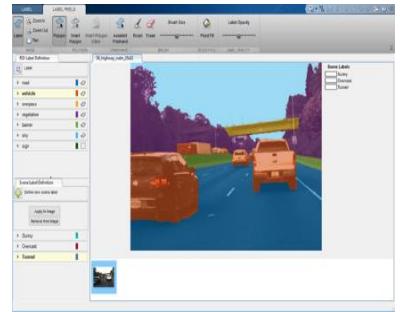


Deploying AI on Jetson Xavier/DRIVE Xavier with TensorRT and MATLAB

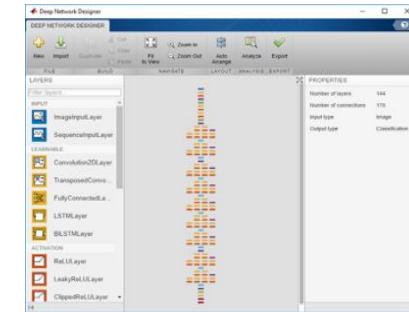
Jaya Shankar, Engineering Manager (Deep Learning Code Generation)

Avinash Nehemiah, Principal Product Manager (Computer Vision, Deep Learning, Automated Driving)

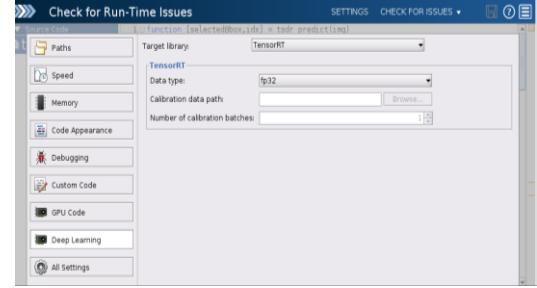
Outline



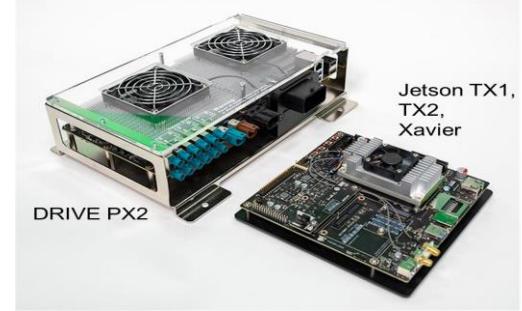
Ground Truth Labeling



Network Design and Training



CUDA and TensorRT Code Generation



Jetson Xavier and DRIVE Xavier Targeting

Key Takeaways

Platform Productivity: Workflow automation, ease of use

Framework Interoperability: ONNX, Keras-TensorFlow, Caffe

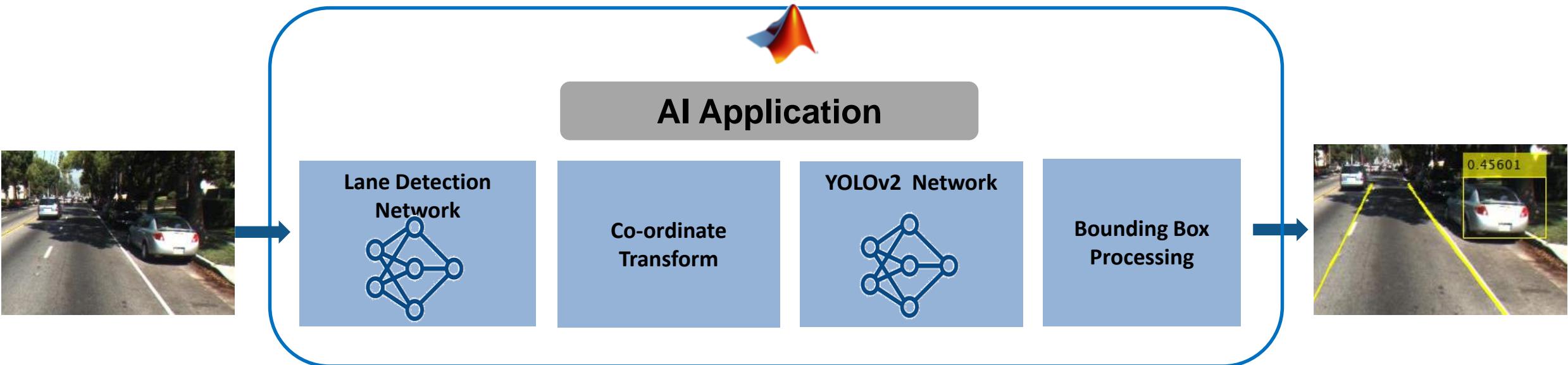
Key Takeaways

Optimized CUDA and TensorRT code generation

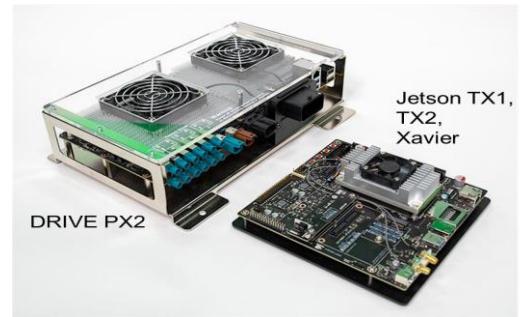
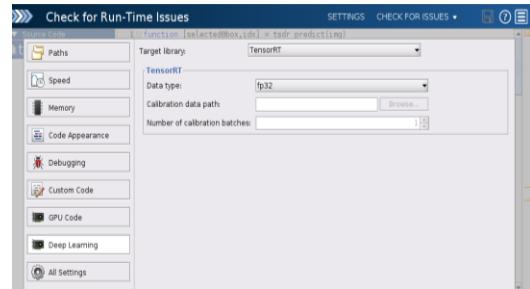
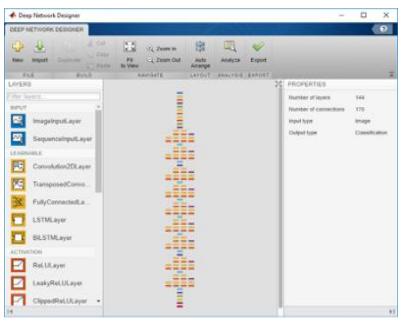
Jetson Xavier and DRIVE Xavier targeting

Processor-in-loop(PIL) testing and system integration

Example Used in Today's Talk



Outline





Unlabeled Training Data



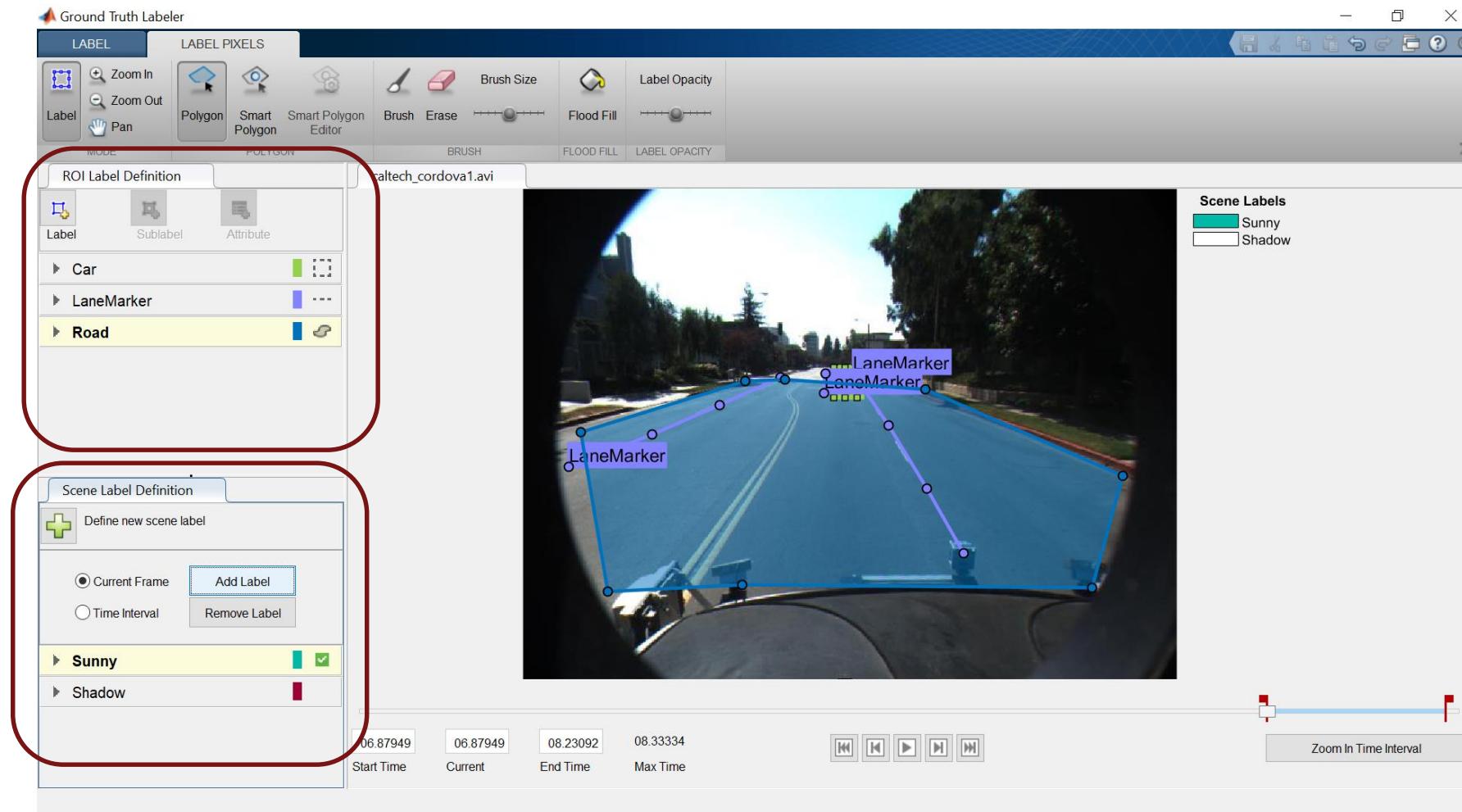
Ground Truth Labeling



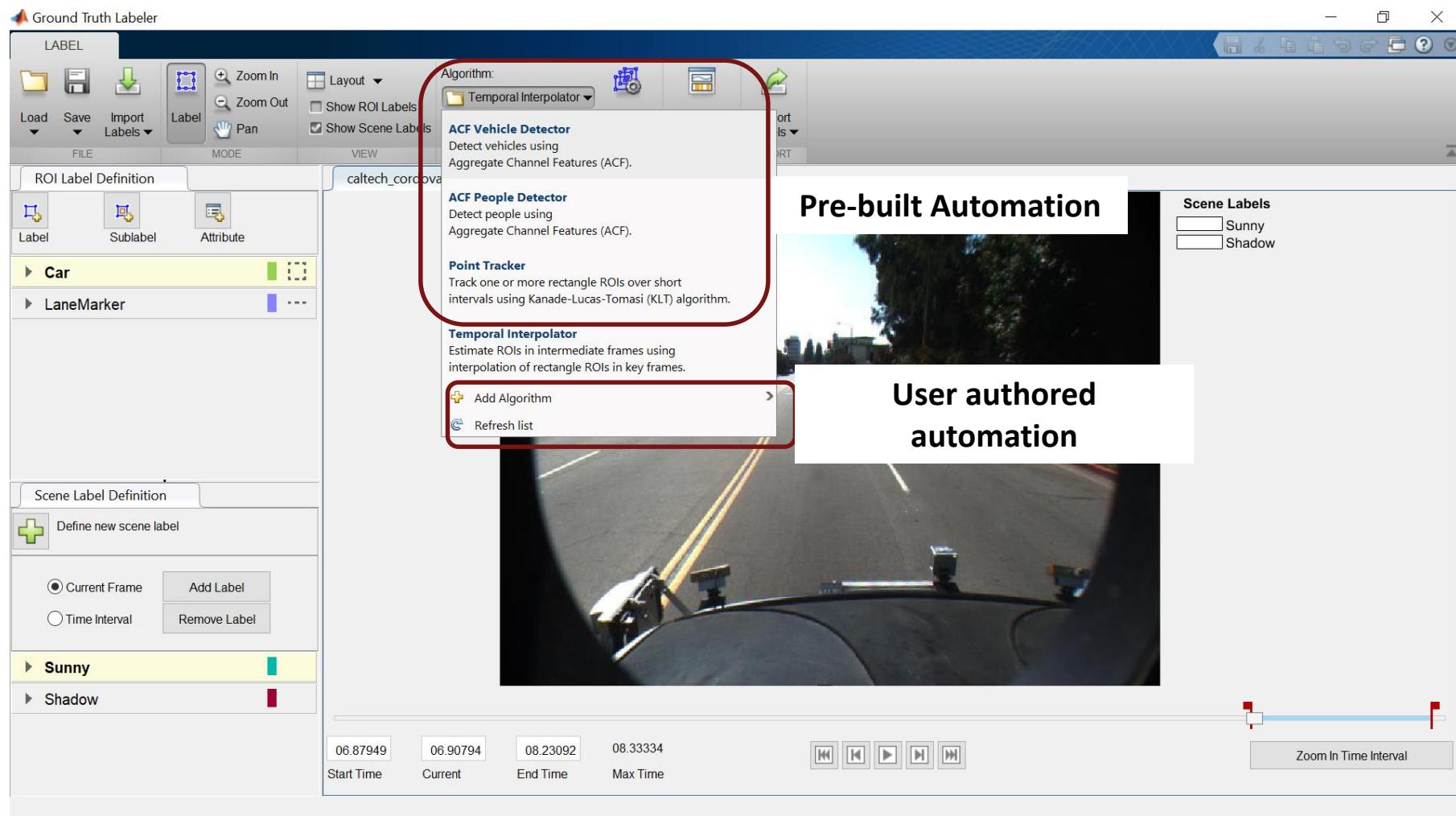
```
>> gTruth  
  
gTruth =  
  
groundTruth with properties:  
  
    DataSource: [1x1 groundTruthDataSource]  
    LabelDefinitions: [4x3 table]  
    LabelData: [250x4 timetable]  
  
>> gTruth.LabelData  
  
ans =  
  
250x4 timetable  
  
Time           Car            LaneMarker        Sunny      Shadow  
_____|_____|_____|_____|_____|_____|  
0 sec          [2x4 double]   {2x1 cell}    true     false  
0.033333 sec  [2x4 double]   {2x1 cell}    true     false  
0.066667 sec  []             []           false    false
```

Labels for Training

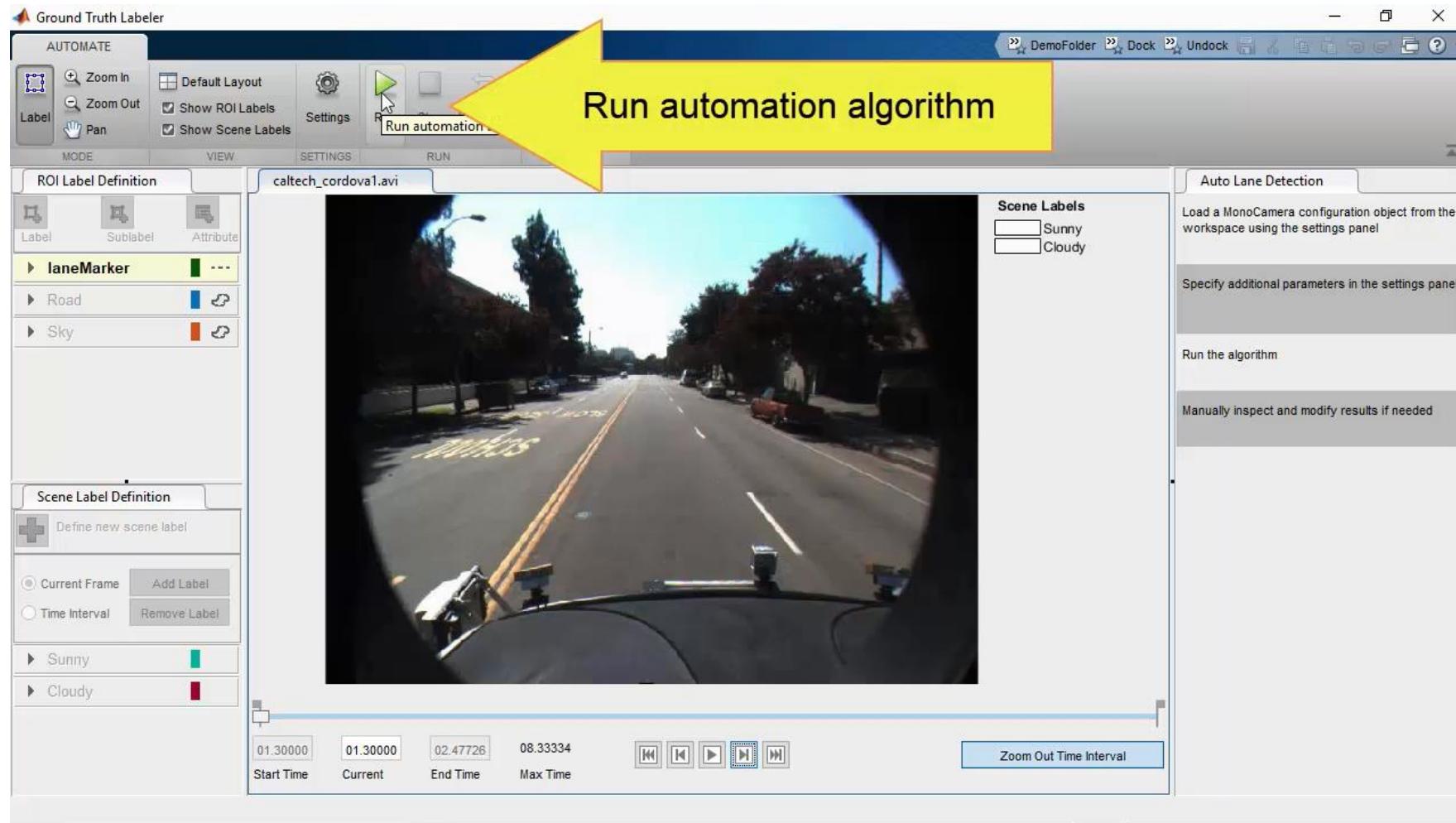
Interactive Tools for Ground Truth Labeling



