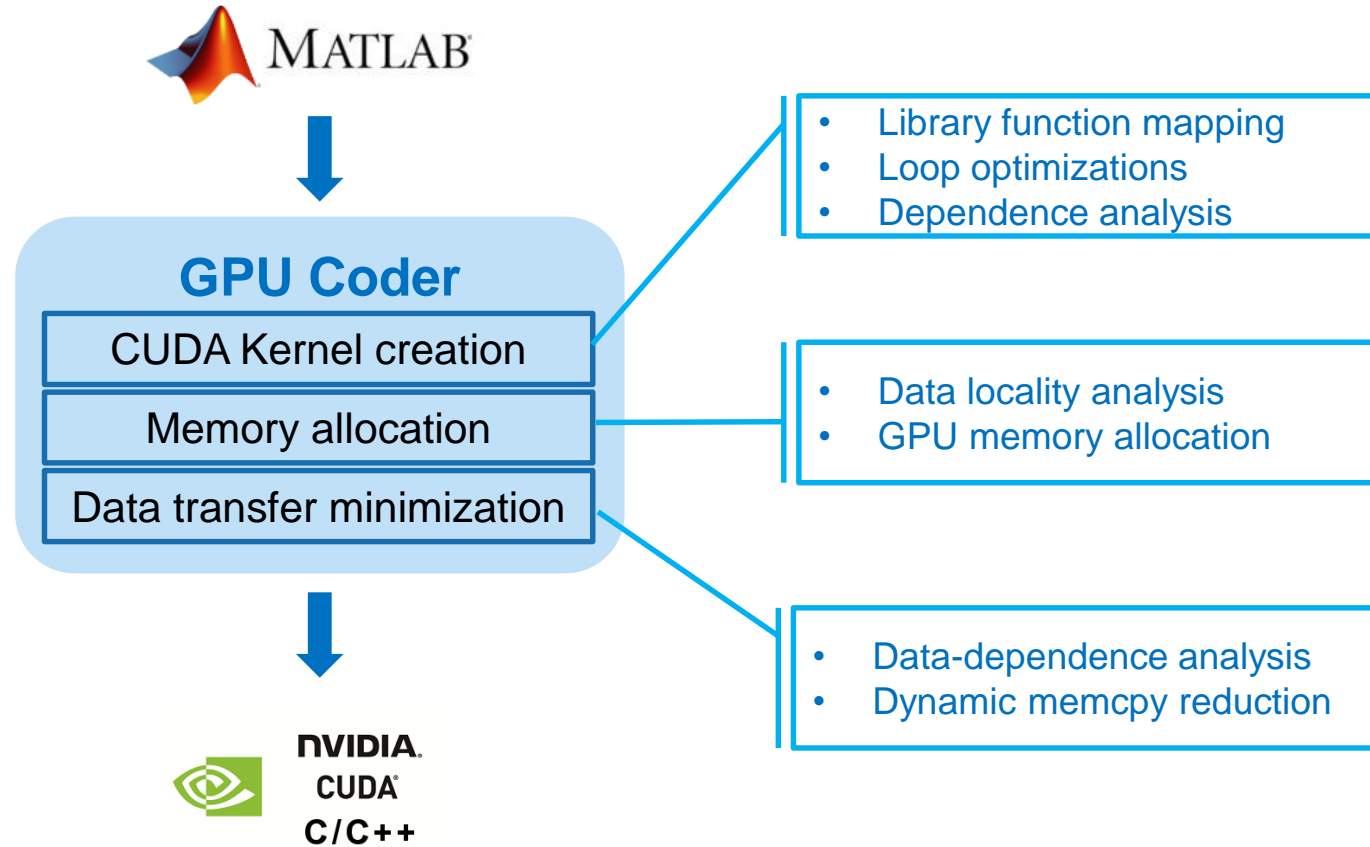
- Learning to program in CUDA
  - Need to rewrite algorithms for parallel processing paradigm
- Creating CUDA kernels
  - Need to analyze algorithms to create CUDA kernels that maximize parallel processing
- Allocating memory
  - Need to deal with memory allocation on both CPU and GPU memory spaces
- Minimizing data transfers
  - Need to minimize while ensuring required data transfers are done at the appropriate parts of your algorithm

# What is GPU Coder?

- Generates **CUDA** code for NVIDIA GPUs
- Also generates code for Deep Neural Networks for Intel CPUs and ARM Cortex-A platforms.



# 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;
```



## 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]);
    }
}
```

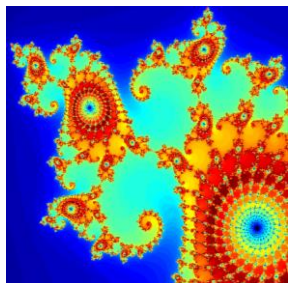
Loops and matrix operations are directly compiled into kernels

# 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

## 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);
    
```

...  
...

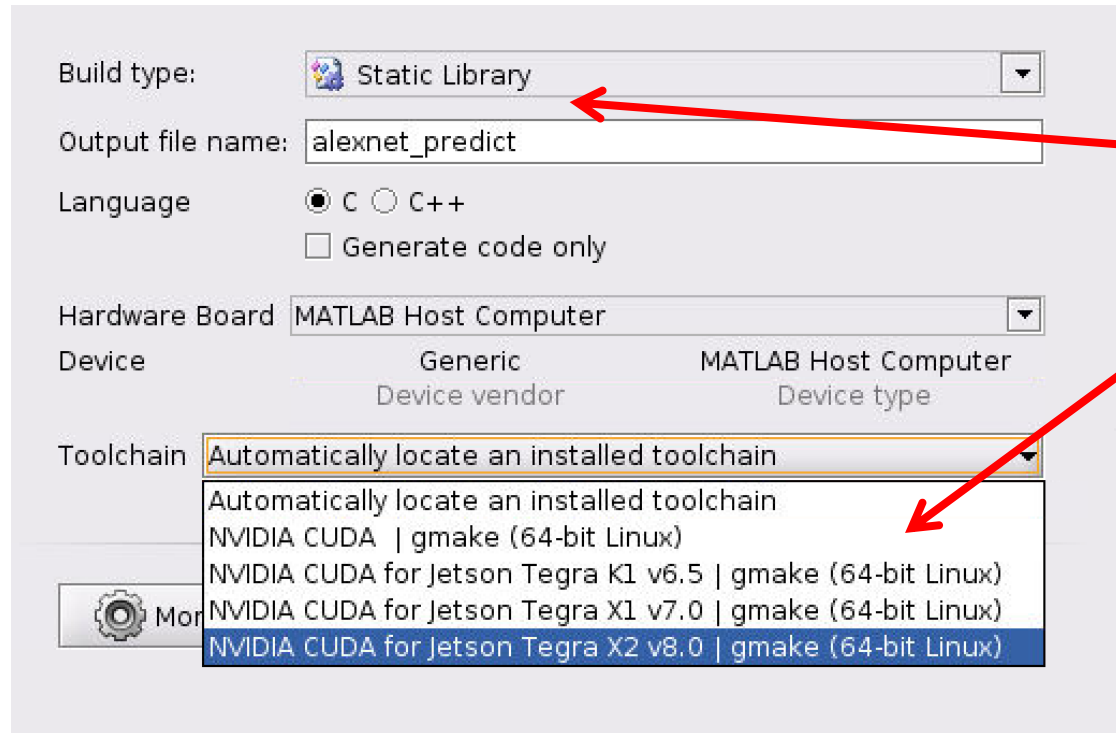


# Demo: Alexnet Deployment with 'mex' Code Generation

The image shows a MATLAB R2017b environment. The top toolbar includes options like 'New Script', 'Open', 'Save', 'Run and Time', and 'Set Path'. The left sidebar shows the 'Current Folder' with files such as 'test\_alexnet\_codegen.m', 'synsetWords.txt', 'peppers\_out.png', 'old\_workspace.mat', 'getAlexnet.m', 'cleanup.m', 'alexnet\_webcam.m', 'alexnet\_predict.prj', 'alexnet\_predict.m', and 'alexnet\_live.m'. The 'Workspace' section is currently empty. The 'Command Window' on the right contains the prompt 'fx >>'.

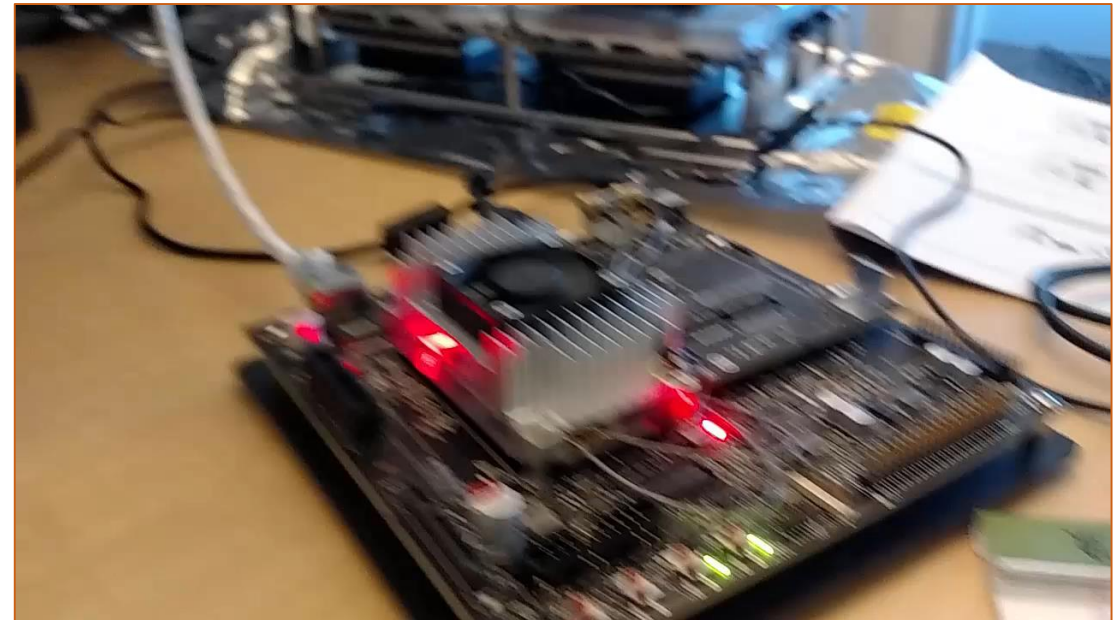