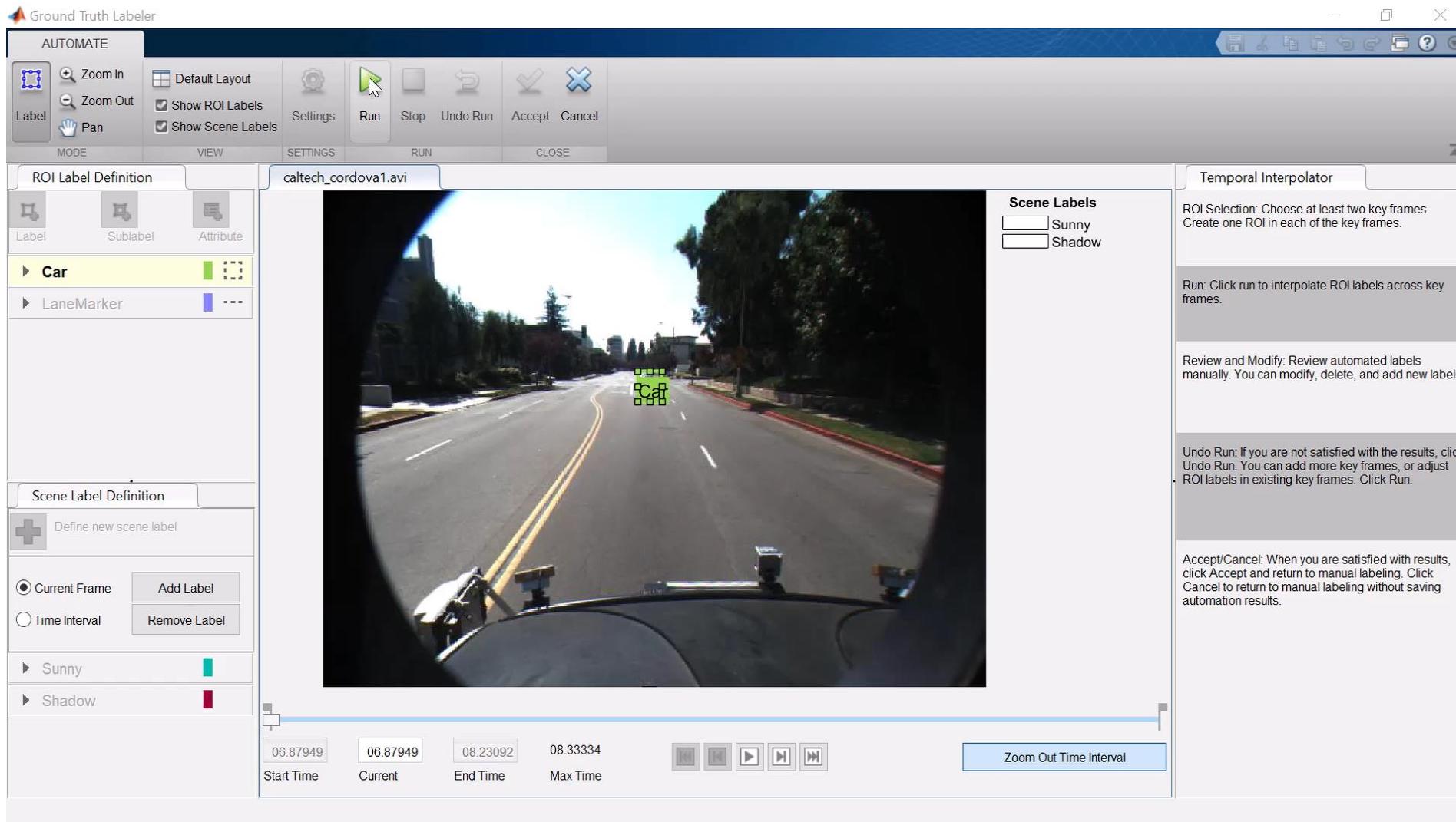
Automate Ground Truth Labeling



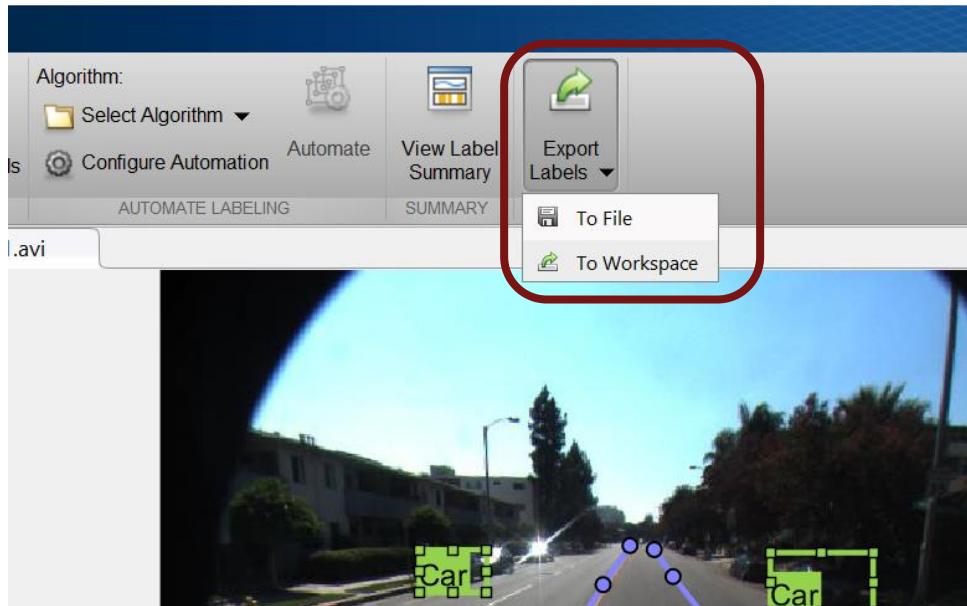
Automating Labeling of Lane Markers



Automate Labeling of Bounding Boxes for Vehicles

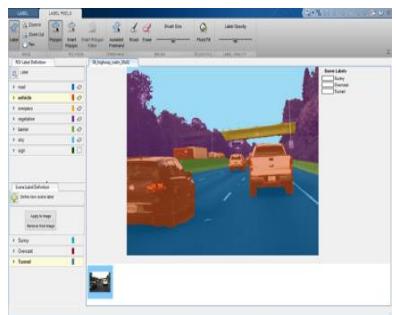


Export Labeled Data for Training



```
>> gTruth  
  
gTruth =  
  
groundTruth with properties:  
  
    DataSource: [1x1 groundTruthDataSource]  
    LabelDefinitions: [4x3 table]  
    LabelData: [250x4 timetable]  
  
>> gTruth.LabelData  
  
ans =  
  
250x4 timetable  
  
Time           Car        LaneMarker  Sunny  Shadow  
_____         _____      _____     ____  ____  
  
0 sec          [2x4 double] {2x1 cell} true   false  
0.033333 sec  [2x4 double] {2x1 cell} true   false  
0.066667 sec  []          []          false  false  
  
↑  
Bounding Boxes Labels  
  
↑  
Polyline Labels
```

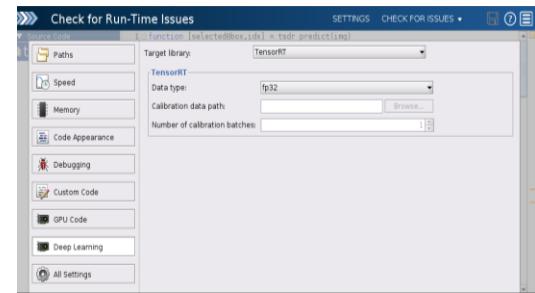
Outline



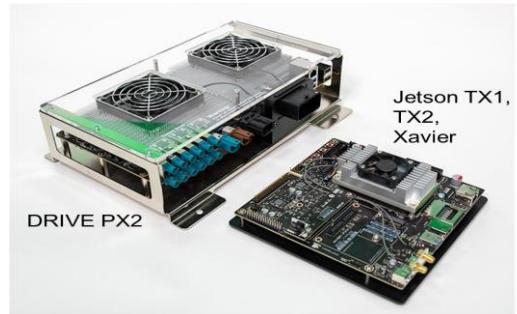
Ground Truth Labeling



Network Design and Training

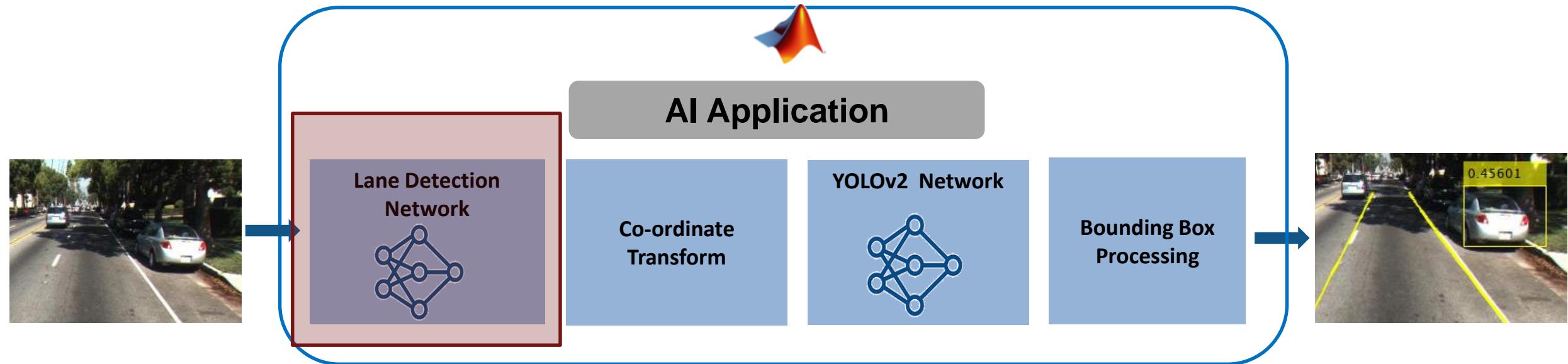


CUDA and TensorRT Code Generation

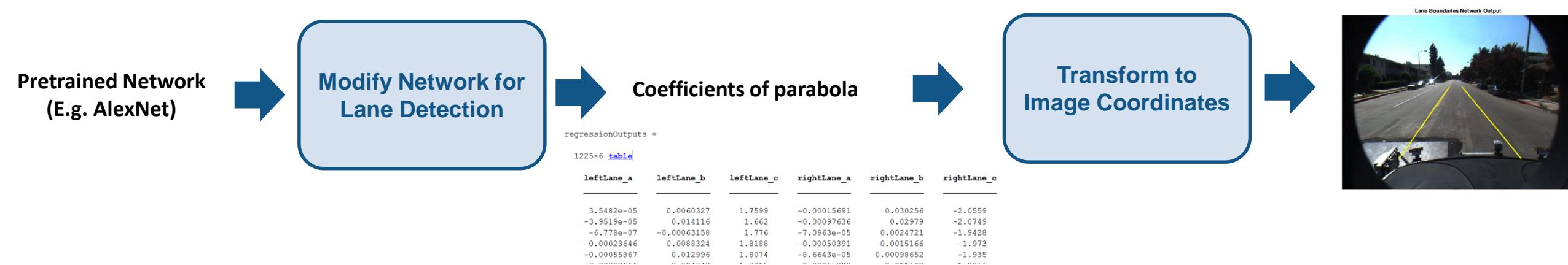


Jetson Xavier and DRIVE
Xavier Targeting

Example Used in Today's Talk



Lane Detection Algorithm



Lane Detection: Load Pretrained Network

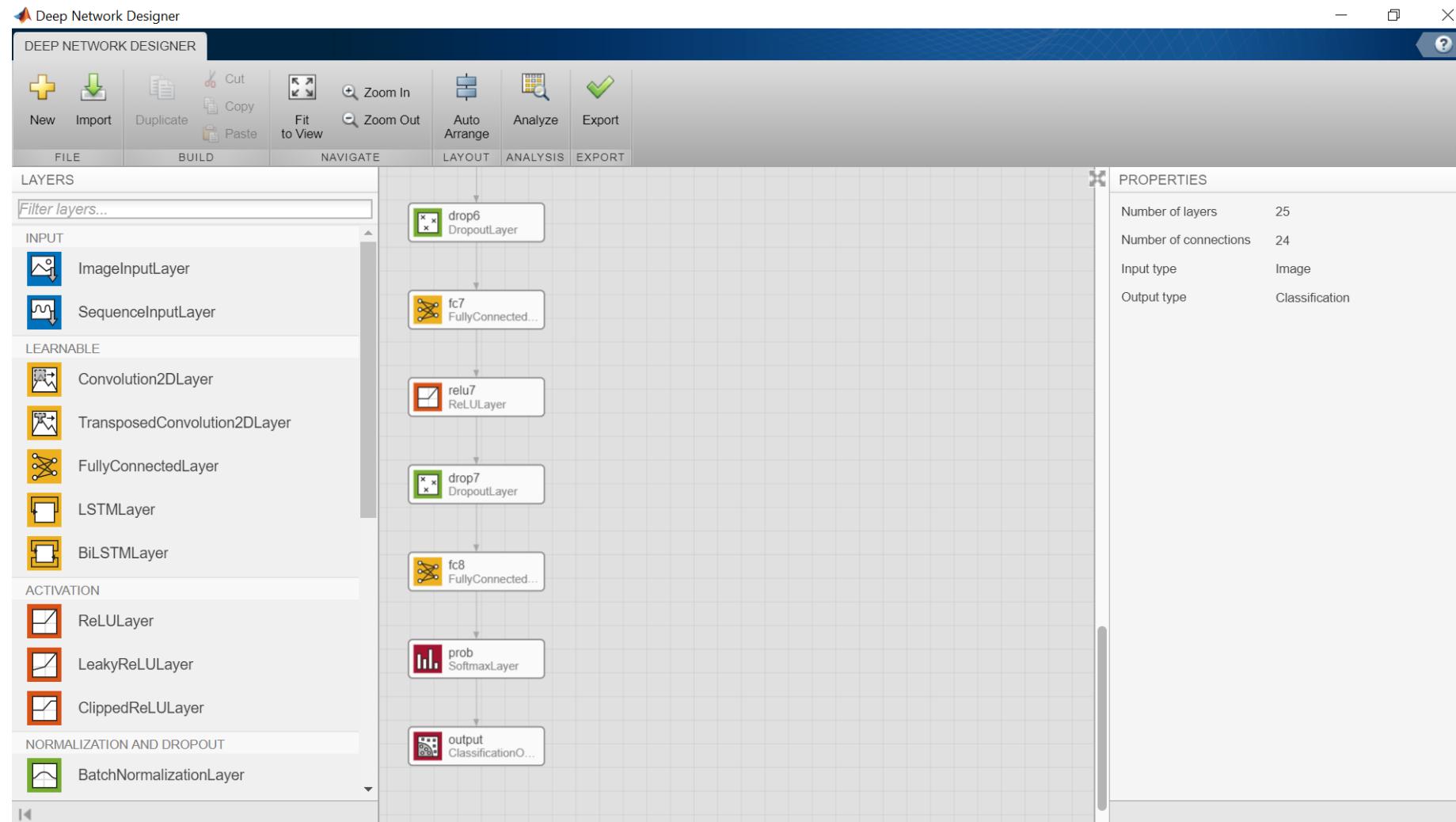


```
>> net = alexnet  
>> deepNetworkDesigner
```

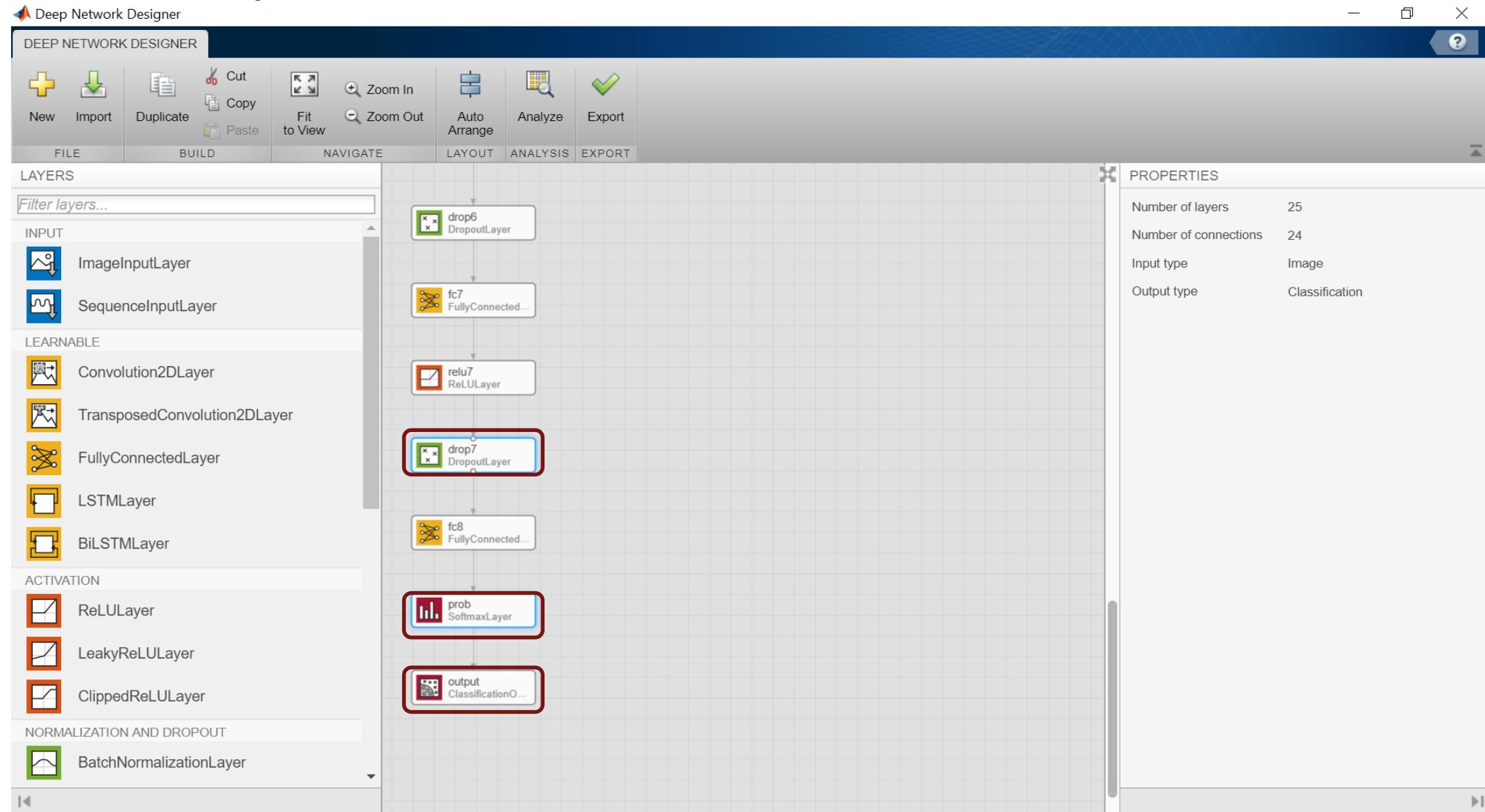
Lane Detection Network

- Regression CNN for lane parameters
- MATLAB code to transform to image co-ordinates

View Network in Deep Network Designer App



Remove Layers from AlexNet



Add Regression Output for Lane Parameters

The screenshot shows the Deep Network Designer app interface. On the left, there's a toolbar with 'New', 'Import', 'Duplicate', 'Cut', 'Copy', 'Paste', 'Fit to View', 'Zoom In', 'Zoom Out', 'Auto Arrange', 'Analyze', and 'Export' buttons. Below the toolbar are sections for 'FILE', 'BUILD', 'NAVIGATE', 'LAYOUT', 'ANALYSIS', and 'EXPORT'. The 'LAYER' section contains a dropdown menu and a list of learnable layers: Convolution2DLayer, TransposedConvolution2DLayer, FullyConnectedLayer, LSTMLayer, BiLSTMLayer. The 'ACTIVATION' section lists ReLU, LeakyReLU, and ClippedReLU layers. The 'NORMALIZATION AND DROPOUT' section includes BatchNormalizationLayer, CrossChannelNormalizationLayer, and DropoutLayer. The main workspace displays a neural network diagram with the following layers connected sequentially: drop6 (DropoutLayer), fcLane1 (FullyConnectedLayer), relu7 (ReLUlayer), fcLane1Relu (FullyConnectedLayer), relu (ReLUlayer), fcLane2 (FullyConnectedLayer), and regressionout... (RegressionOut...). The 'regressionout...' layer is highlighted with a red border. To the right of the workspace is a 'PROPERTIES' panel showing the following details:

PROPERTIES	
Number of layers	25
Number of connections	24
Input type	Image
Output type	Regression

Regression Output for Lane Coefficients

Transparently Scale Compute for Training

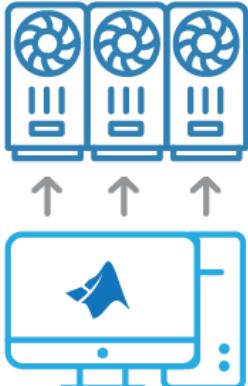
Specify Training on:



'CPU'



'gpu'



'multi-gpu'

Works on Windows
(no additional setup)

Quickly change training hardware

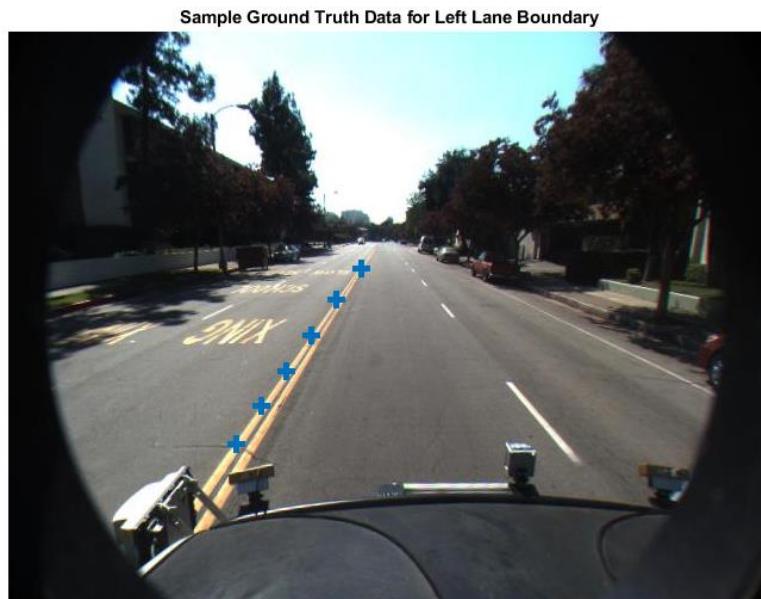
```
opts = trainingOptions('sgdm', ...
    'BatchSize', 100, ...
    'MinibatchSize', 250, ...
    'InitialLearnRate', 0.00005, ...
    'ExecutionEnvironment', 'auto');
```