# Alexnet Deployment to Tegra: Cross-Compiled with 'lib'

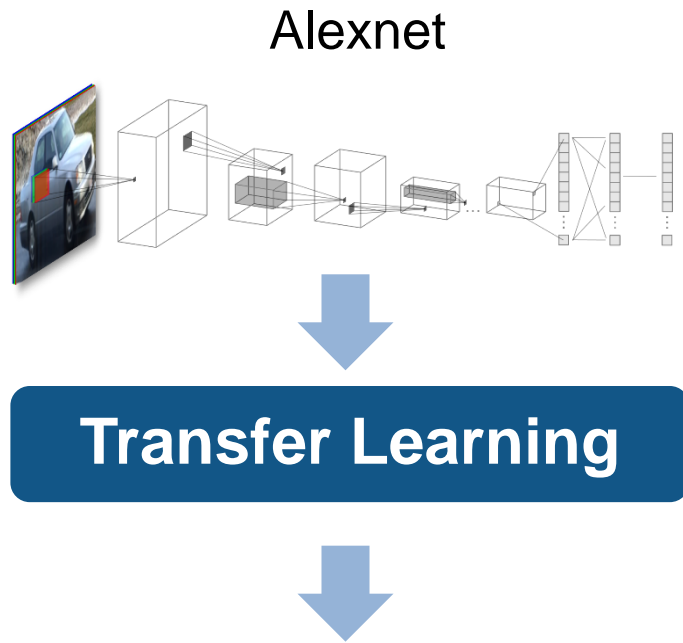


## Two small changes

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



# End-to-End Application: Lane Detection



**nvidia** ACCELERATED COMPUTING Downloads Training Ecosystem

**PARALLEL FOR ALL** Features Pro Tips Spotlights CUDACasts

← Previous Next →

## Deep Learning for Automated Driving with MATLAB

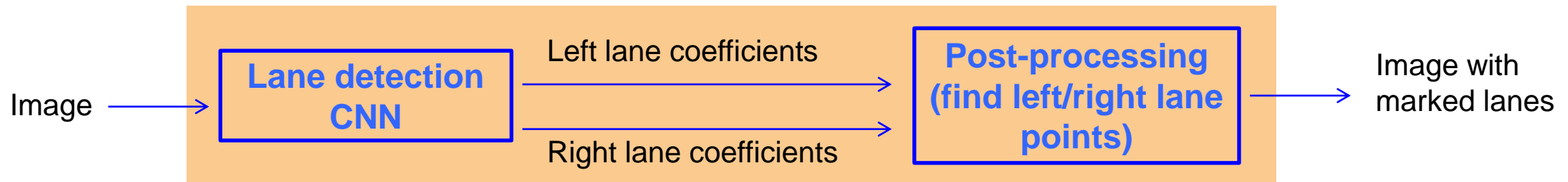
Share: [Twitter](#) [LinkedIn](#) [Facebook](#) [Google+](#) [Email](#)

Posted on July 20, 2017 by Avinash Nehemiah and Arvind Jayaraman | 0 Comments

Tagged Autonomous Vehicles, Deep Learning, MATLAB

You've probably seen headlines about innovation in automated driving now that there are several cars available on the market that have some level of self-driving capability. I often get questions from colleagues on how automated driving systems perceive their environment and make "human-like"

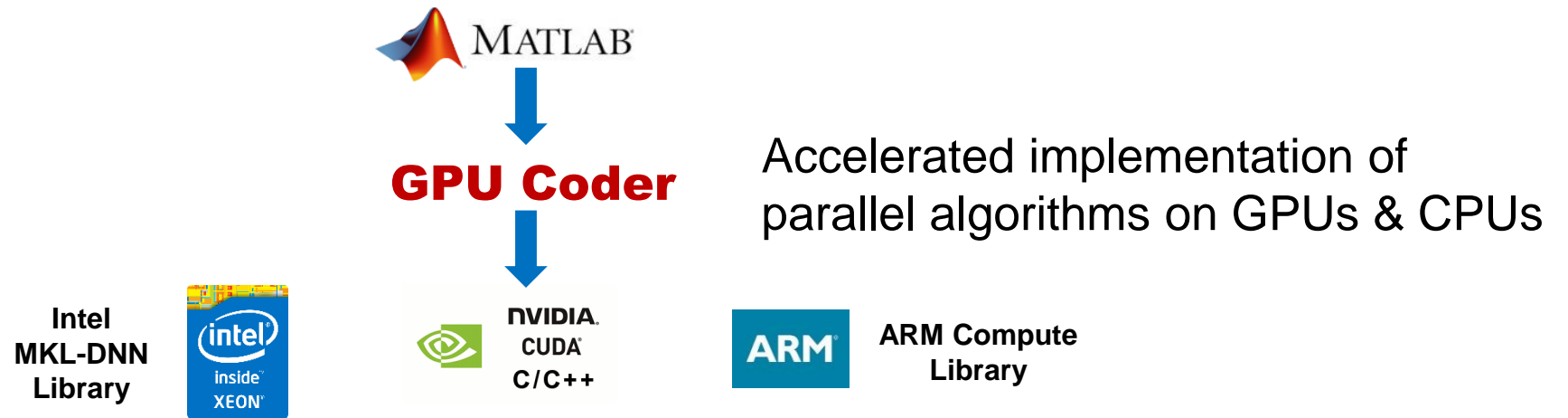
Output of CNN is lane parabola coefficients according to:  $y = ax^2 + bx + c$



**GPU coder generates code for whole application**

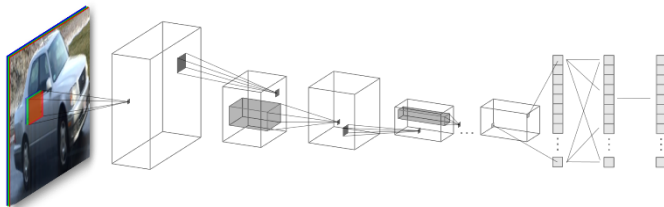


# GPU Coder for Deployment



## Deep Neural Networks

Deep Learning, machine learning