NVIDIA NGC & DGX Supports MATLAB for Deep Learning

- GPU-accelerated MATLAB Docker container for deep learning
 - Leverage multiple GPUs on NVIDIA DGX Systems and in the Cloud
 - Cloud providers include: AWS, Azure, Google, Oracle, and Alibaba
- NVIDIA DGX System / Station
 - Interconnects 4/8/16 Volta GPUs in one box
- Containers available for R2018a and R2018b
 - New Docker container with every major release (a/b)
- Download MATLAB container from NGC Registry
 - <https://ngc.nvidia.com/registry/partners-matlab>



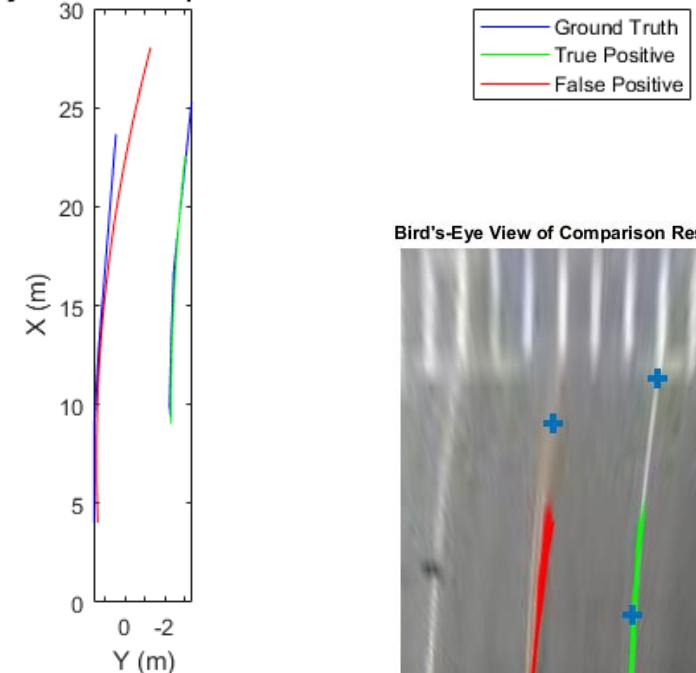
Evaluate Lane Boundary Detections vs. Ground Truth



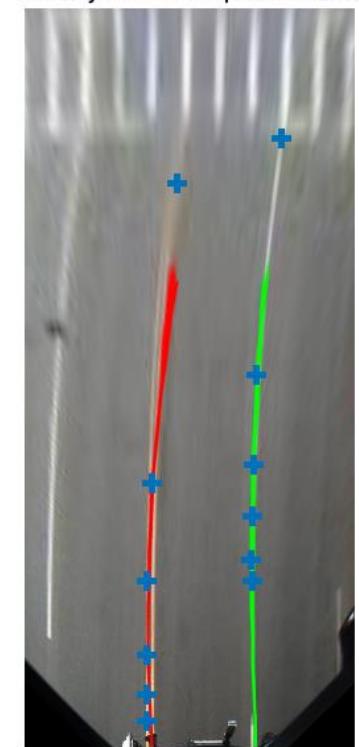
`evaluateLaneBoundaries`



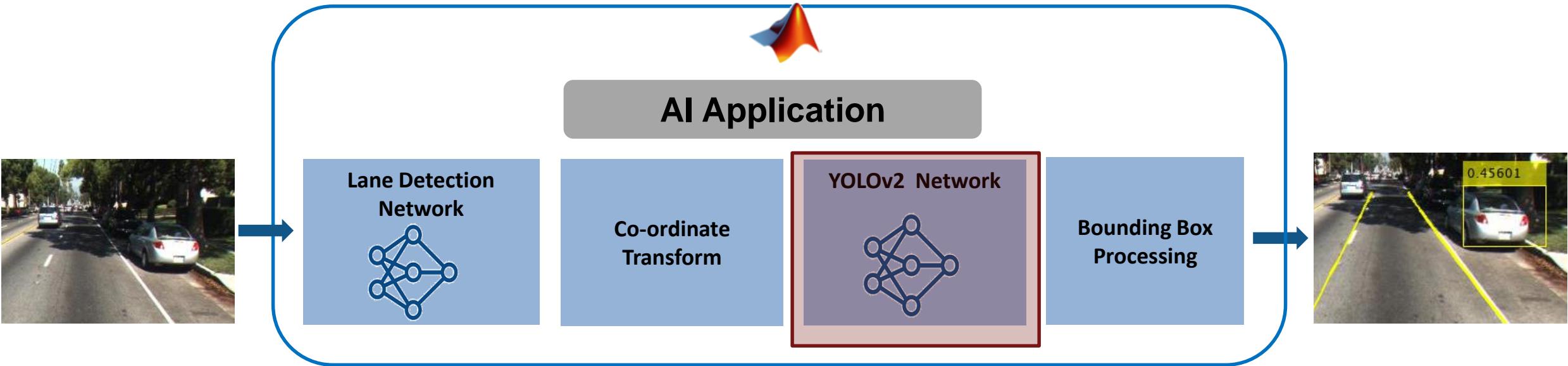
Bird's-Eye Plot of Comparison Results



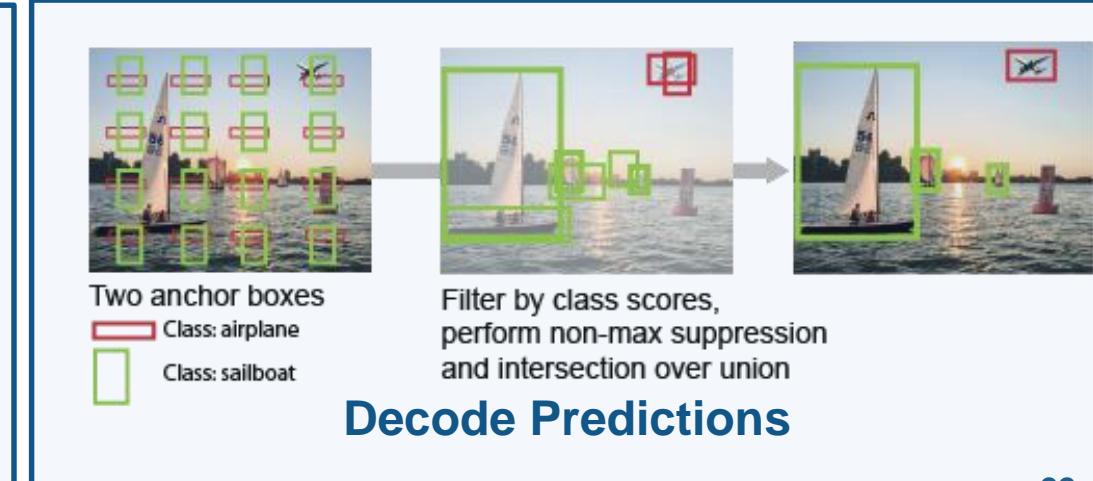
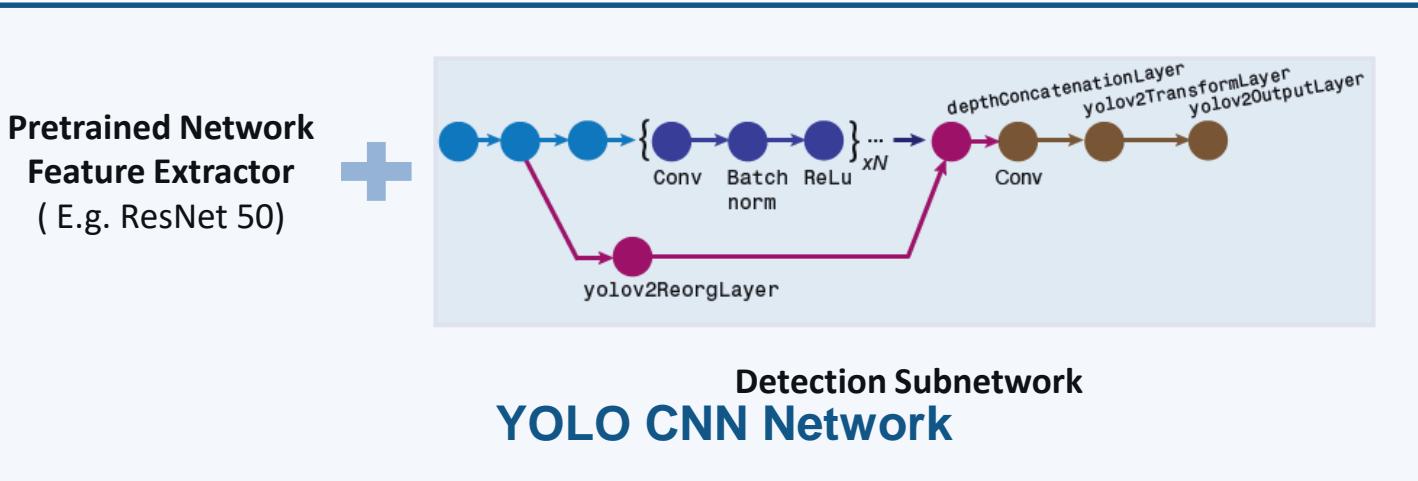
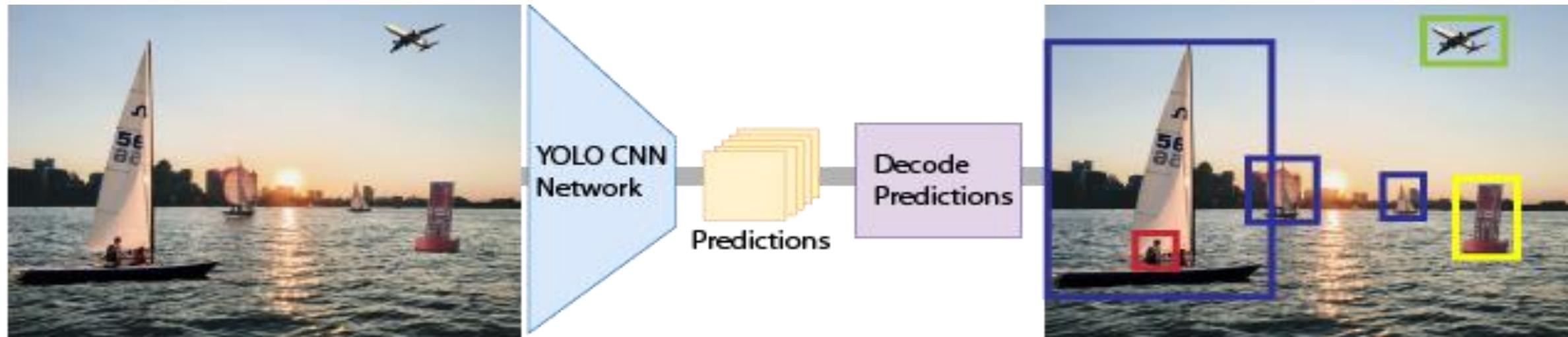
Bird's-Eye View of Comparison Results



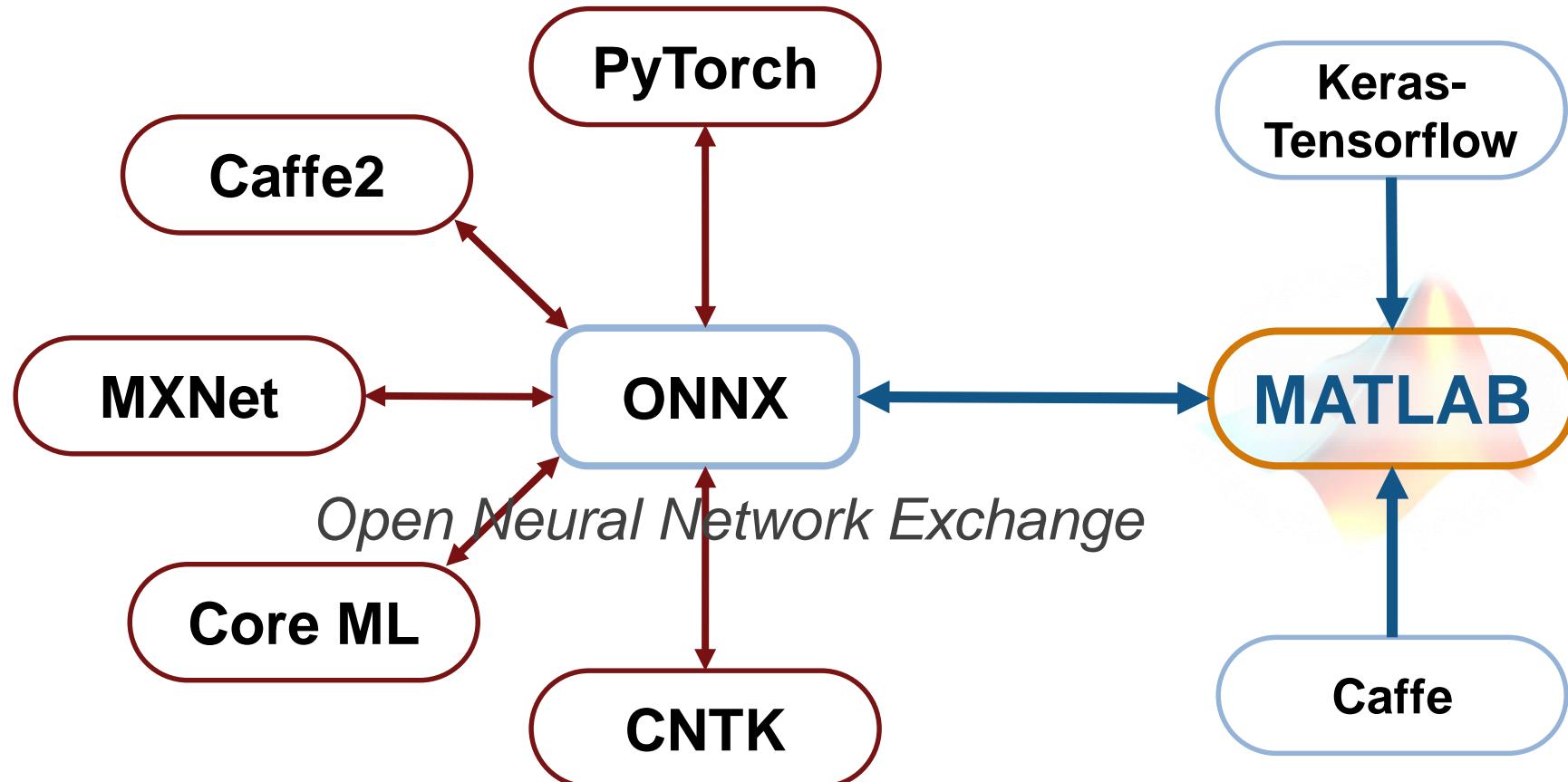
Example Used in Today's Talk



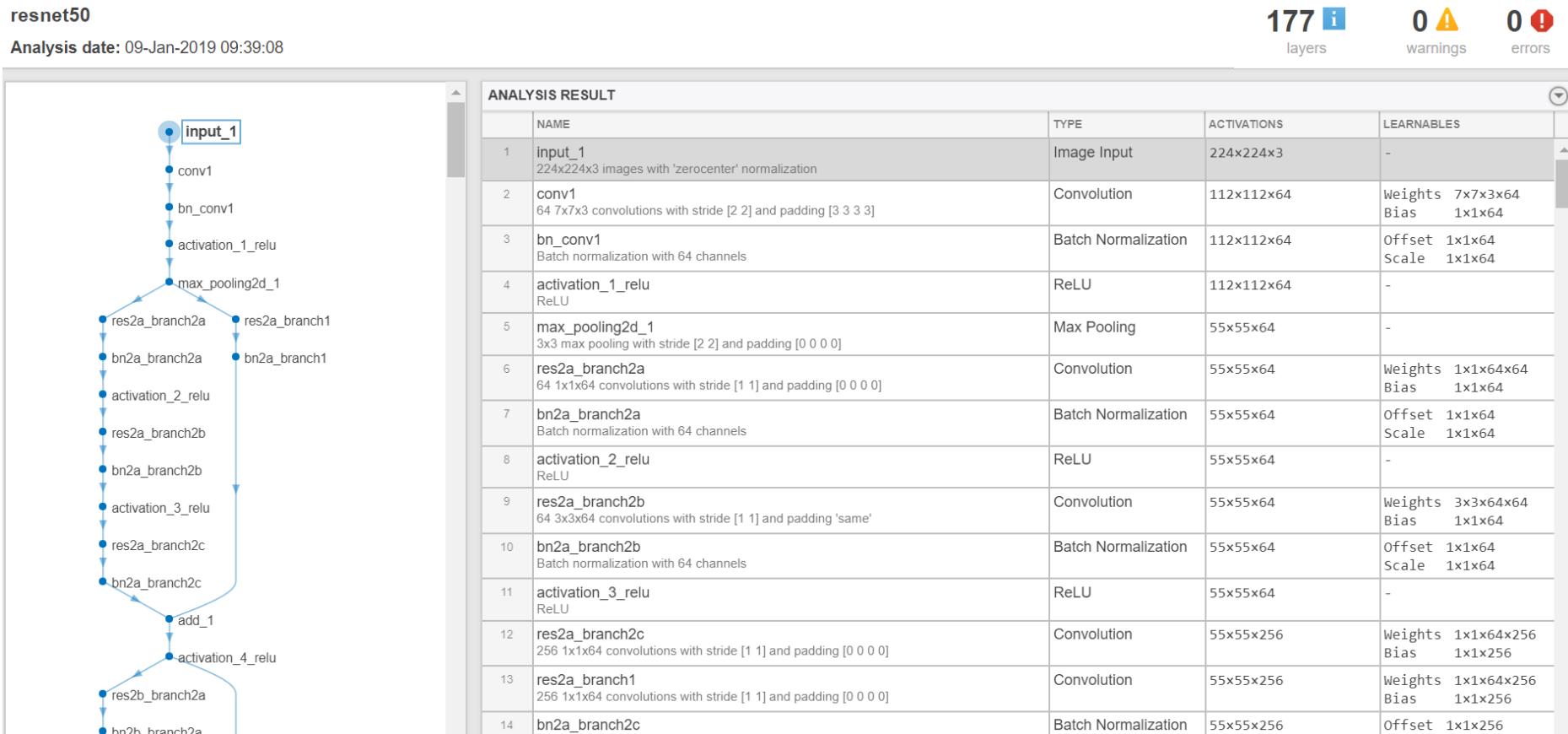
YOLO v2 Object Detection



Model Exchange with MATLAB



Import Pretrained Network in ONNX Format



```
load resnetClassNames.mat
net = importONNXNetwork('resnet50.onnx', ...
    'OutputLayerType', 'classification', ...
    'ClassNames', classnames);
analyzeNetwork(net)
```

Modify Network

```
lgraph = layerGraph(net);
lgraph = removeLayers(lgraph, 'Input_input_1');
lgraph = removeLayers(lgraph, 'fc1000_Flatten1');
lgraph = connectLayers(lgraph, 'avg_pool', 'fc1000');

avgImgBias = -1*(lgraph.Layers(1).Bias);

%Create new input layer and incorporate average image bias
larray = imageInputLayer([224 224 3], ...
    'Name','input',...
    'AverageImage',avgImgBias);

lgraph = replaceLayer(lgraph, 'input_1_Sub',larray);

netModified = assembleNetwork(lgraph);

save('resnet50_model.mat','netModified');
```