**5x faster** than TensorFlow  
**2x faster** than MXNet

## Image Processing and Computer Vision

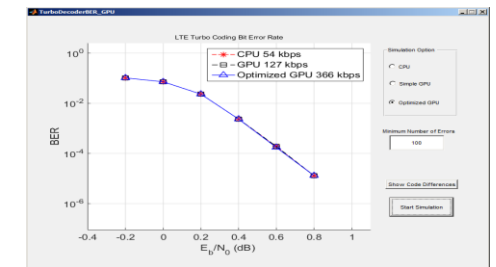
Image filtering, feature detection/extraction



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

## Signal Processing and Communications

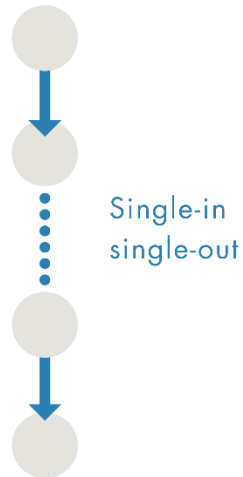
FFT, filtering, cross correlation,



**20x faster** than CPUs for FFTs

# Deep Learning Network Support (with Neural Network Toolbox)

## SeriesNetwork



GPU Coder: **R2017b**

Networks: MNist  
 Alexnet  
 YOLO  
 VGG  
 Lane detection  
 Pedestrian detection

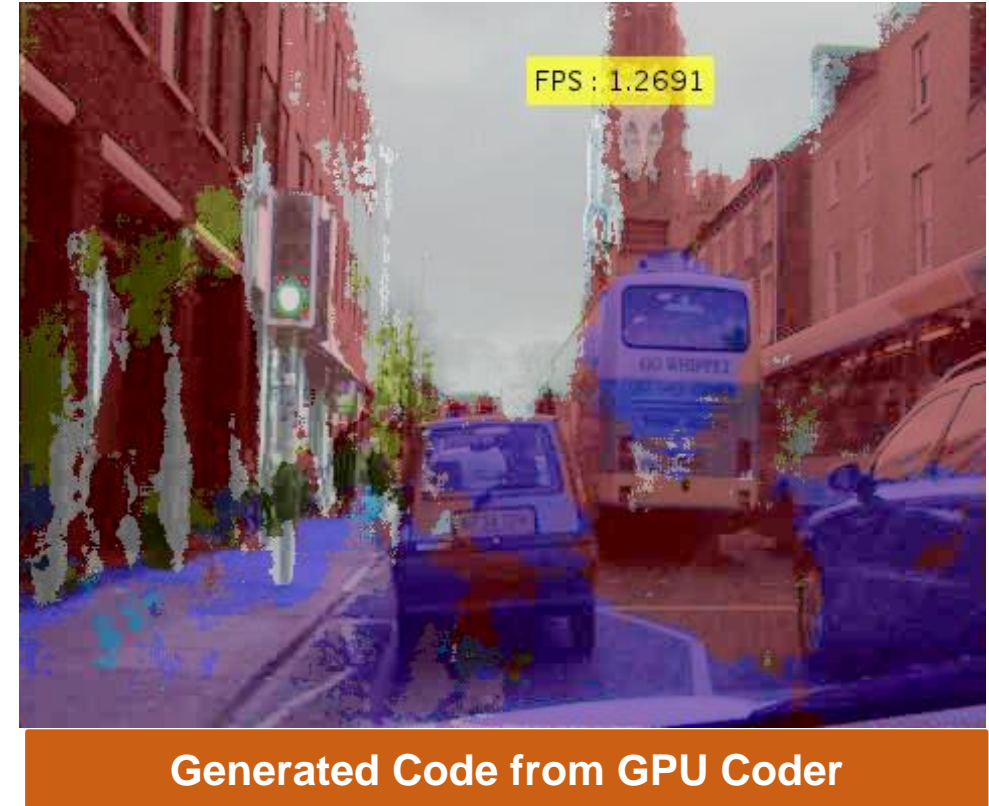
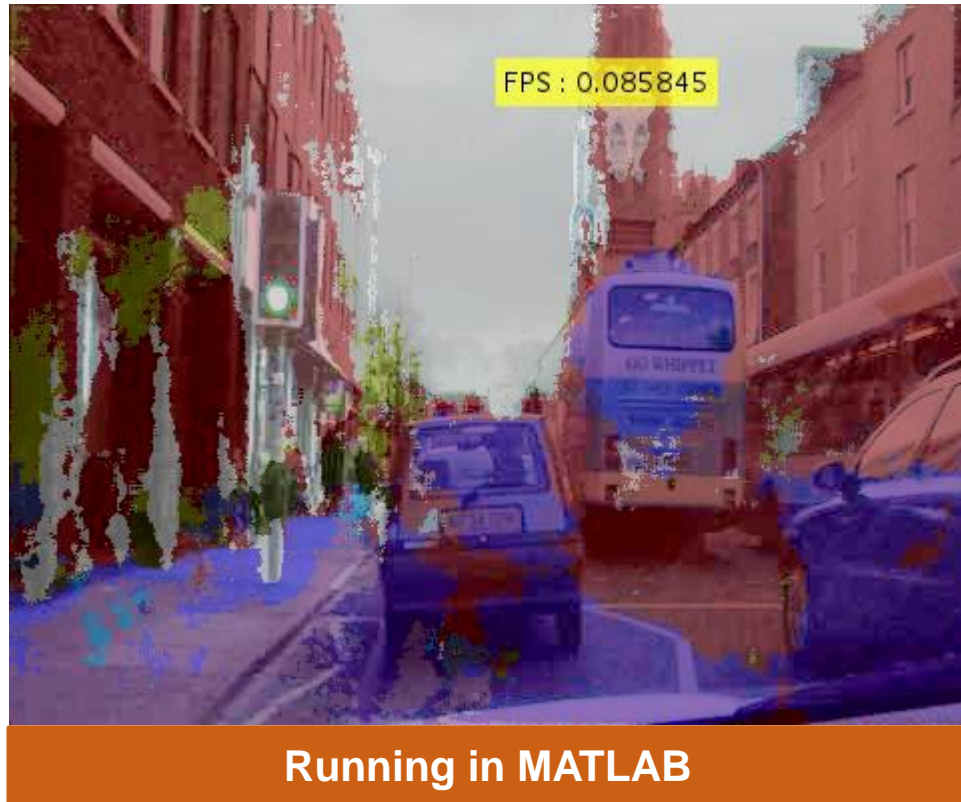
## DAGNetwork



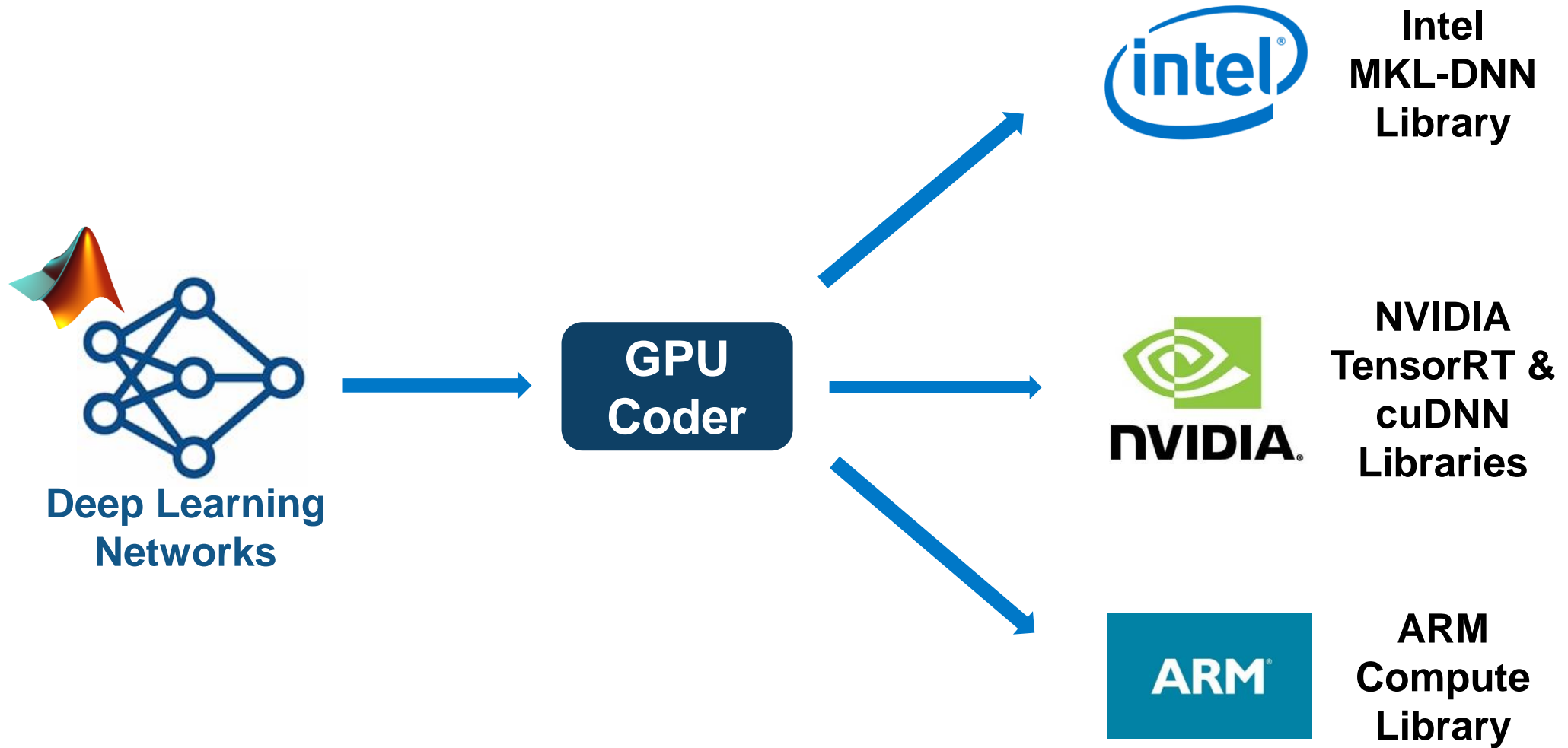