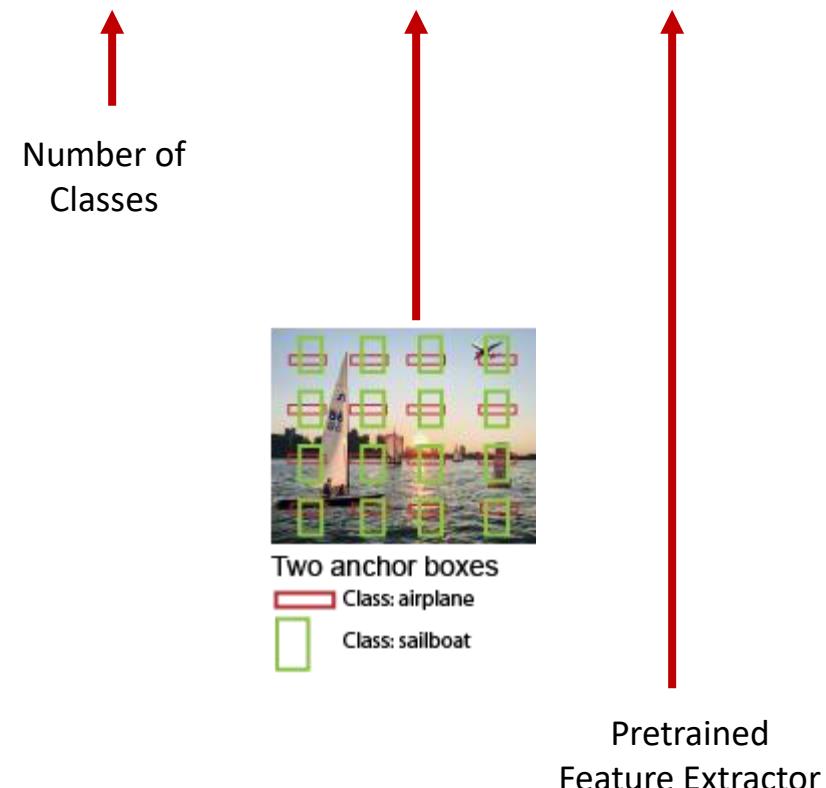
Removing the 2
ResNet-50 layers

imageInputLayer replaces the
input and subtraction layer

YOLOv2 Detection Network

- **yolov2Layers:** Create network architecture

```
>> lgraph = yolov2Layers(imageSize, numClasses, anchorBoxes, network, featureLayer)
```

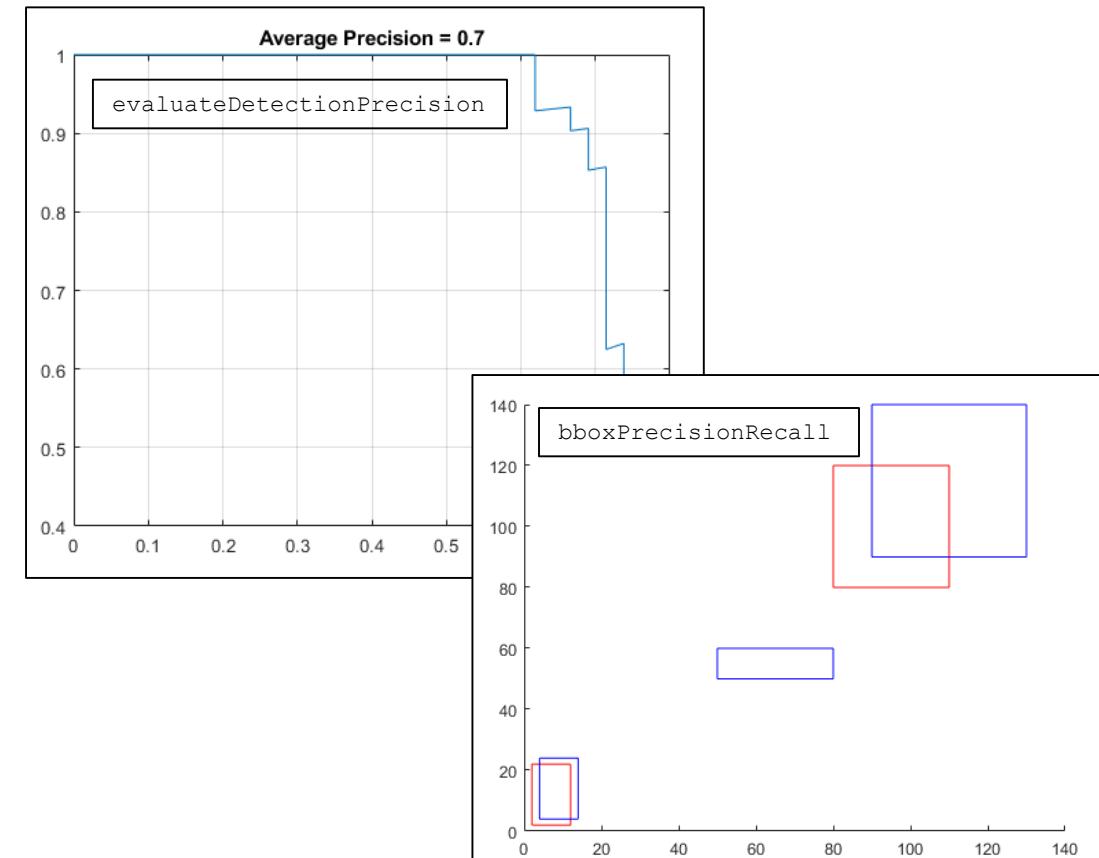


```
>> detector = trainYOLov2ObjectDetector(trainingData,lgraph,options)
```

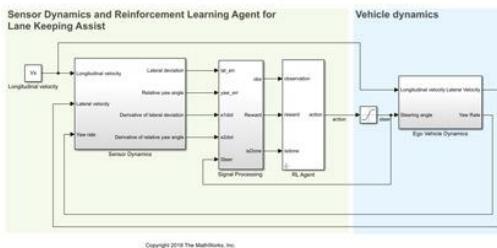
Evaluate Performance of Trained Network

- **Set of functions** to evaluate trained network performance
 - evaluateDetectionMissRate
 - **evaluateDetectionPrecision**
 - bboxPrecisionRecall
 - bboxOverlapRatio

```
>> [ap,recall,precision] =  
evaluateDetectionPrecision(results,vehicles(:,2));
```



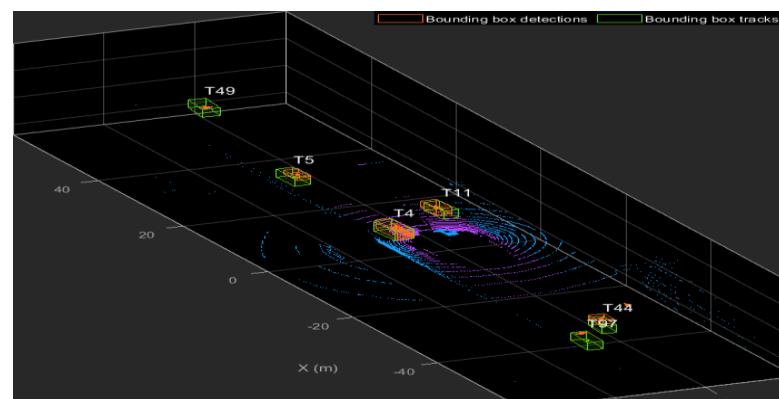
Example Applications using MATLAB for AI Development



Lane Keeping Assist using Reinforcement Learning

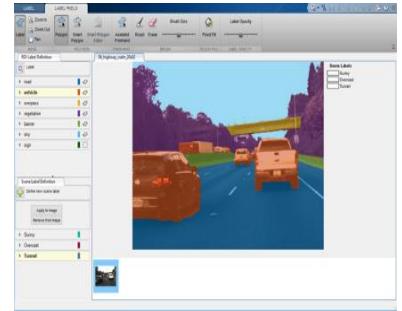


Occupancy Grid Creation using Deep Learning

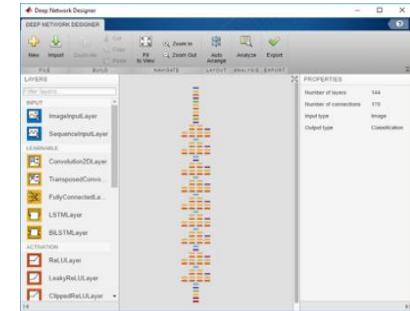


Lidar Segmentation with Deep Learning

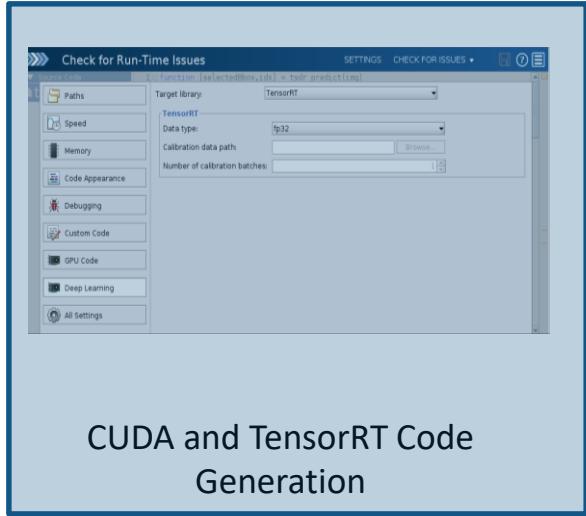
Outline



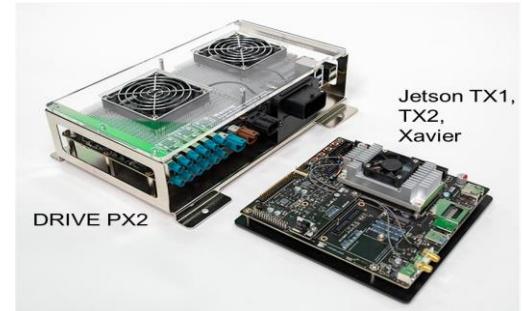
Ground Truth Labeling



Network Design and Training



CUDA and TensorRT Code Generation



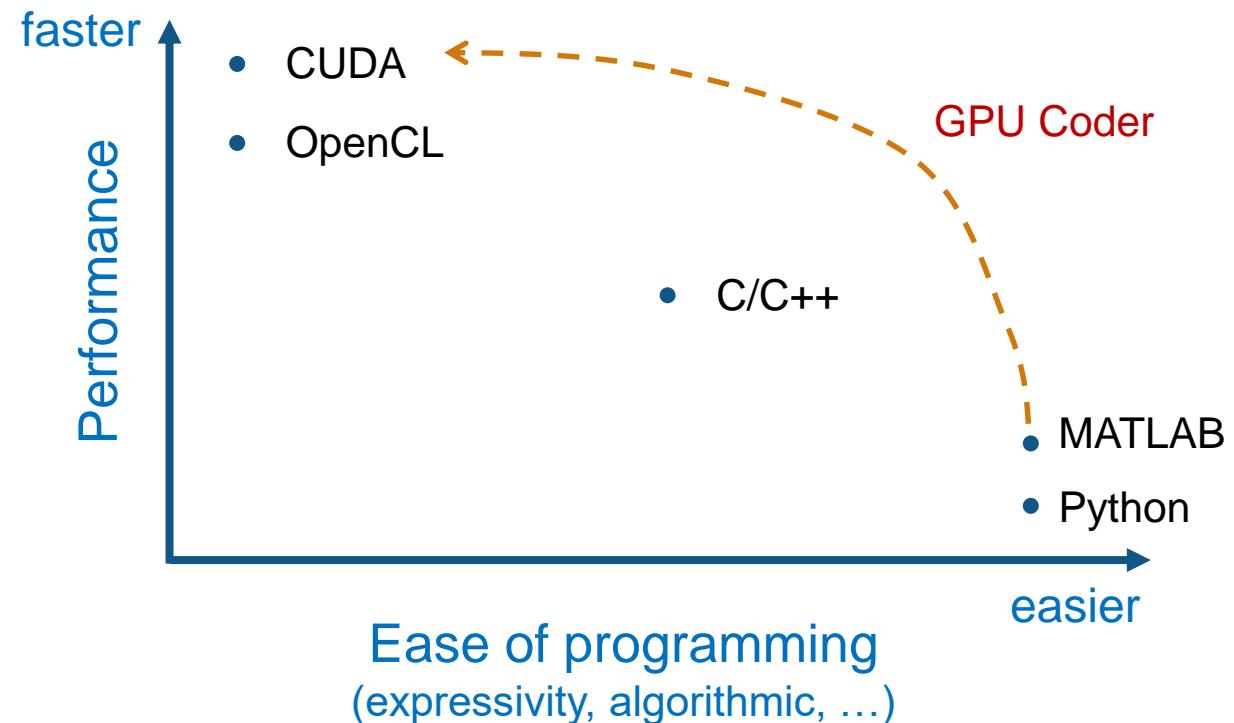
Jetson Xavier and DRIVE
Xavier Targeting

Key Takeaways

Platform Productivity: Workflow automation, ease of use

Framework Interoperability: ONNX, Keras-TensorFlow, Caffe

CUDA code generation with GPU Coder



GPUs are “hardware on steroids”, ... but, programming them is hard

Consider an example: saxpy

Simple MATLAB Loop

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



GPU Coder



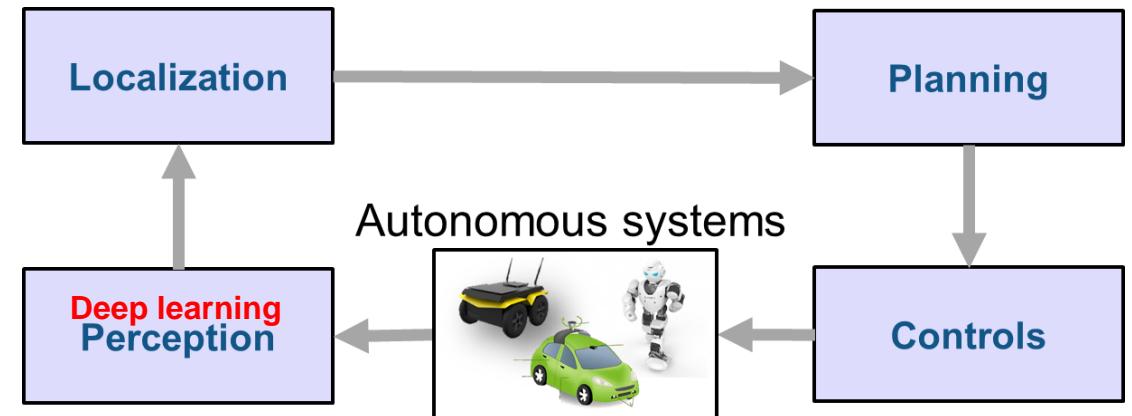
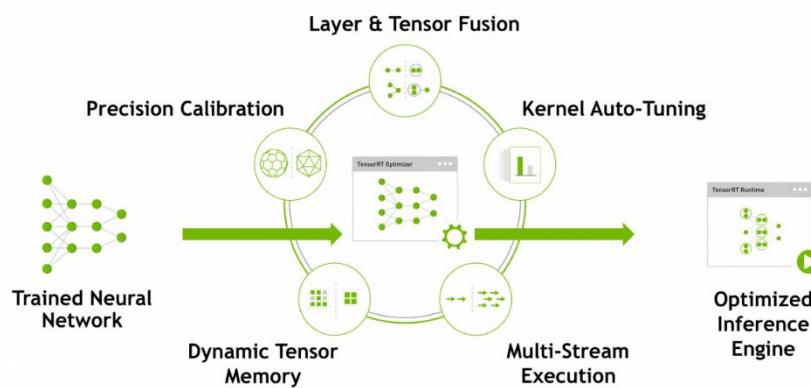
```
static __global__ launch_bounds_(512, 1) void saxpy_kernel1(const real32_T *y,
    const real32_T *x, real32_T a, real_T *z)
{
    int i = blockIdx.x*blockDim.x + threadIdx.x;
    if (!(i >= 1048576)) {
        z[i] = (real_T)(a * x[i] + y[i]);
    }
}

void saxpy(real32_T a, const real32_T x[1048576], const real32_T y[1048576],
           real_T z[1048576])
{
    real32_T *gpu_y;
    real32_T *gpu_x;
    real_T *gpu_z;
    cudaMalloc(&gpu_z, 8388608UL);
    cudaMalloc(&gpu_x, 4194304UL);
    cudaMalloc(&gpu_y, 4194304UL);
    cudaMemcpy((void *)gpu_y, (void *)&y[0], 4194304UL, cudaMemcpyHostToDevice);
    cudaMemcpy((void *)gpu_x, (void *)&x[0], 4194304UL, cudaMemcpyHostToDevice);
    saxpy_kernel1<<<dim3(2048U, 1U, 1U), dim3(512U, 1U, 1U)>>>(gpu_y, gpu_x, a,
        gpu_z);
    cudaMemcpy((void *)&z[0], (void *)gpu_z, 8388608UL, cudaMemcpyDeviceToHost);
    cudaFree(gpu_y);
    cudaFree(gpu_x);
    cudaFree(gpu_z);
}
```

Automatic compilation from an expressive language to a high-performance language

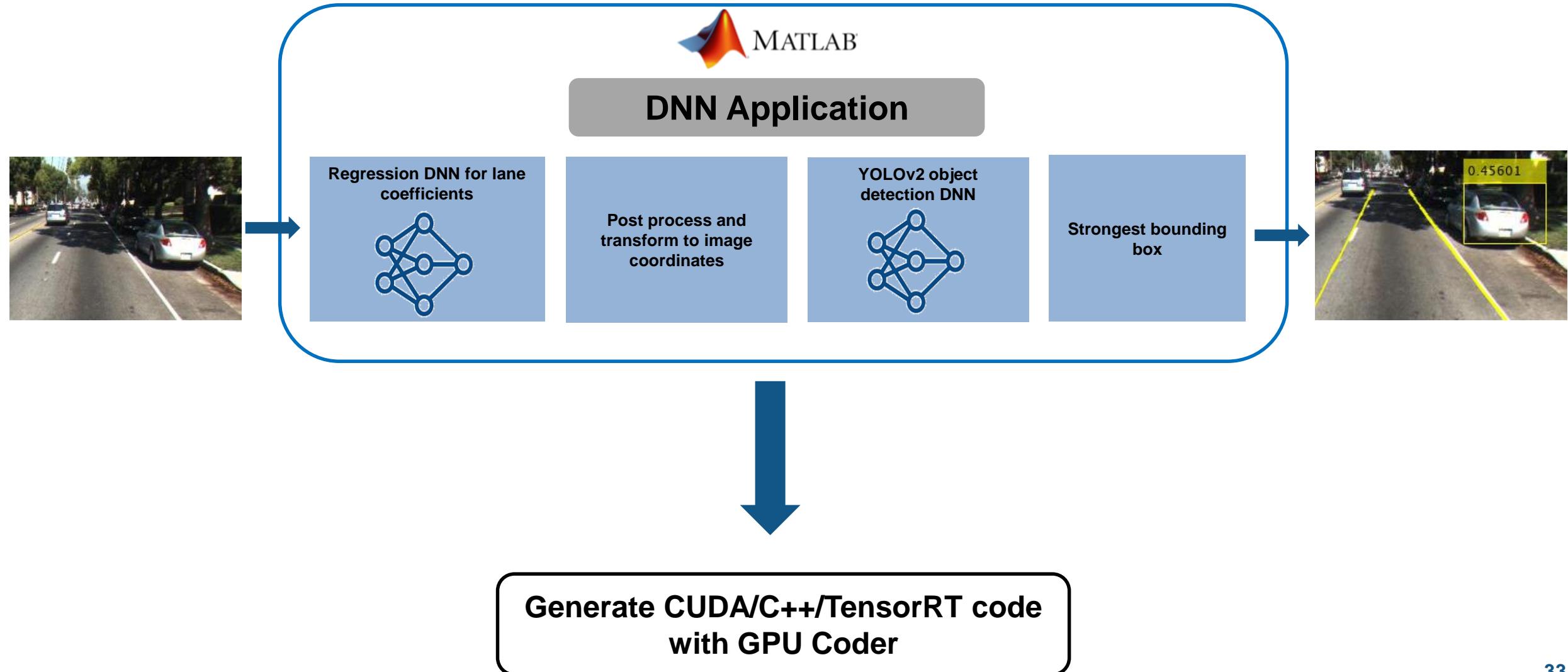
Deep Learning code generation with GPU Coder

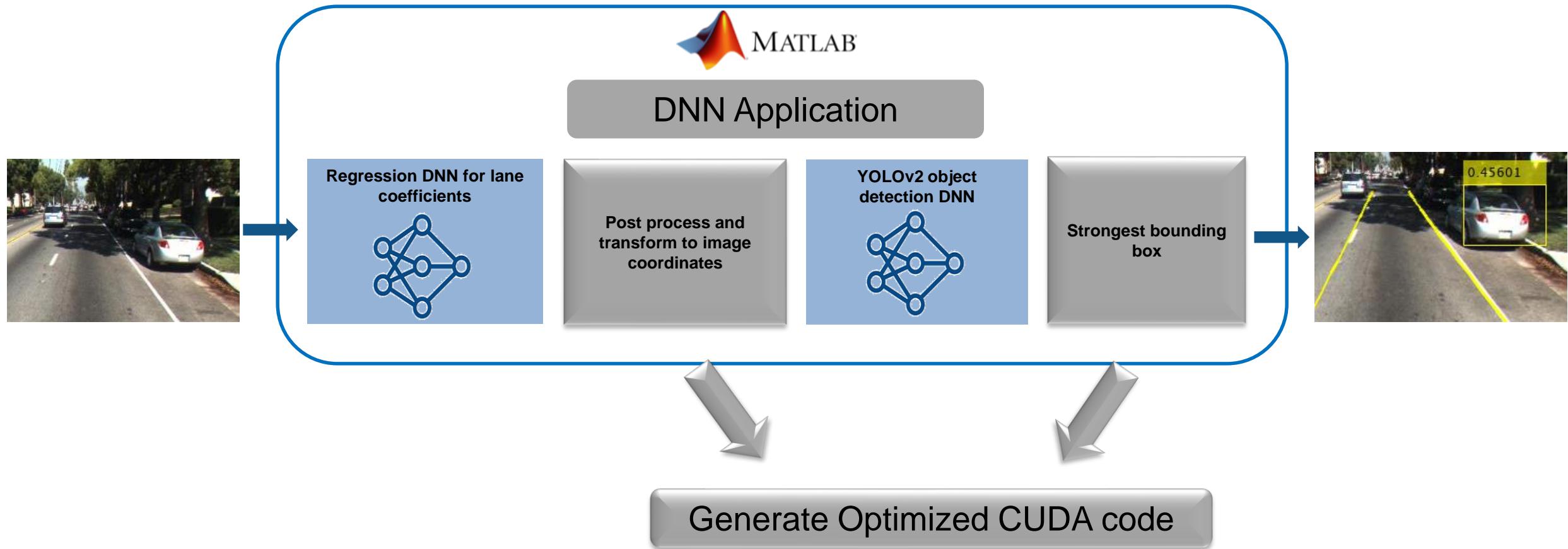
TensorRT



TensorRT is great for inference, ... but, applications require more than inference

Generate code for your application with GPU Coder





GPU Coder automatically extracts parallelism from MATLAB

1. Scalarized MATLAB ("for-all" loops)

```
%% Pixel processing on the height/width of an image
for i = 1:height
    for j = 1:width
        tmpVal = (width*height);
        for x = 1:width
            disValTmp = temp(x,i);
            dist = ((j-x)^2 + disValTmp^2 );
            if (dist < tmpVal)
                tmpVal = single(dist);
            end
        end
        out(i,j) = tmpVal;
    end
end
```

Infer CUDA kernels from MATLAB loops

2. Vectorized MATLAB (math operators and library functions)

```
%% Parallel element-wise math to compute
% Restoration with inverse Koschmieder's law
factor=1.0./(1.0-(diff_im));
restoreOut(:,:,1)= (input(:,:,1)-diff_im).*factor;
restoreOut(:,:,2)= (input(:,:,2)-diff_im).*factor;
restoreOut(:,:,3)= (input(:,:,3)-diff_im).*factor;
```

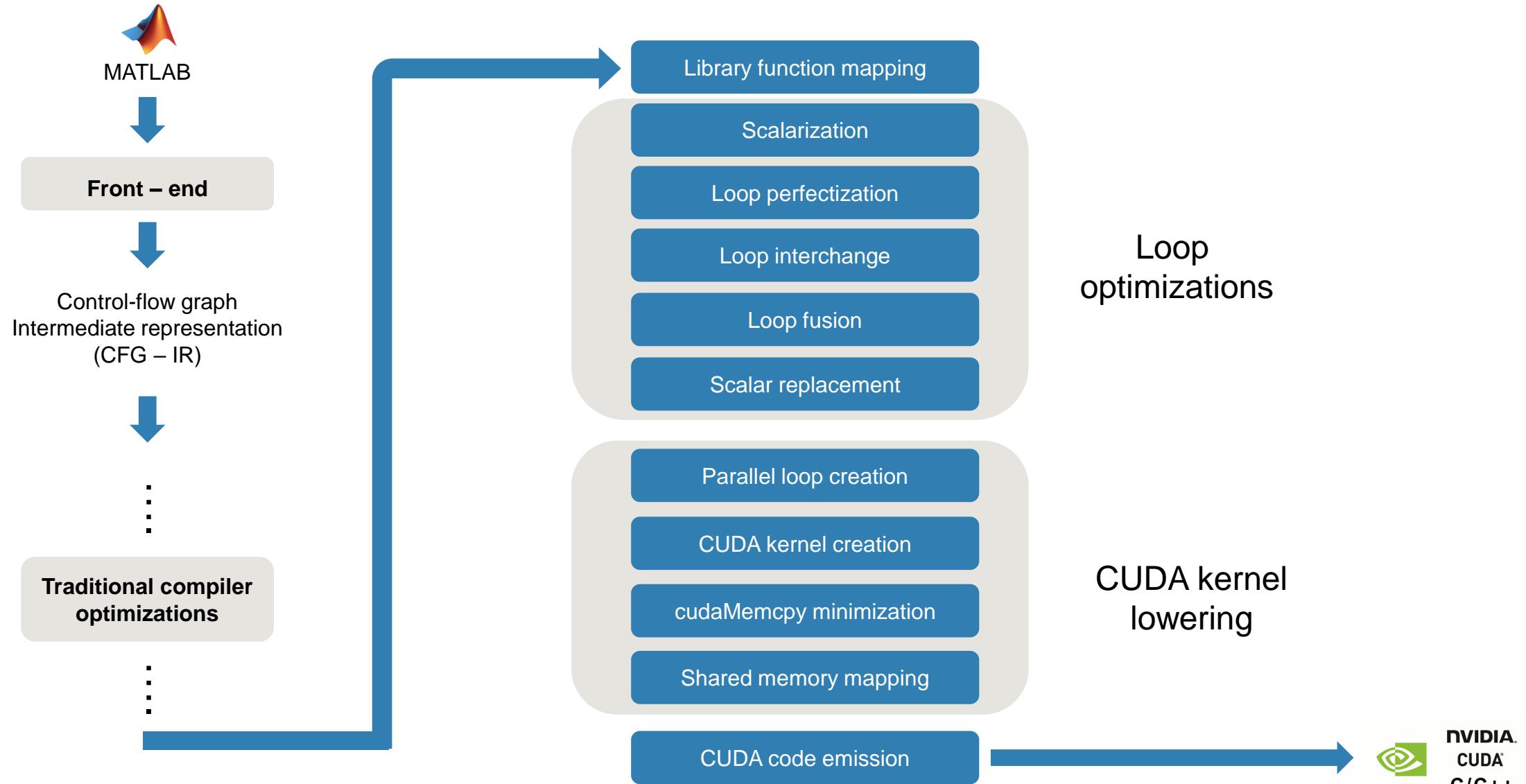
3. Composite functions in MATLAB (maps to cuBlas, cuFFT, cuSolver, cuDNN, TensorRT)

```
C = A * B; % cuBLAS
y = fft(in); % cuFFT
z = A \ B; % cuSolver
```

Library replacement

```
predictions = detectionnet.activations(img_rz,56,'OutputAs','channels'); % cuDNN or TensorRT
```

GPU Coder runs a host of compiler transforms to generate CUDA



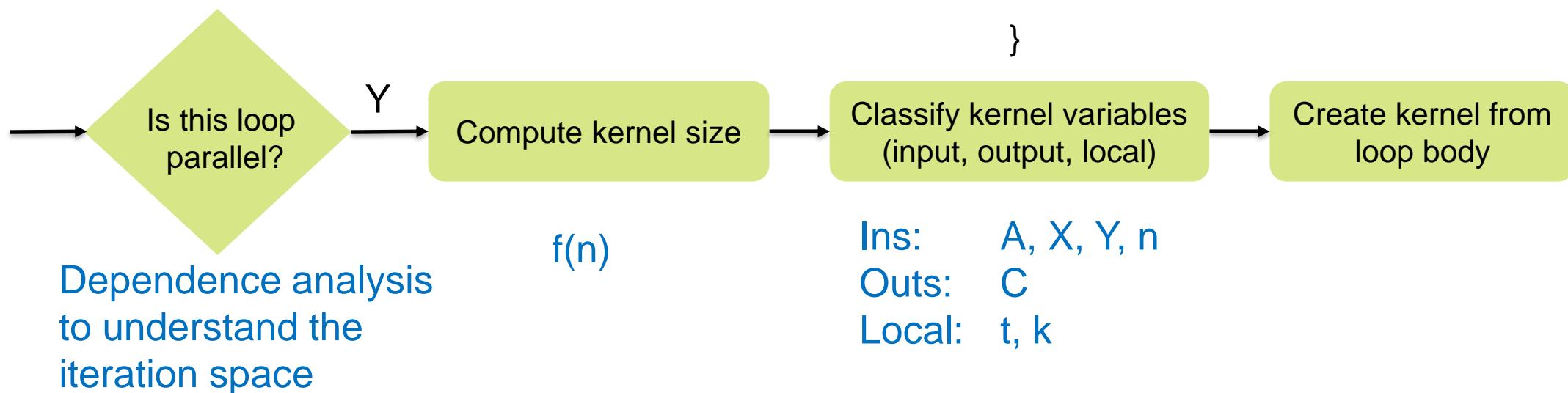
From a loop to a CUDA kernel

Extracting parallelism from loops

1. **Scalarized MATLAB (for loops)**
2. Vectorized MATLAB

```
for k = 1:n
    t = A(k) .* X(k);
    C(k) = t + Y(k);
end
```

```
{ ...
    mykernel<<< f(n) >>>(
    ...
)
static __global__ mykernel(A, X, Y, C, n)
{
```



From a loop to a CUDA kernel

Extracting parallelism from loops

1. **Scalarized MATLAB (for loops)**
2. Vectorized MATLAB

Imperfect Loops

```
for i = 1:p
    ...(outer prologue code)...
    for j = 1:m
        for k = 1:n
            ...(inner loop)...
        end
        ...(outer epilogue code)...
    end
end
```

Perfect Loops

```
for i = 1:p
    for j = 1:m
        for k = 1:n
            ...(inner loop)...
        end
    end
end
```

From a loop to a CUDA kernel

- Extracting parallelism from loops
1. **Scalarized MATLAB (for loops)**
 2. Vectorized MATLAB

Imperfect Loops

```
for i = 1:p
    ... (outer prologue code) ...
    for j = 1:m
        for k = 1:n
            ... (inner loop) ...
        end
        ... (outer epilogue code) ...
    end
end
```

Perfect Loops

```
for i = 1:p
    for j = 1:m
        for k = 1:n
            ... (outer prologue code) ...
            ... (inner loop) ...
            if k == n
                ... (outer epilogue code) ...
            end
        end
    end
end
```



From vectorized MATLAB to CUDA kernels

Extracting parallelism from loops

1. Scalarized MATLAB (for loops)
2. **Vectorized MATLAB**

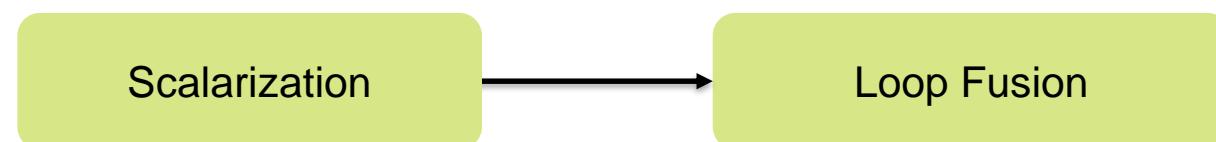
```
output(:, 1) = (input(:, 1) - x_im) .* factor;
```

```
for i = 1:M
    diff(i) = input(i, 1) - x_im(i);
end
for a = 1:M
    output(i, 1) = diff(i) * factor(i);
end
```

```
for i = 1:M
    diff(i) = input(i, 1) - x_im(i);
    output(i, 1) = diff(i) * factor(i);
end
```

Assume the following sizes:

'output'	: M x 3
'input'	: M x 3
'x_im'	: M x 1
'factor'	: M x 1



Reduce to
for-loops

Create larger parallel
loops (and hence
CUDA kernels)

From vectorized MATLAB to CUDA kernels

Extracting parallelism from loops

1. Scalarized MATLAB (for loops)
2. **Vectorized MATLAB**

```
output(:, 1) = (input(:, 1) - x_im) .* factor;
```

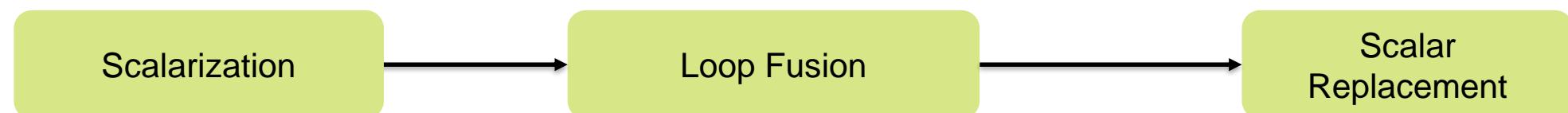
```
for i = 1:M
    diff(i) = input(i, 1) - x_im(i);
end
for a = 1:M
    output(i, 1) = diff(i) * factor(i);
end
```

```
for i = 1:M
    diff(i) = input(i, 1) - x_im(i);
    output(i, 1) = diff(i) * factor(i);
end
```

Assume the following sizes:

'output'	: M x 3
'input'	: M x 3
'x_im'	: M x 1
'factor'	: M x 1

```
for i = 1:M
    tmp = input(i, 1) - x_im(i);
    output(i, 1) = tmp * factor(i);
end
```



Reduce to
for-loops

Create larger parallel
loops (and hence
CUDA kernels)

Reduce temp matrices
to temp scalars

Optimizing CPU-GPU data movement is a challenge

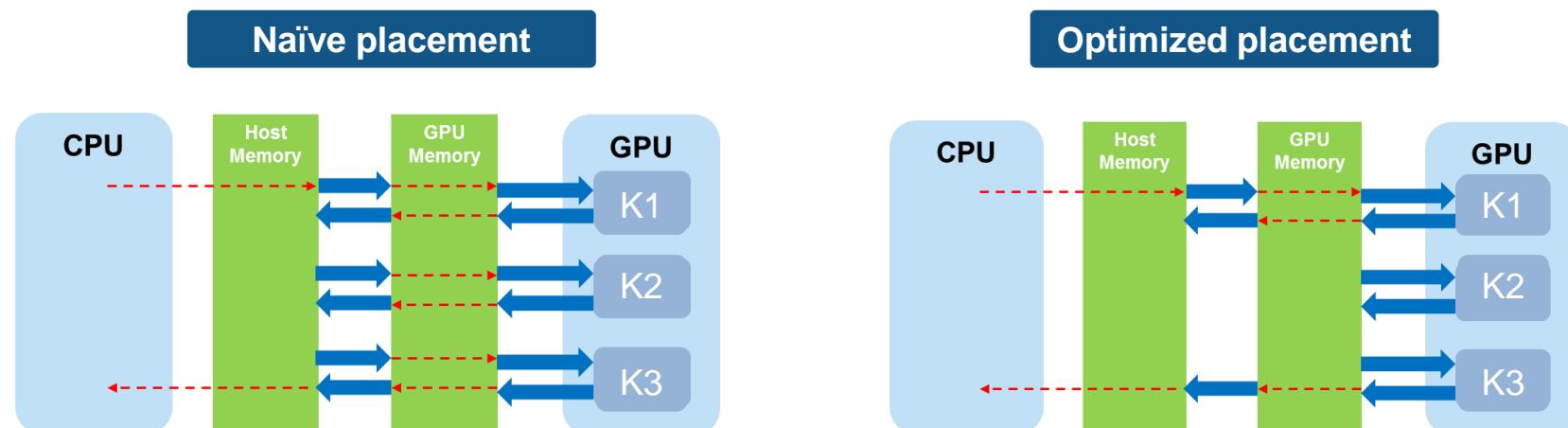
[Memory Optimizations](#)

```
A = ...  
...  
for i = 1:N  
    ... A(i)  
end  
...  
  
imfilter  
...
```



```
A = ...  
...  
cudaMemcpyHtoD(gA, a);  
kernel1<<<...>>>(gA)  
cudaMemcpyDtoH(...)  
...  
cudaMemcpyHtoD(...)  
imfilter_kernel(...)  
cudaMemcpyDtoH(...)  
...
```

Where is the ideal placement of memcpy?



GPU Coder optimizes memcpy placement

```

A(:) = ...
C(:) = ...

for i = 1:N
    ...
    gB = kernel1(gA);
    gA = kernel2(gB);
    if (some_condition)
        gC = kernel3(gA, gB);
    end
    ...
end
...
... = C;

```

cudaMemcpy *not* needed

cudaMemcpy *definitely* needed

cudaMemcpy *may be* needed

Assume gA, gB and gC are mapped to GPU memory

Observations

- Equivalent to Partial redundancy elimination (PRE)
- Dynamic strategy – track memory location with a status flag per variable
- Use-Def to determine where to insert memcpy

Generated (pseudo) code

```

A(:) = ...
A_isDirtyOnCpu = true;
...
for i = 1:N
    if (A_isDirtyOnCpu)
        cudaMemcpy(gA, A);
        A_isDirtyOnCpu = false;
    end
    gB = kernel1(gA);
    gA = kernel2(gB);
    if (somecondition)
        gC = kernel3(gA, gB);
        C_isDirtyOnGpu = true;
    end
    ...
end
...
if (C_isDirtyOnGpu)
    cudaMemcpy(C, gC);
    C_isDirtyOnGpu = false;
end
... = C;

```

From composite functions to optimized CUDA

Core math

- Matrix multiply (cuBLAS)
- Linear algebra (cuSolver)
- FFT functions (cuFFT)
- Convolution
- ...

Image processing Computer vision

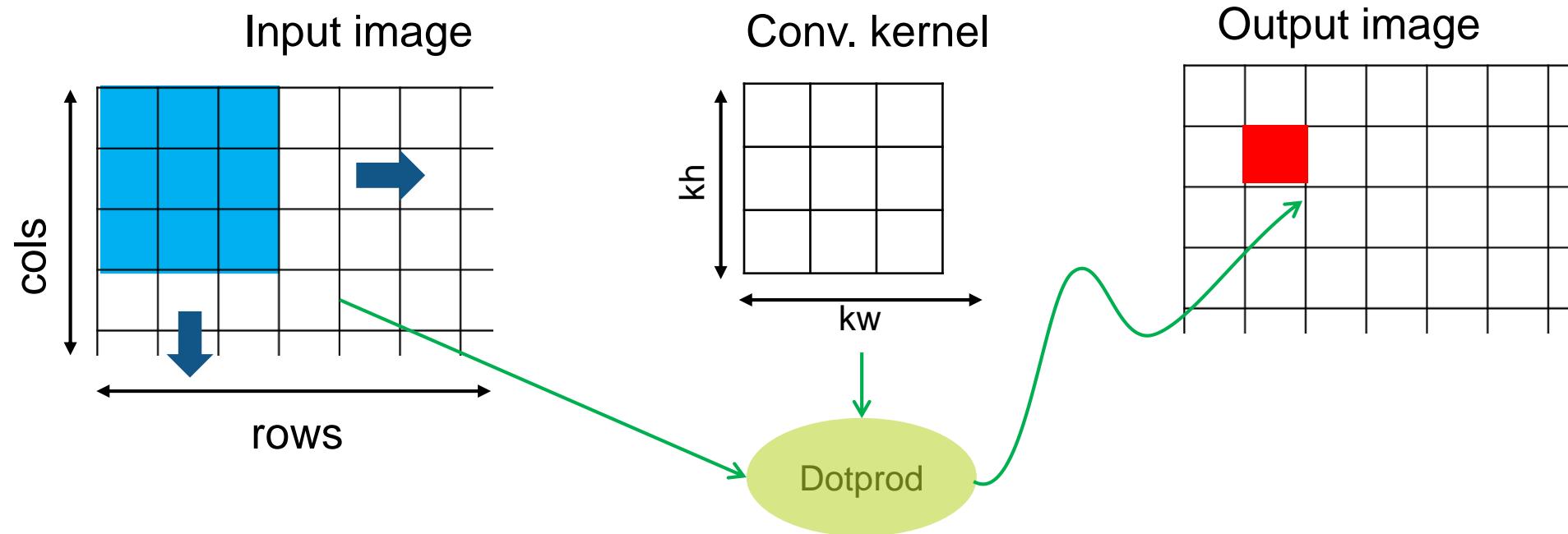
- imfilter
- imresize
- imerode
- imdilate
- bwlabel
- imwarp
- ...