GPU Coder: **R2018a**

Networks: GoogLeNet } Object  
 ResNet } detection  
 SegNet } Semantic  
 DeconvNet } segmentation

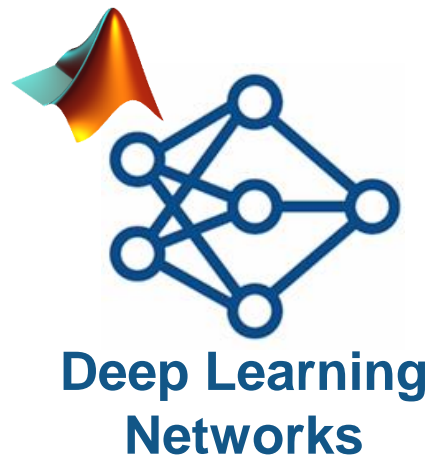
# Semantic Segmentation



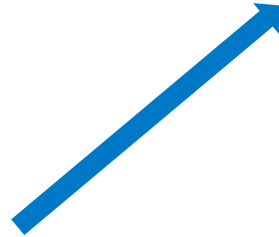
# Deploying to CPUs



# Deploying to CPUs



**GPU Coder**

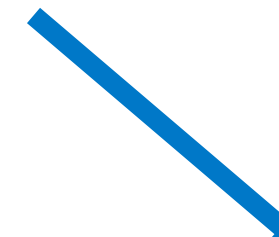


23.88 FPS
89.7% computer keyboard
8.6% space bar
1.7% typewriter keyboard
0.0% mouse
0.0% notebook

**Desktop CPU**



**NVIDIA**  
TensorRT & cuDNN Libraries



**Raspberry Pi board**

# How Good is Generated Code Performance

- Performance of image processing and computer vision
- Performance of CNN inference (Alexnet) on Titan XP GPU
- Performance of CNN inference (Alexnet) on Jetson (Tegra) TX2

# GPU Coder for Image Processing and Computer Vision



Fog removal



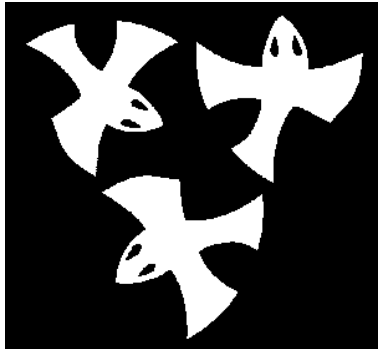
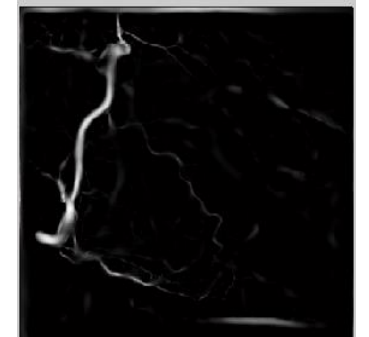
5x speedup



Frangi filter



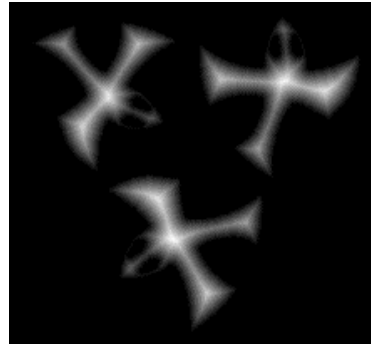
3x speedup



Distance transform



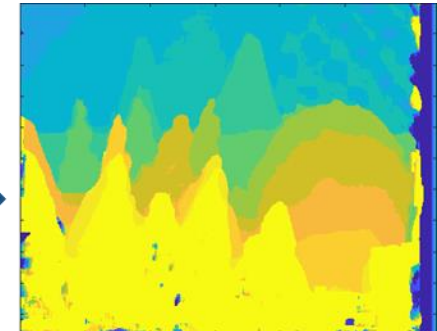
8x speedup



Stereo disparity



50x speedup



Ray tracing



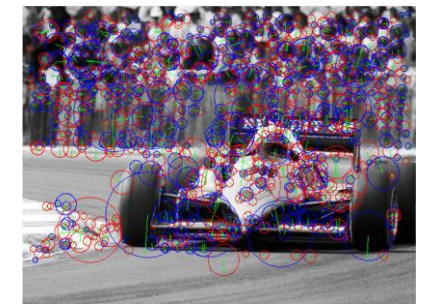
18x speedup



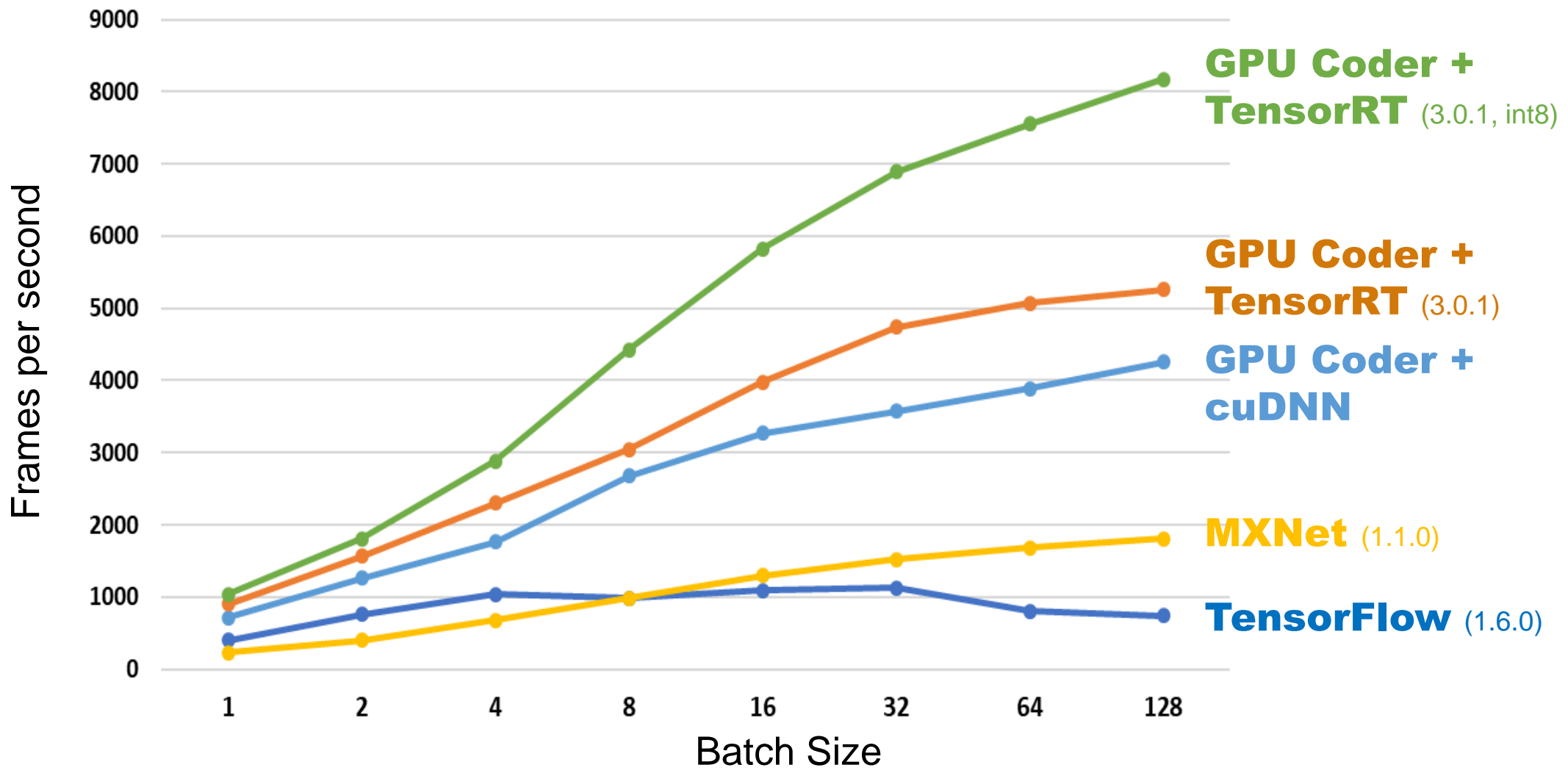
SURF feature extraction



700x speedup



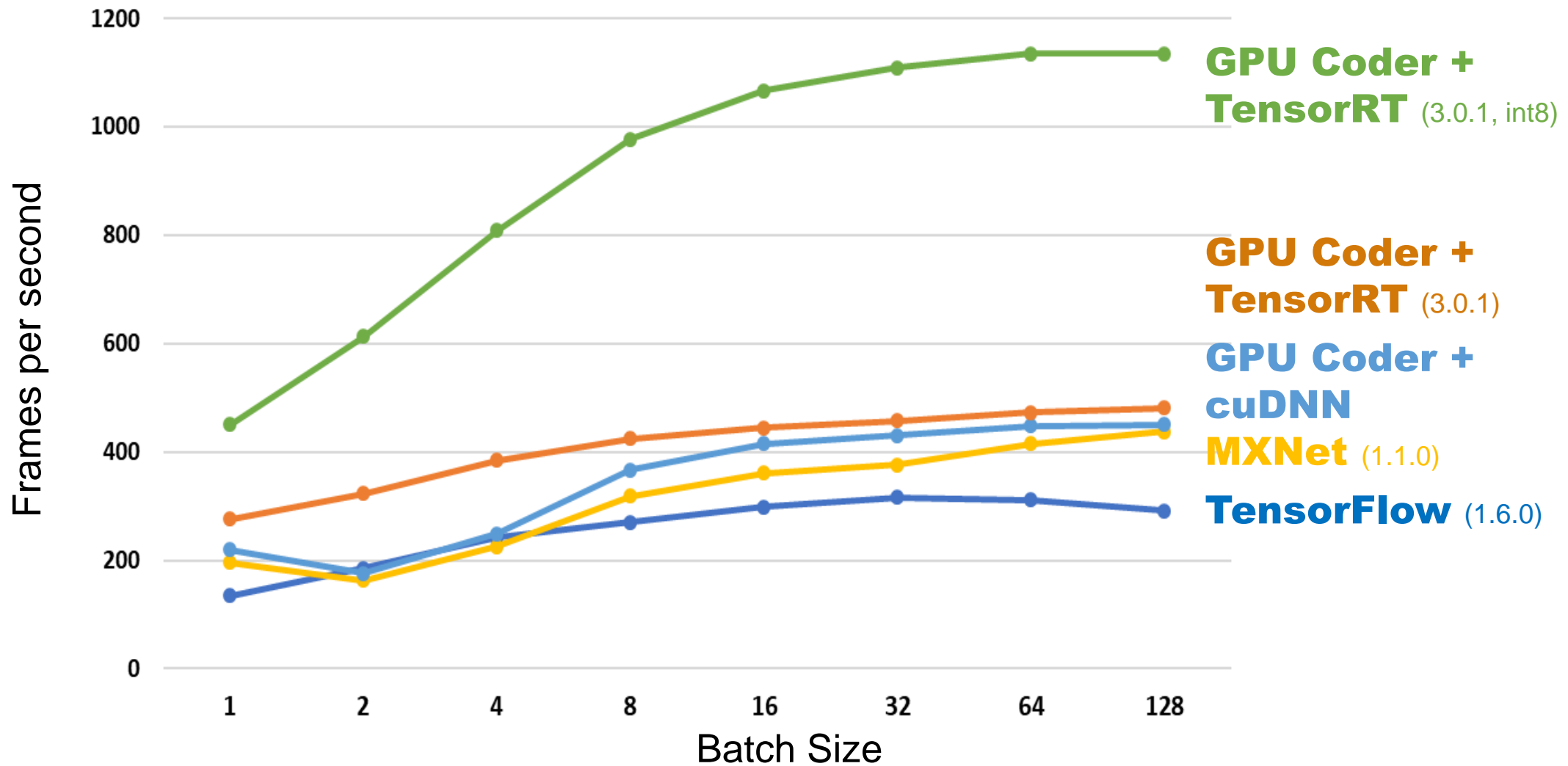