Neural Networks

- Deep learning inference (cuDNN, TensorRT)
- ...

Over 300+ MATLAB functions are optimized for CUDA code generation

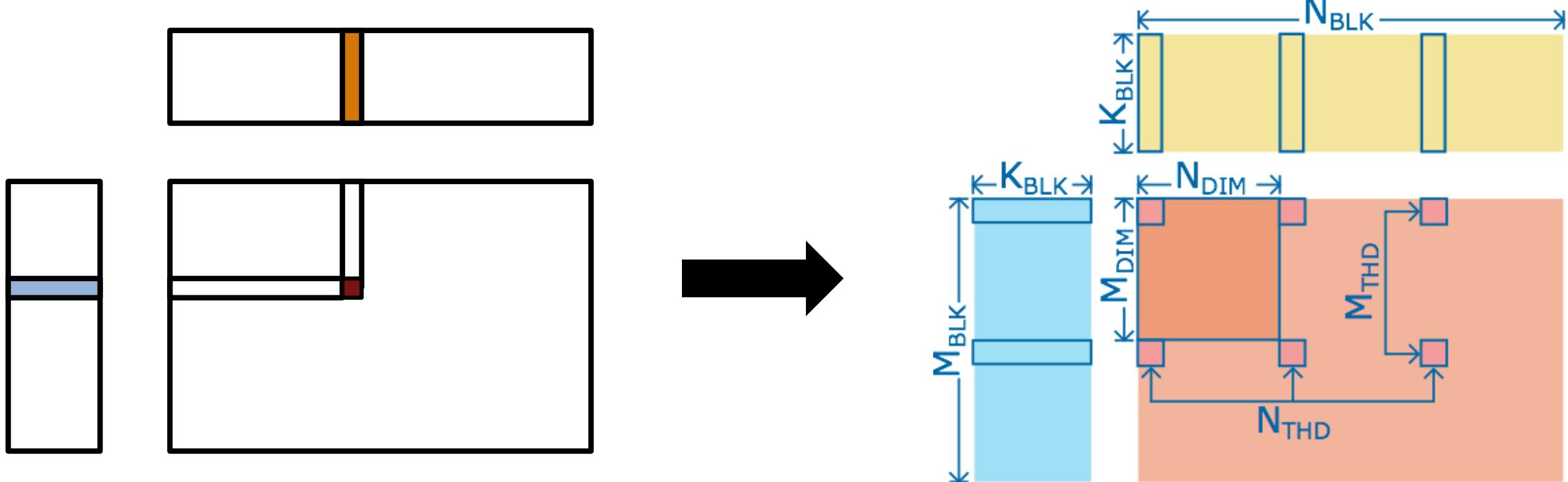
GPU Coder automatically maps data to shared memory



gpuCoder.stencilKernel() automatically generates shared memory algorithm

Used when generating code from any image processing functions in MATLAB
imfilter, imerode, imdilate, conv2, ...

Generates optimal matrix-matrix operations

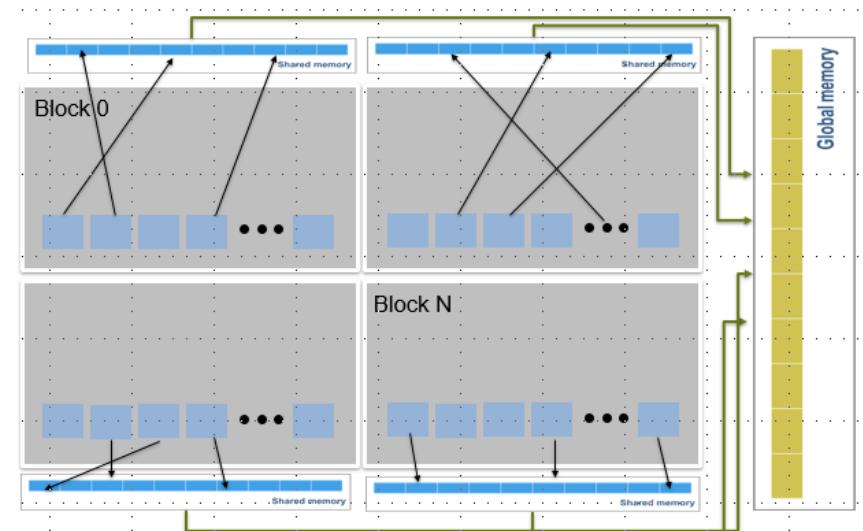
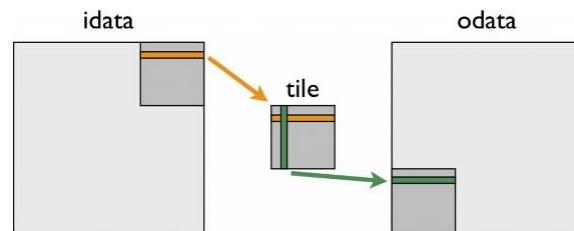
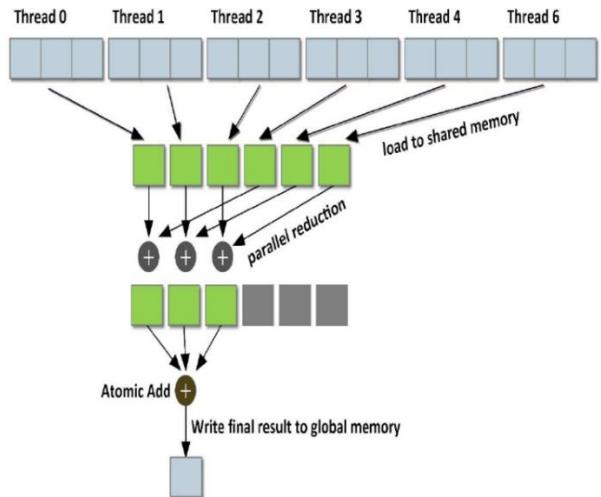


gpuCoder.matrixMatrixKernel automatically generates shared memory algorithm

Used when generating code from many standard MATLAB functions

matchFeatures SAD, SSD, pdist, ...

Generate CUDA friendly algorithms from MATLAB functions



Reductions

Transpose

Histogram

Sort

Min/Max

CumSum

Integrating with legacy CUDA code

Calling a GPU device function

```
__device__ float foo(float a, float b);
```



```
function out = foo(in)
```



Caller needs to be a device
function

```
.....  
out = coder.ceval('-gpudevicefcn,
```

```
    'foo', % your function name
```

```
    a, % arguments
```



Arguments allocated in
device memory

```
    b);
```

```
.....  
end
```

Integrating with legacy CUDA code

Calling a kernel function

```
__global__  
void foo(const float *a[100], const float *b[100], float *c[100]);
```



```
function out = foo(in)  
.....  
out = coder.ceval('foo<1024>',  
                  coder.rref(a, '-gpu'), % value is read on GPU  
                  coder.rref(b, '-gpu'),  
                  coder.wref(c, '-gpu'));% value is written on GPU  
.....  
end
```

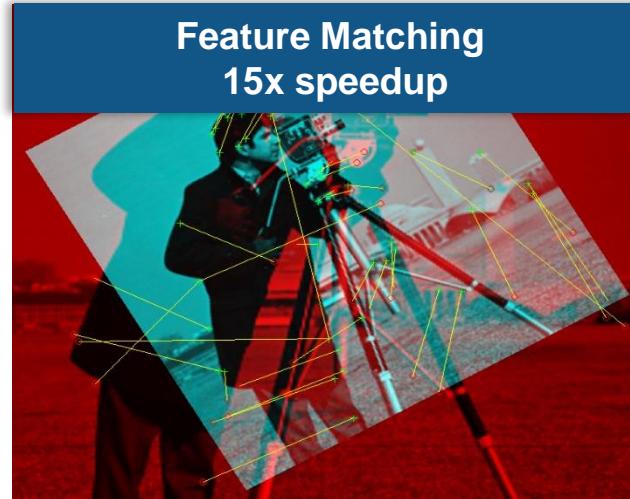


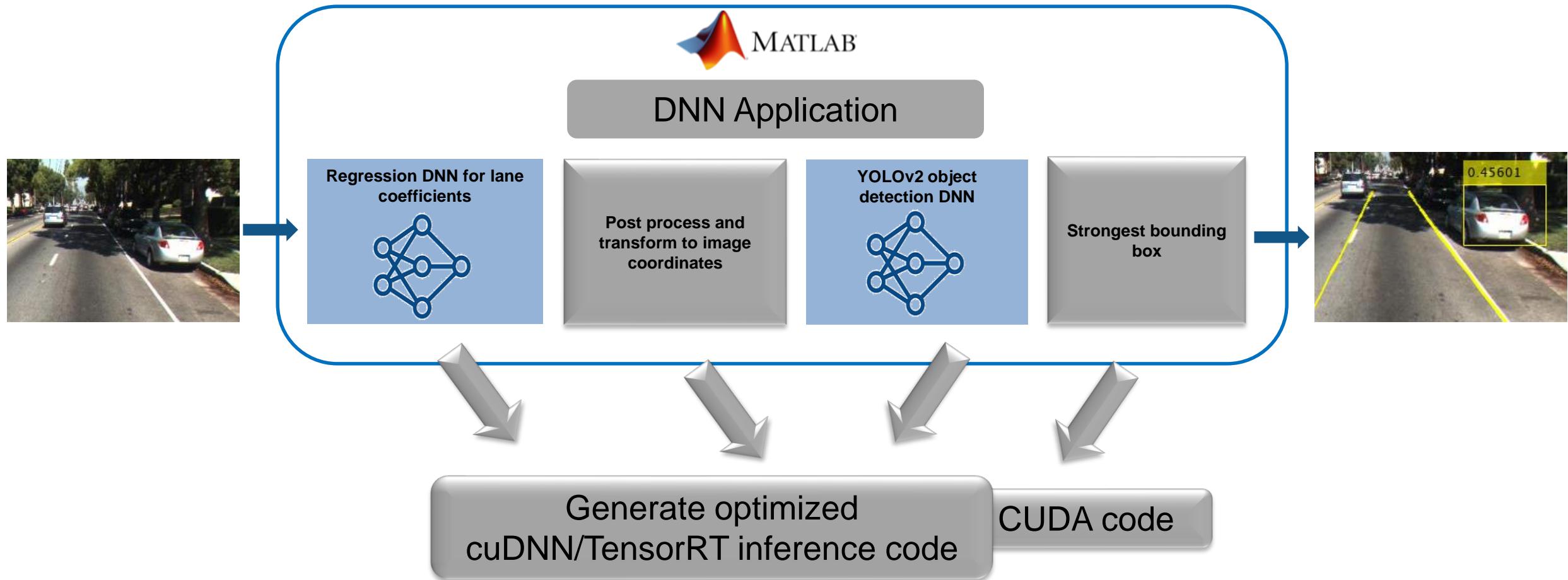
Caller may be a host or device function

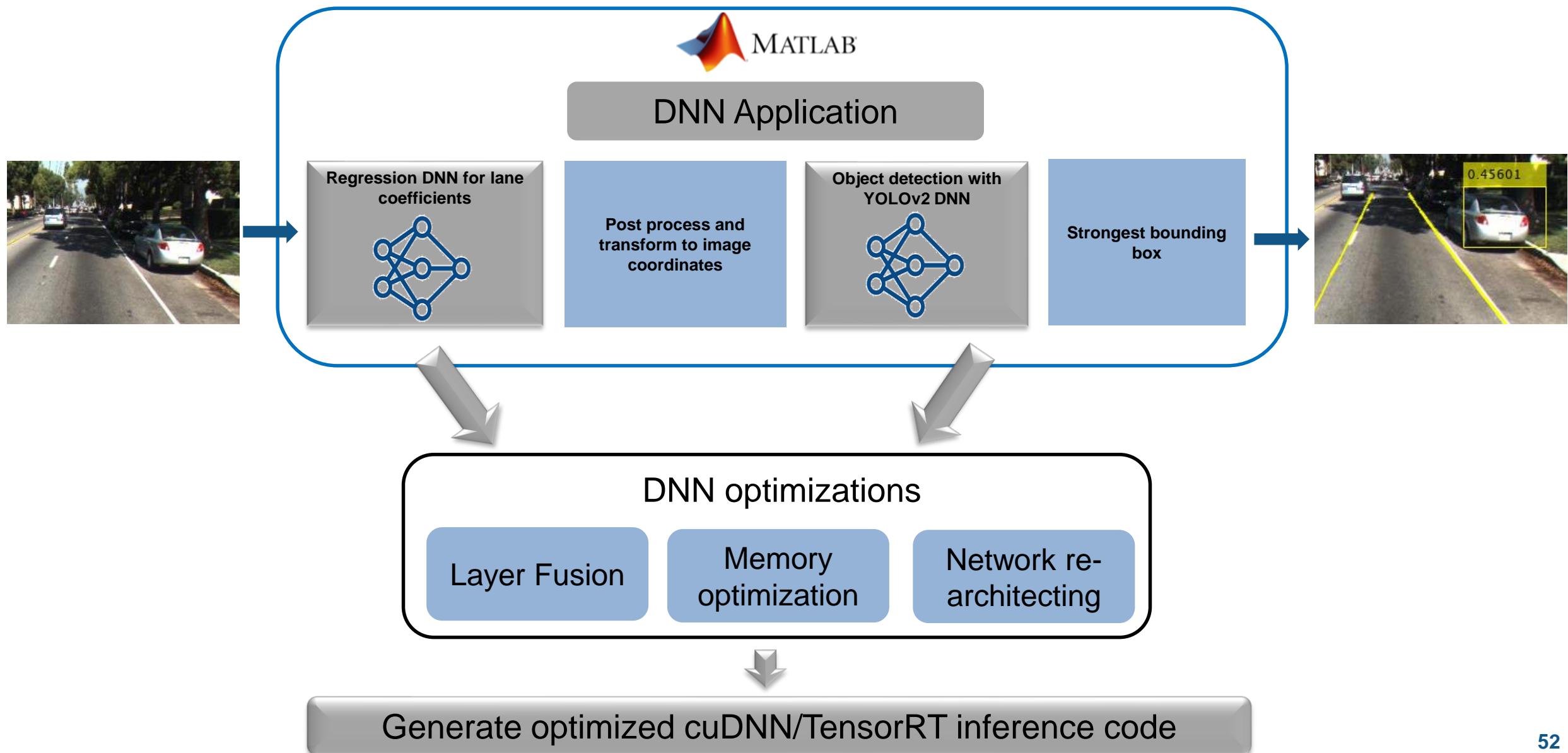


Arguments allocated or copied to device memory as required.

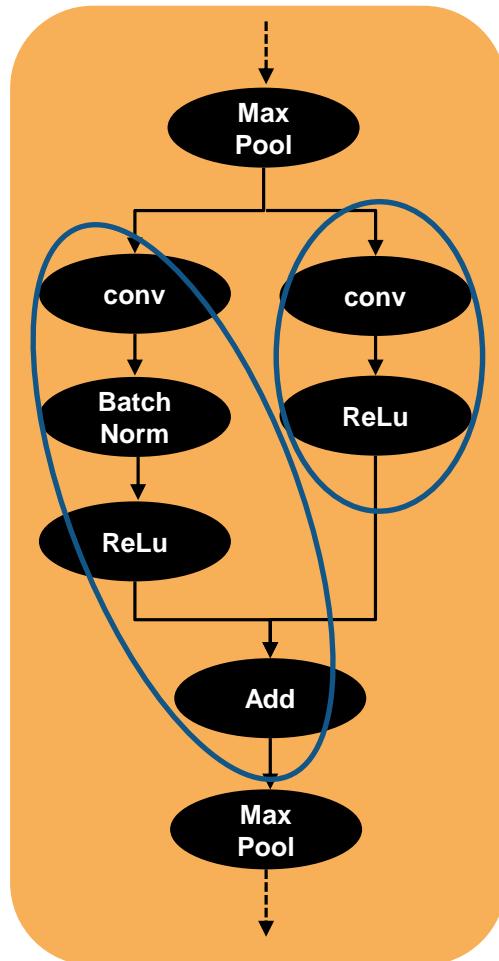
Compiled image processing applications



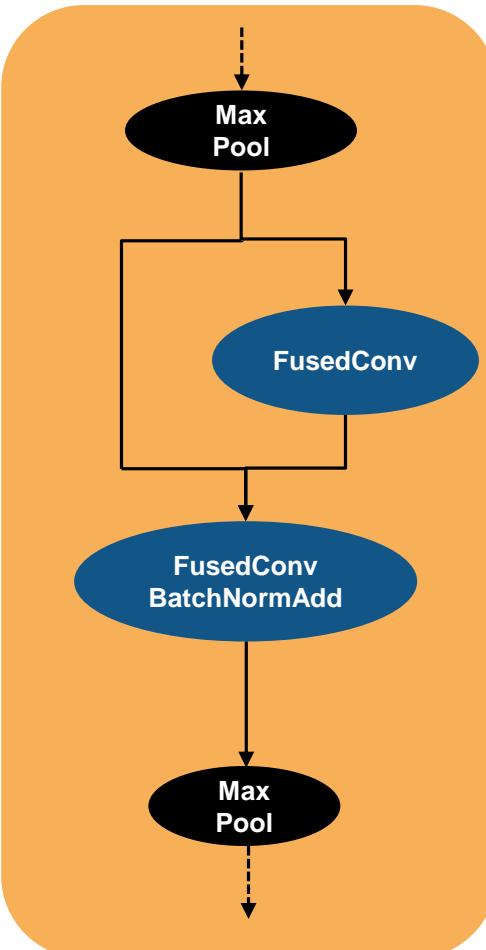




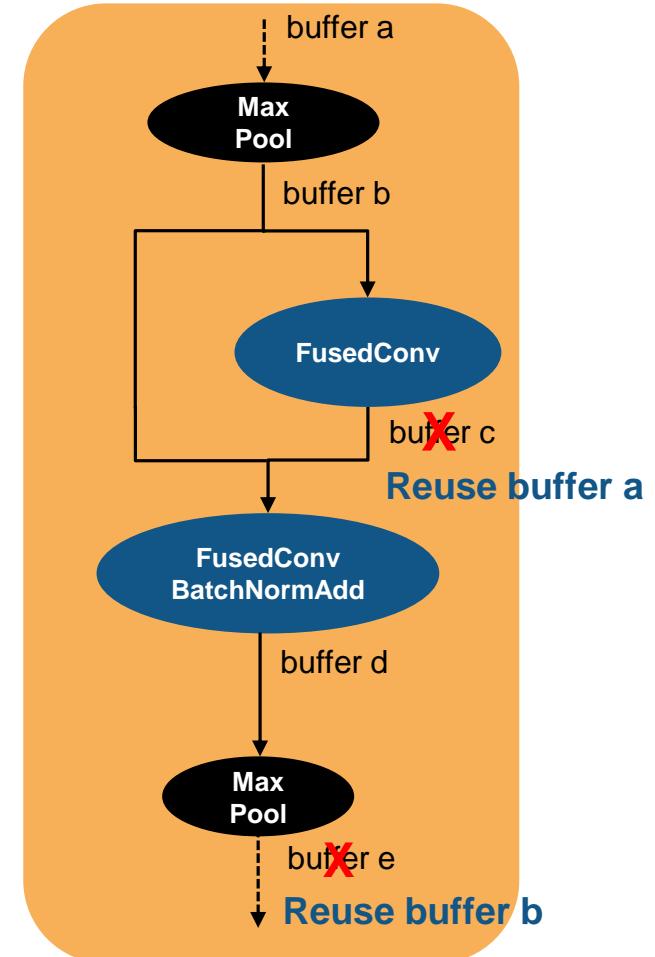
Deep learning network optimizations



Network

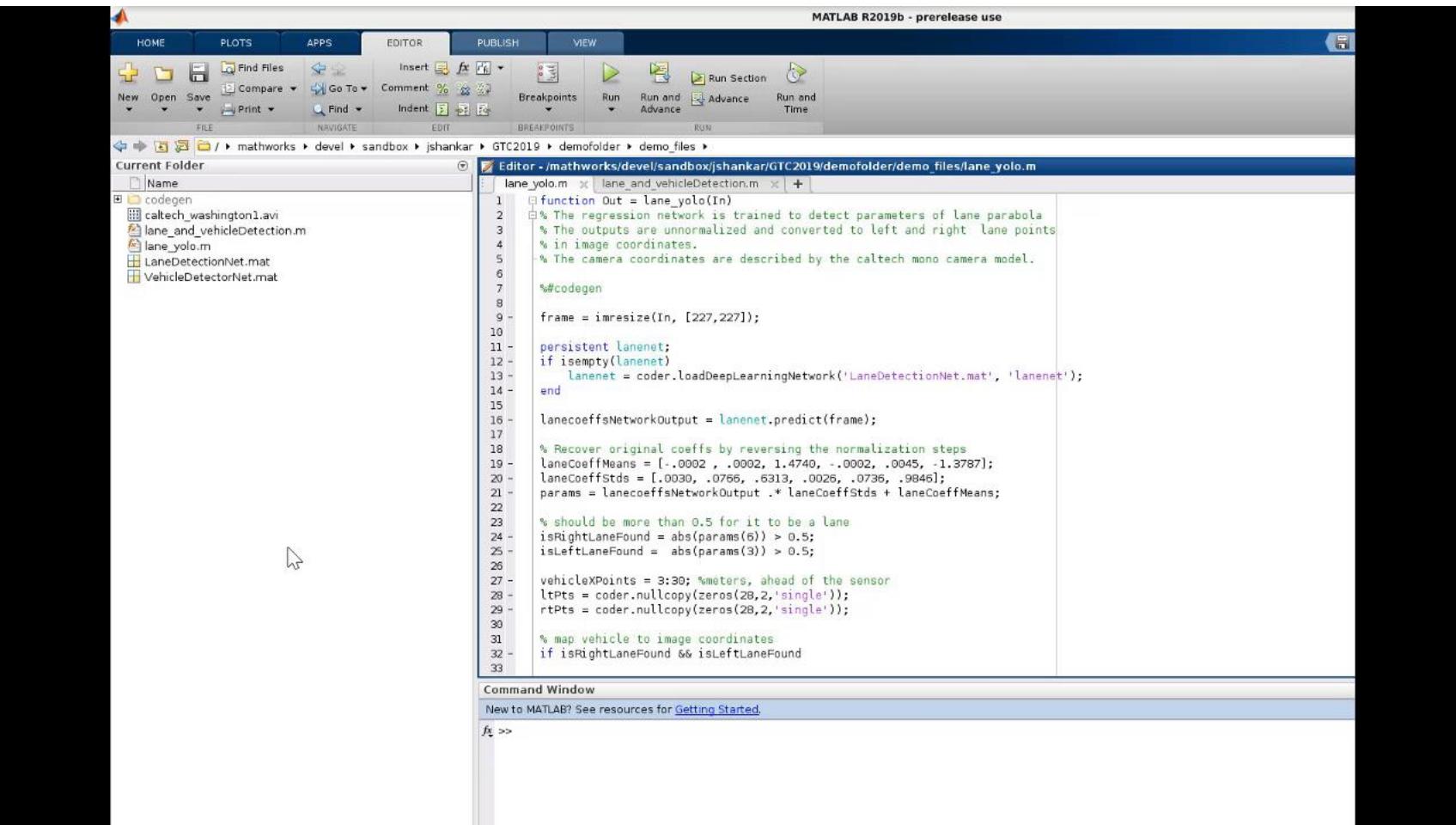


Layer fusion
Optimized computation

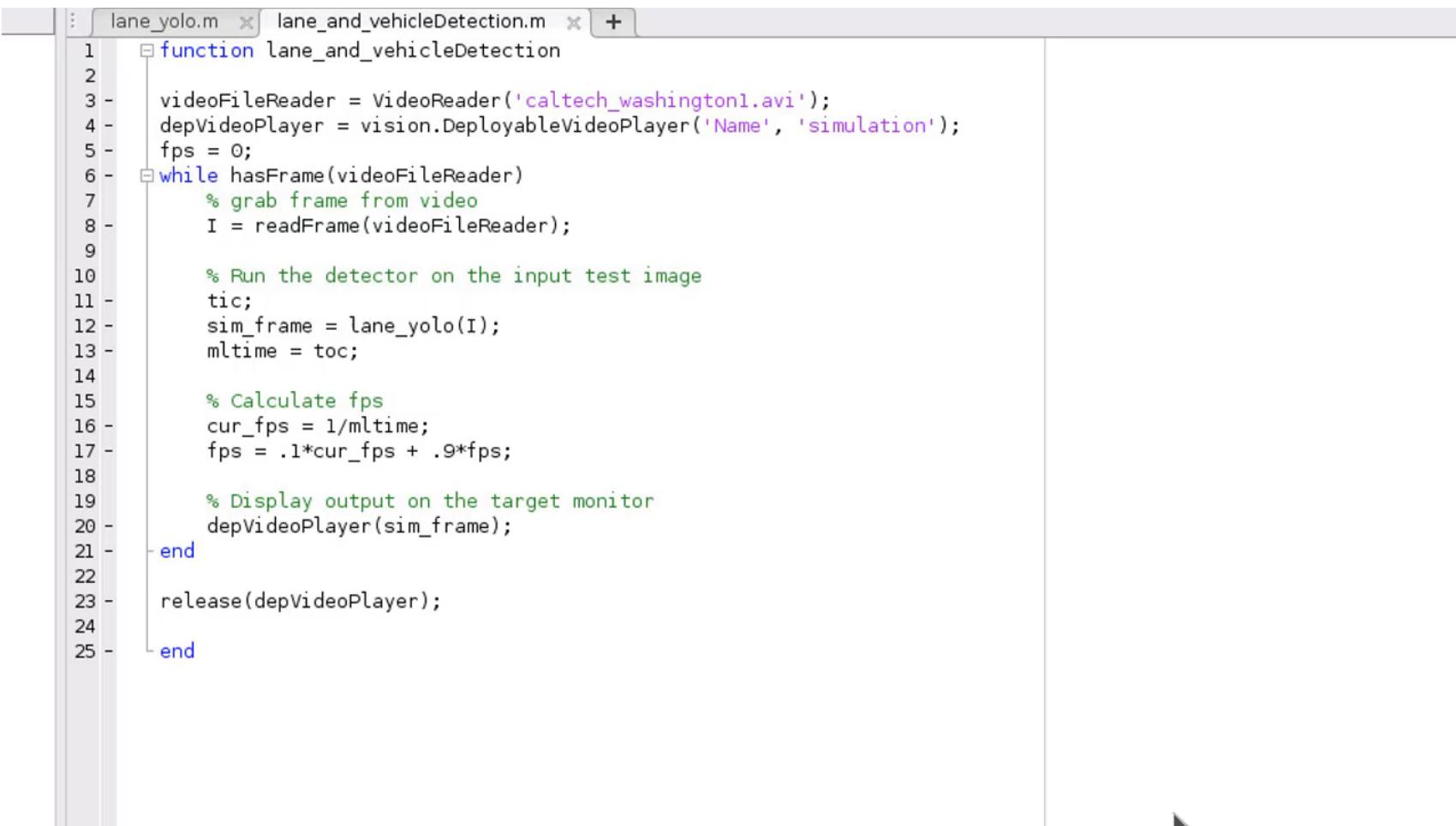


Buffer minimization
Optimized memory

Code generation workflow

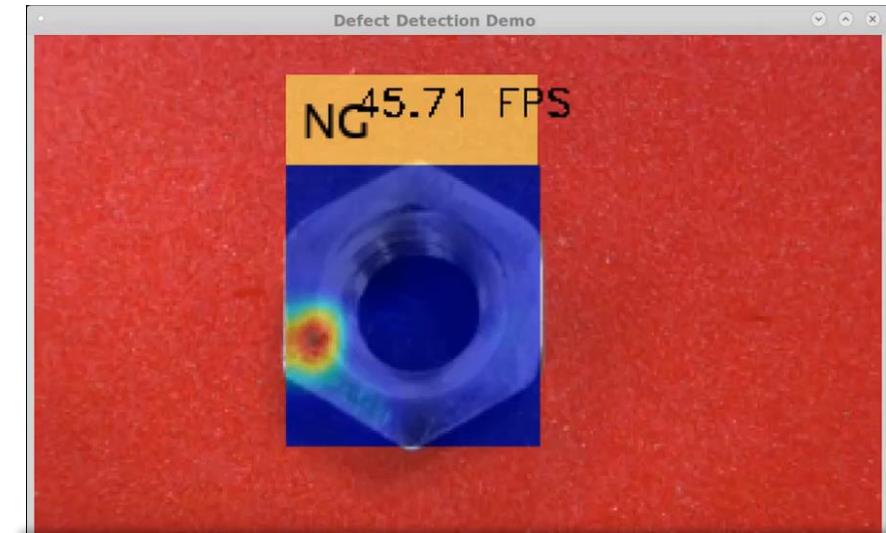
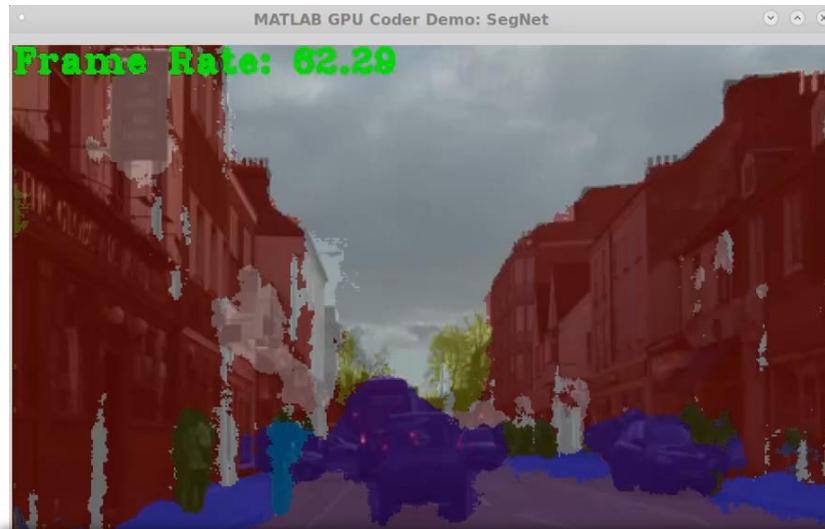


Code generation workflow

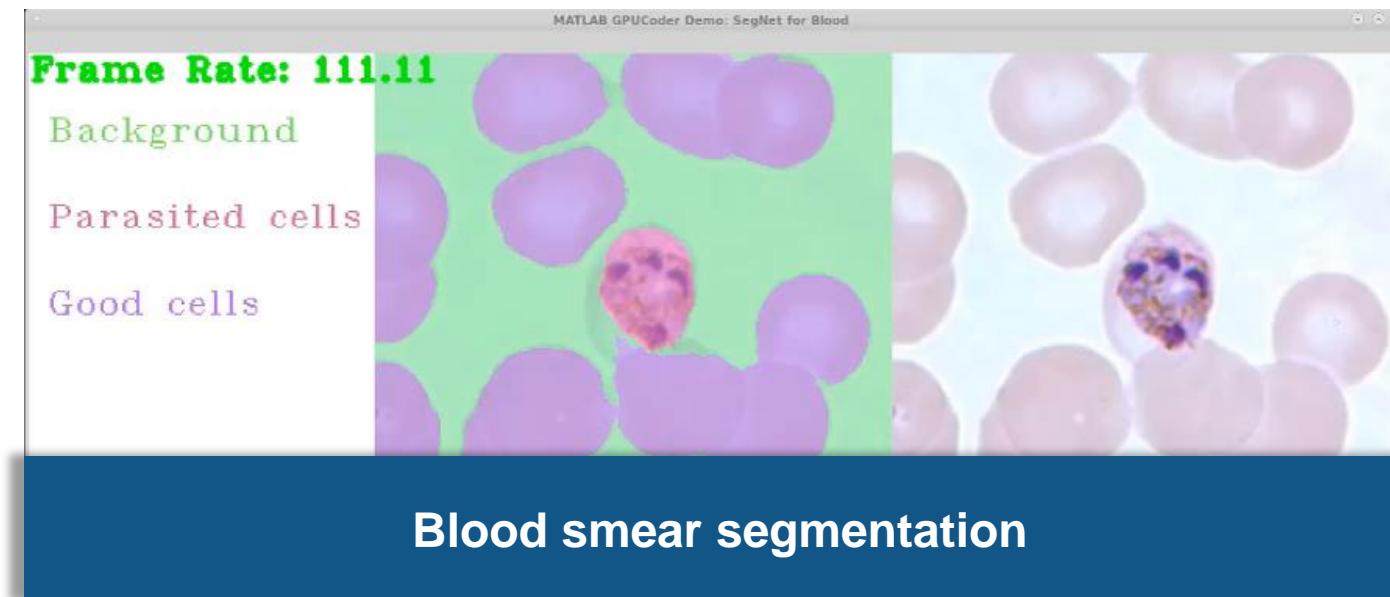


The screenshot shows a MATLAB code editor window with two tabs: 'lane_yolo.m' and 'lane_and_vehicleDetection.m'. The 'lane_and_vehicleDetection.m' tab is active, displaying the following MATLAB script:

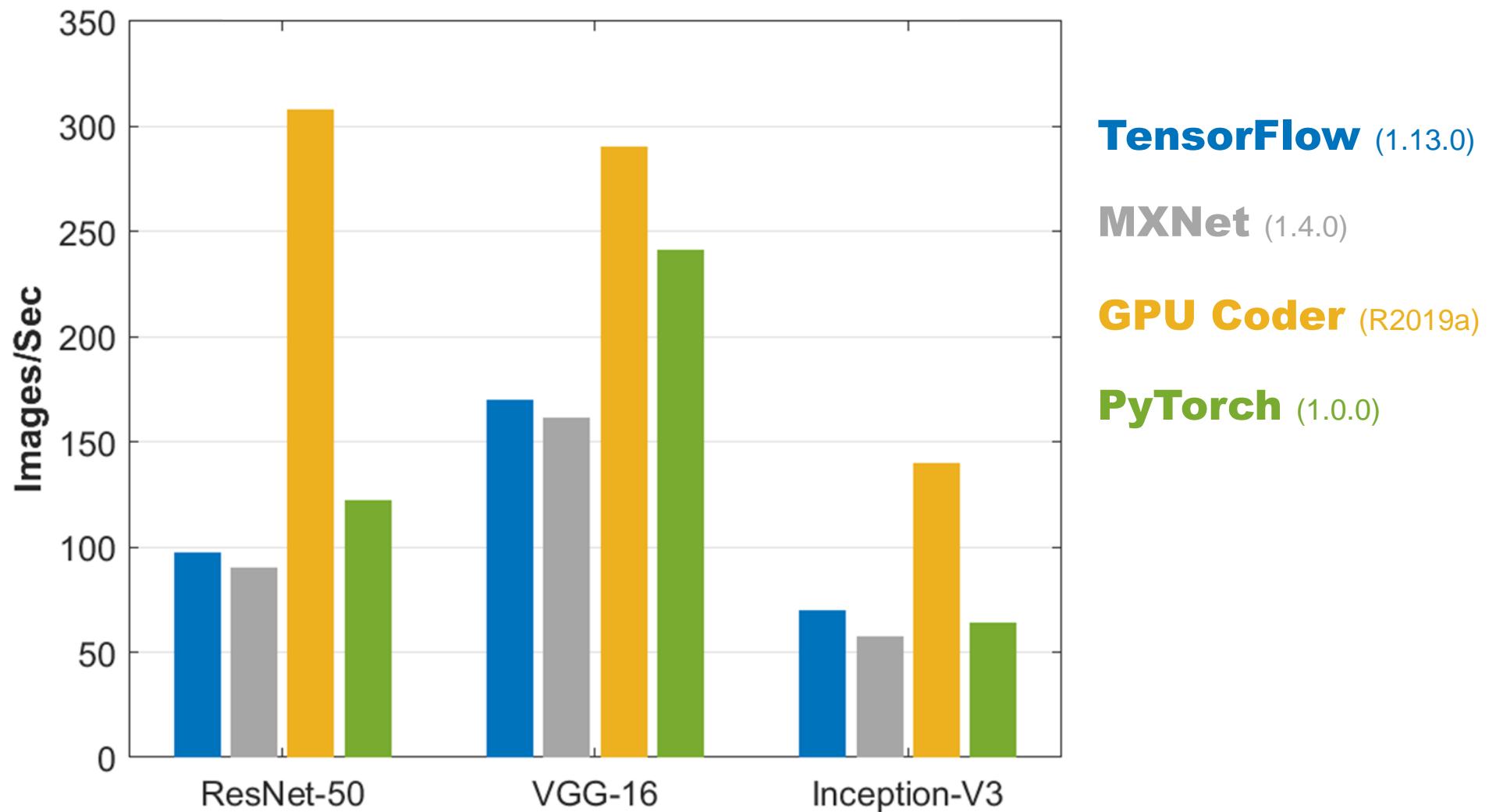
```
1 function lane_and_vehicleDetection
2
3 videoFileReader = VideoReader('caltech_washington1.avi');
4 depVideoPlayer = vision.DeployableVideoPlayer('Name', 'simulation');
5 fps = 0;
6 while hasFrame(videoFileReader)
7     % grab frame from video
8     I = readFrame(videoFileReader);
9
10    % Run the detector on the input test image
11    tic;
12    sim_frame = lane_yolo(I);
13    mltime = toc;
14
15    % Calculate fps
16    cur_fps = 1/mltime;
17    fps = .1*cur_fps + .9*fps;
18
19    % Display output on the target monitor
20    depVideoPlayer(sim_frame);
21 end
22
23 release(depVideoPlayer);
24
25 end
```