# Alexnet Inference on NVIDIA Titan Xp



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

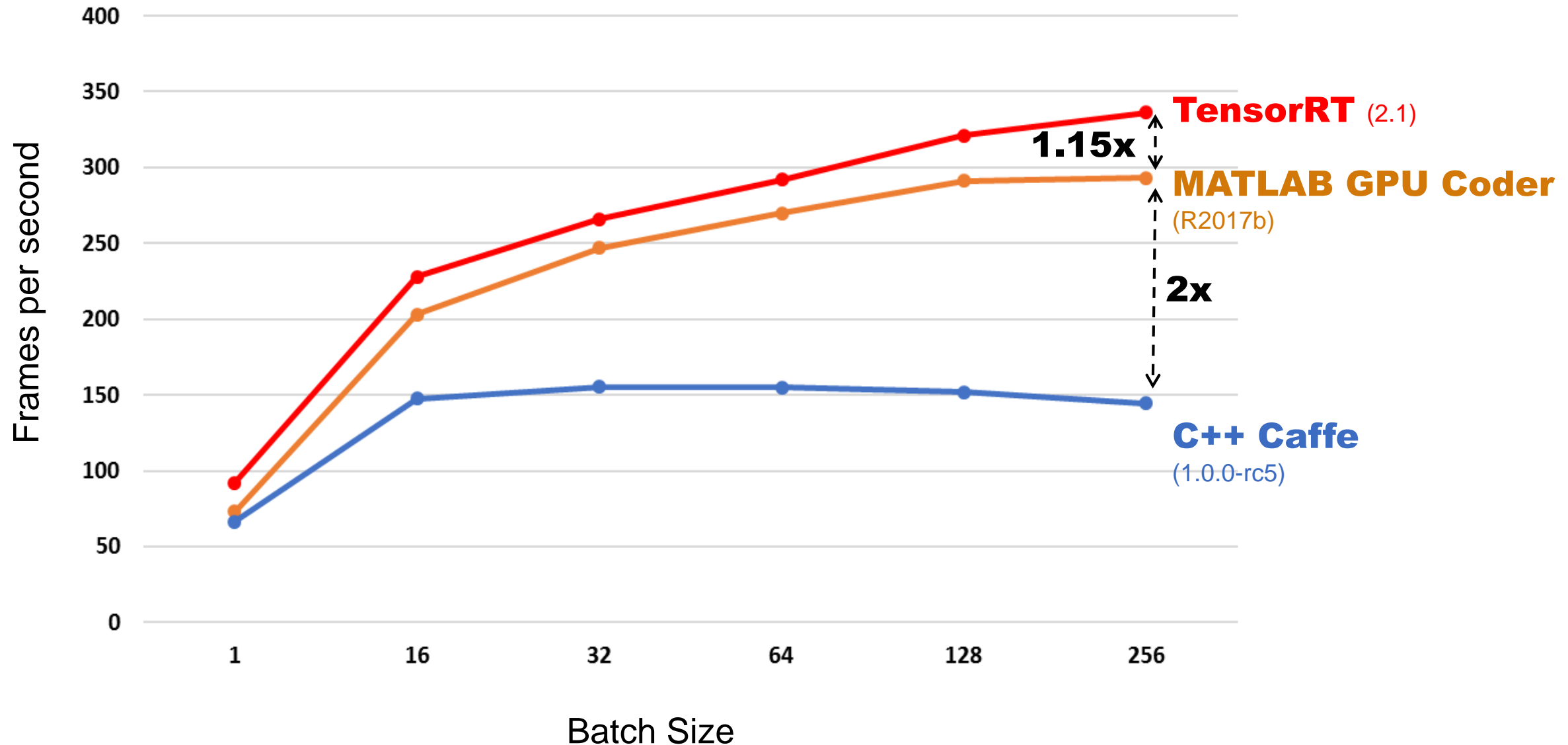


# VGG-16 Inference on NVIDIA Titan Xp

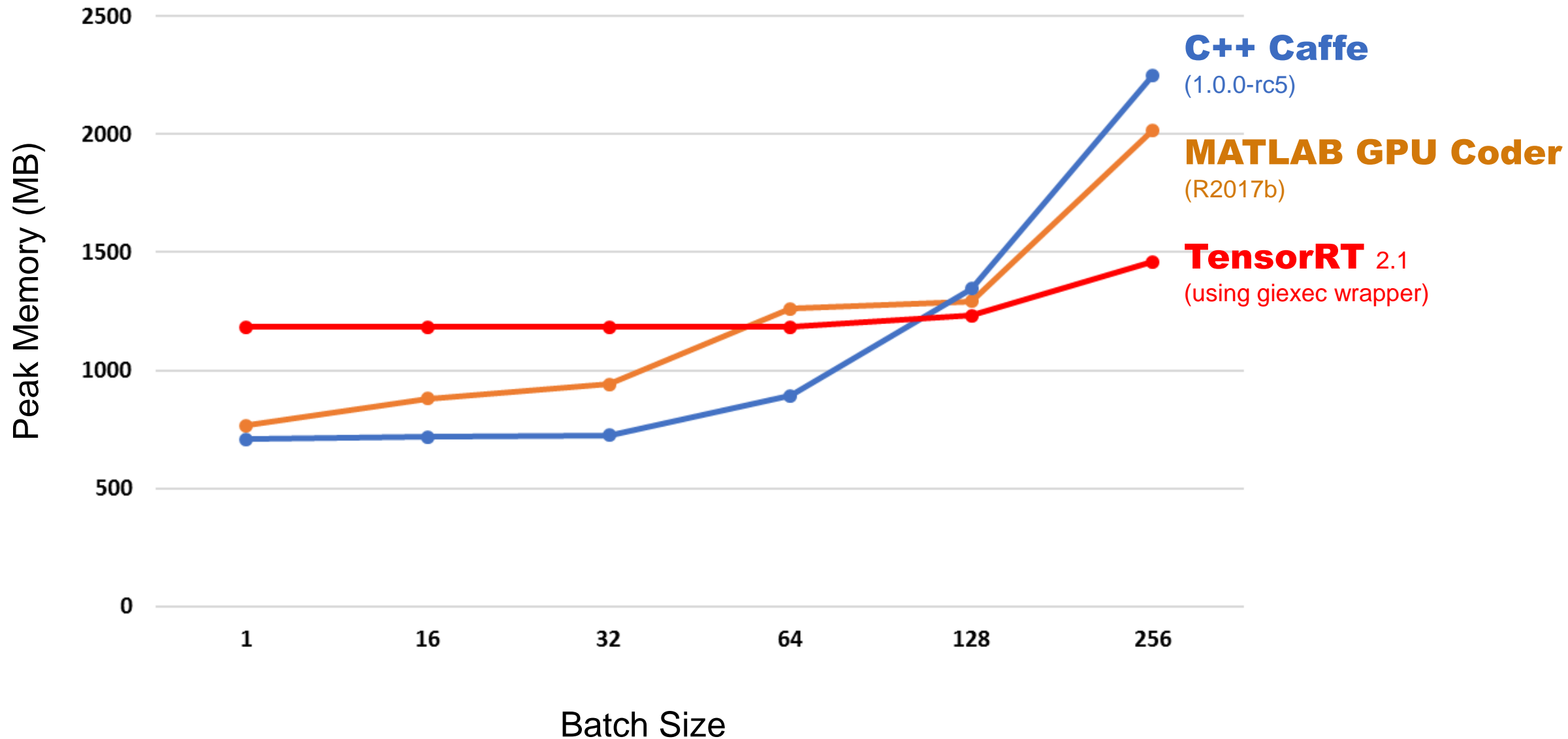


CPU	Intel(R) Xeon(R) CPU E5-1650 v4 @ 3.60GHz
GPU	Pascal Titan Xp
cuDNN	v7

# Alexnet Inference on Jetson TX2: Frame-Rate Performance



# Alexnet Inference on Jetson TX2: Memory Performance



## Key Takeaways

- GPU Coder automates the process of writing CUDA code for general algorithms – not only Deep Learning
- GPU Coder generates code for DNN for multiple platforms
- GPU Coder performs in most times better than other common Deep Learning platforms