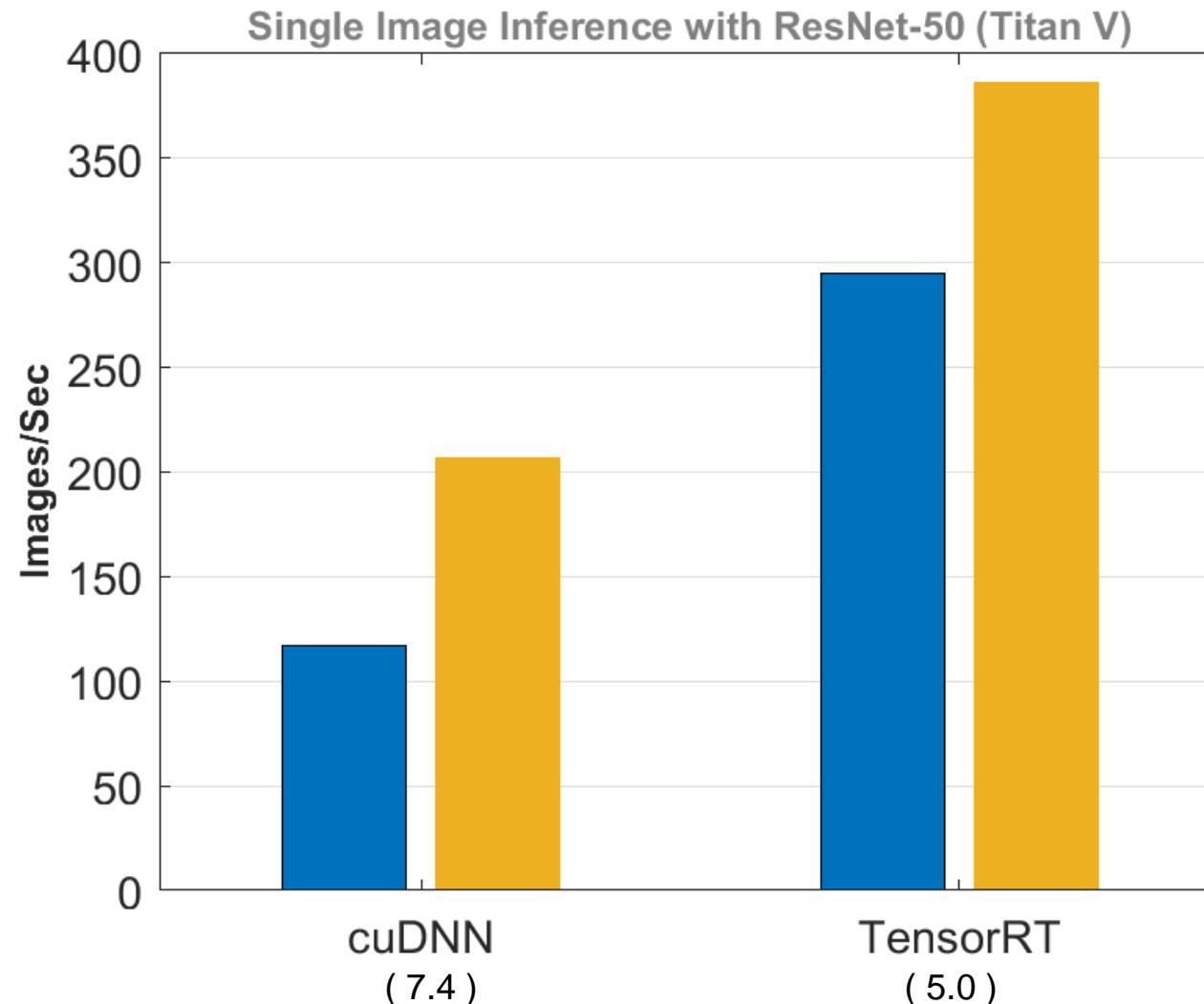
Live Demos
at our GTC
booth #1343



Single Image Inference on Titan V using cuDNN



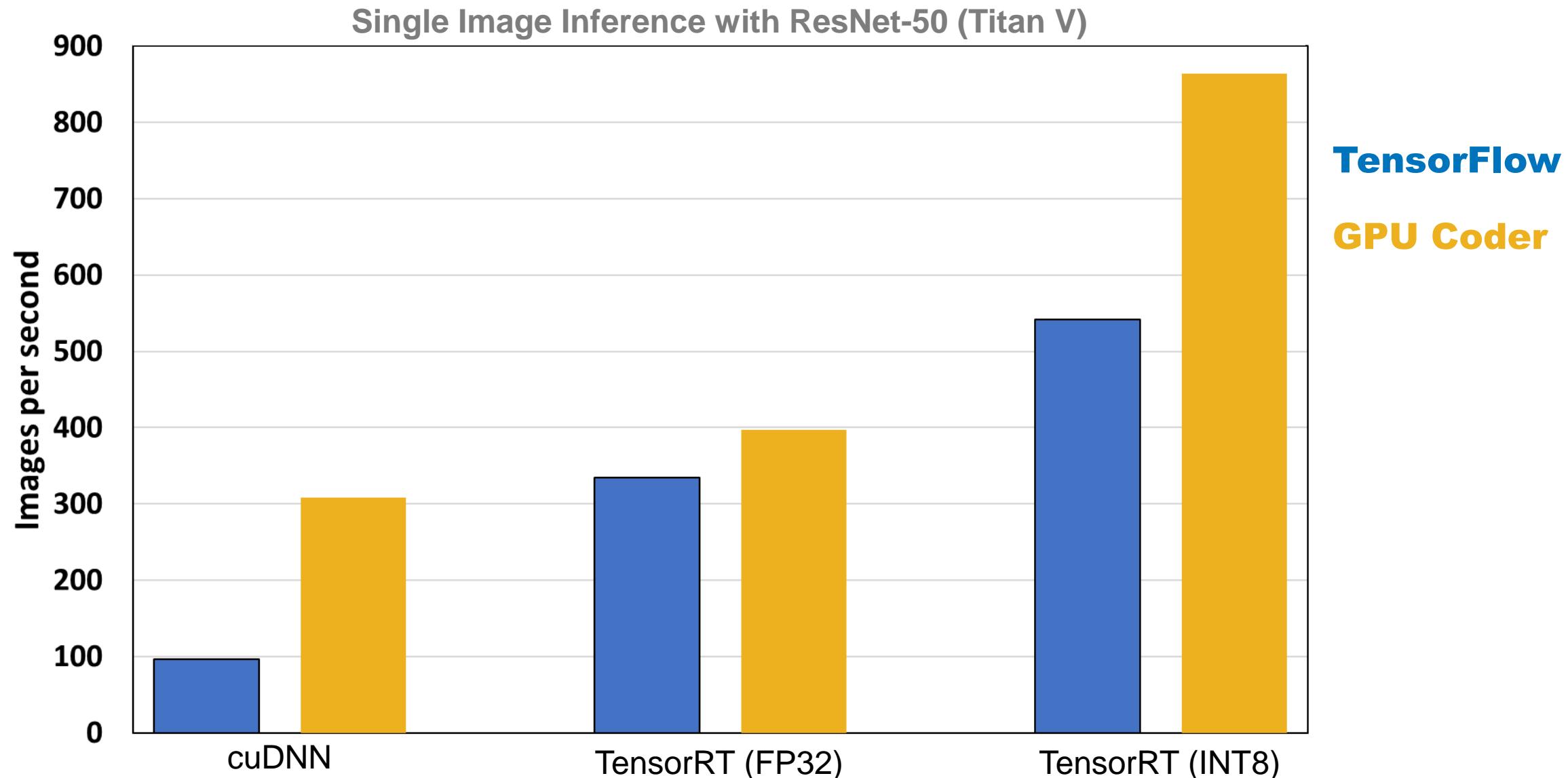
TensorRT Accelerates Inference Performance on Titan V



TensorFlow

MATLAB GPU Coder

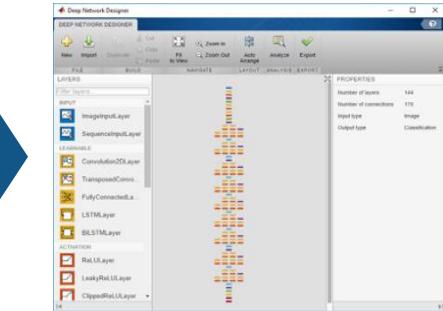
TensorRT Accelerates Inference Performance on Titan V



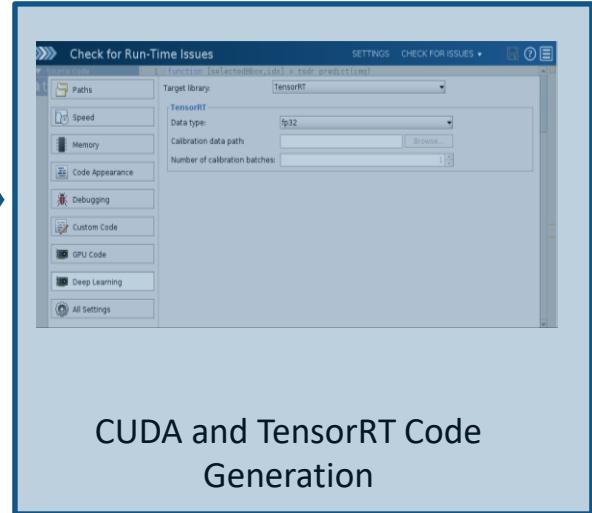
Outline



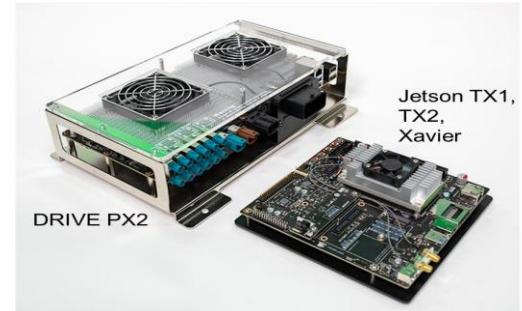
Ground Truth Labeling



Network Design and Training

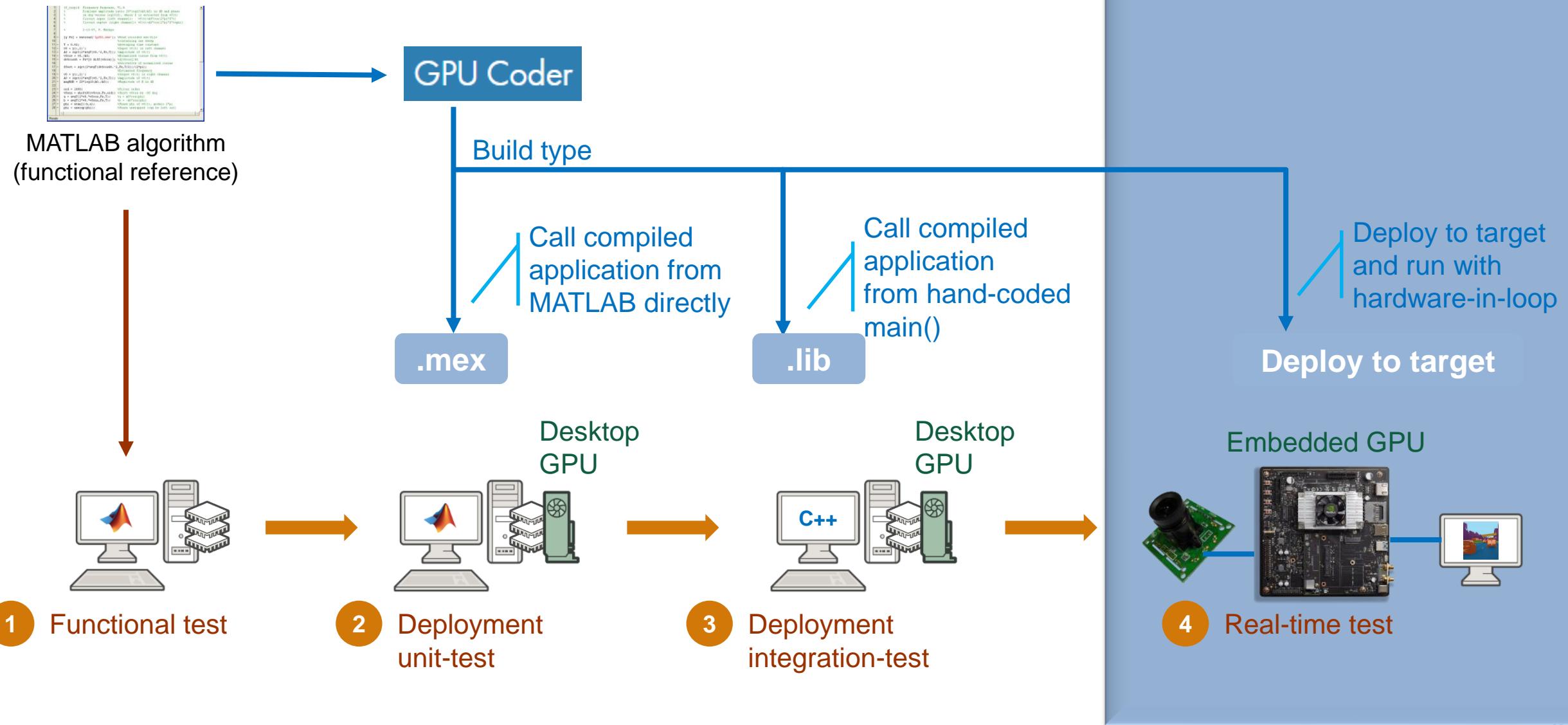


CUDA and TensorRT Code Generation

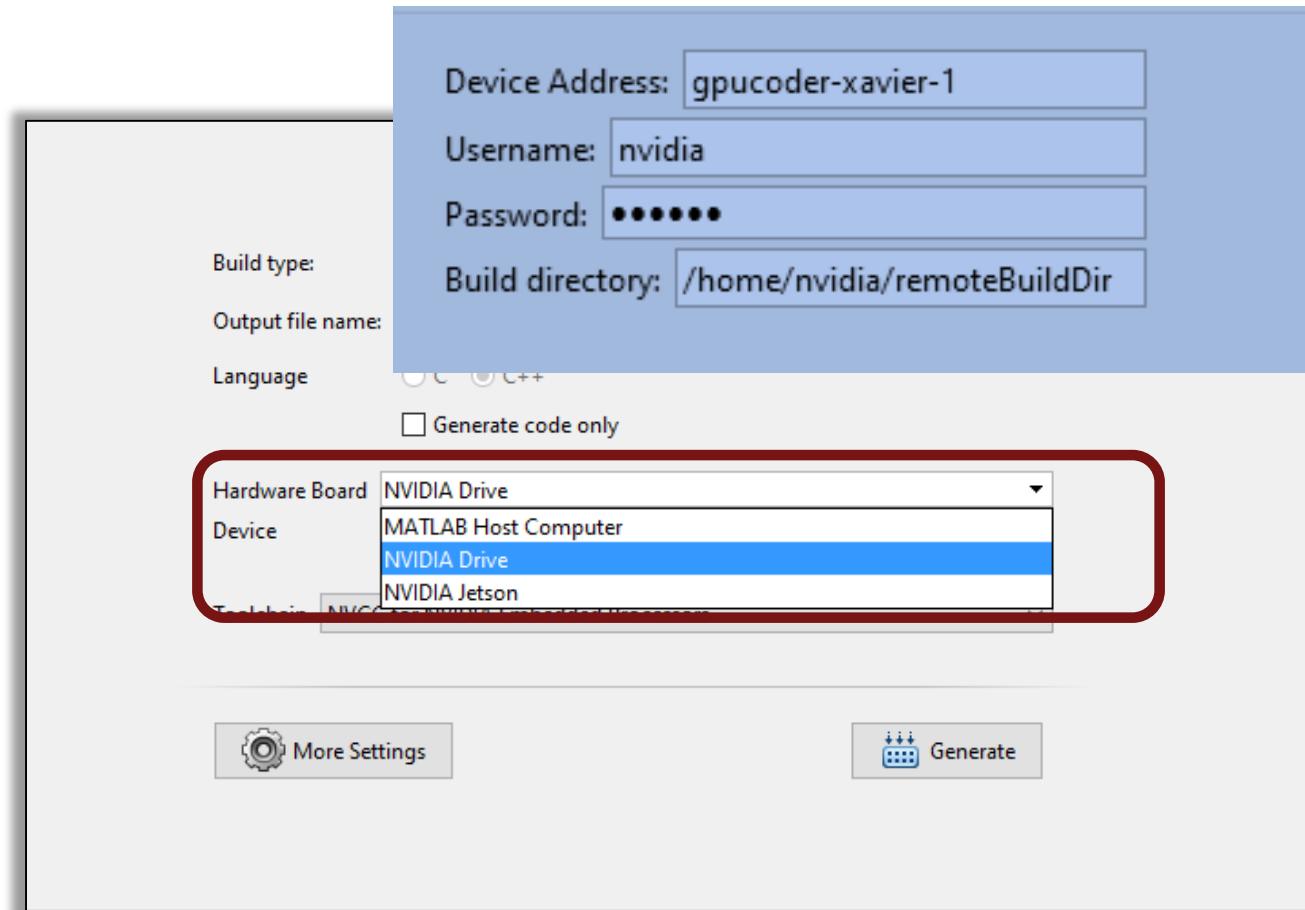
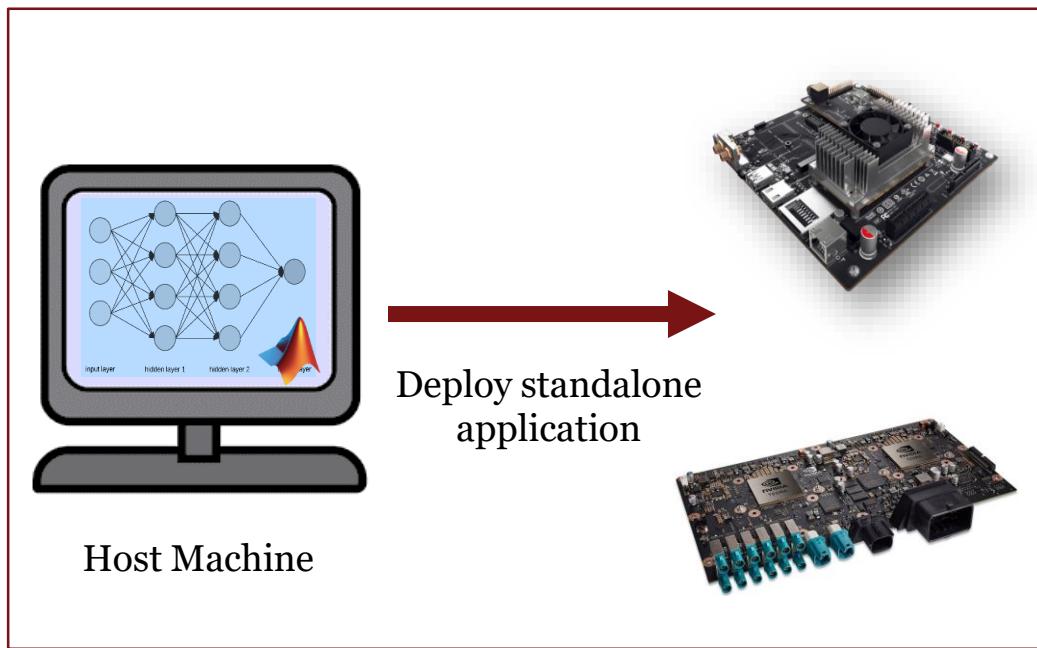


Jetson Xavier and DRIVE
Xavier Targeting

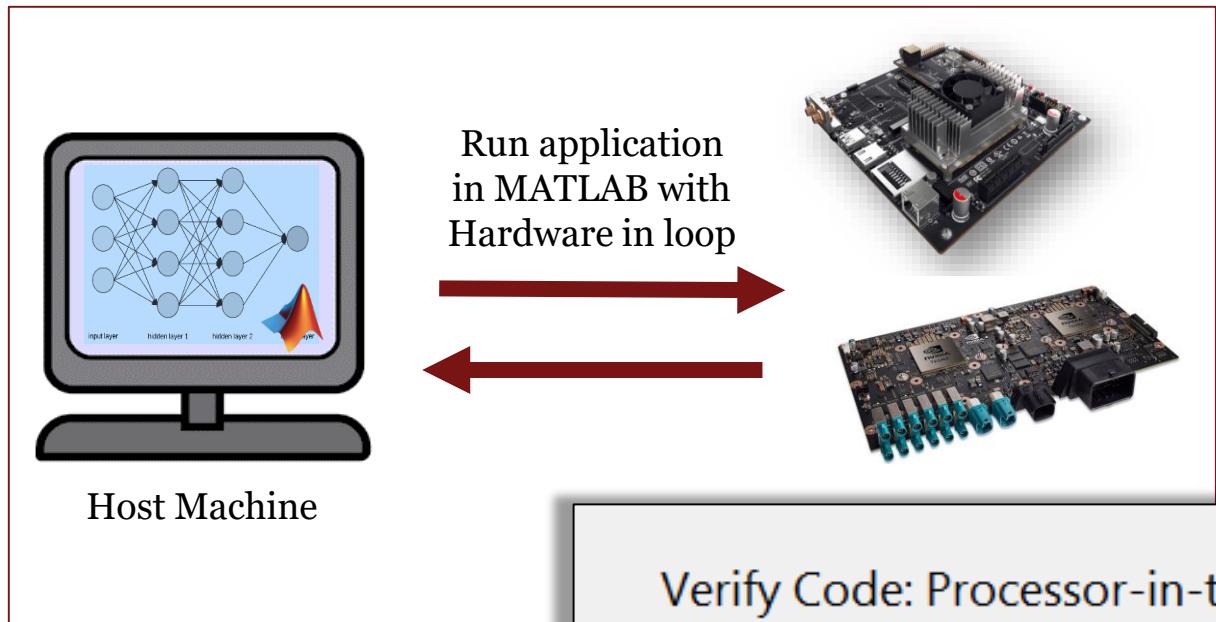
Deploy to Jetson and Drive



Streamlined deployment to Jetson or Drive GPUs.



Closed loop testing and verification on embedded GPUs.



Verify Code: Processor-in-the-Loop Execution (PIL)

Selected output type: Dynamic Library

Hardware: ARM Compatible ARM 64-bit (LP64) (NVIDIA Jetson)

Interface: lane_and_vehicleDetection_pil

```
>> lane_and_vehicleDetection_pil(im)
```

Enable entry point execution profiling

Run using: MATLAB code Generated code

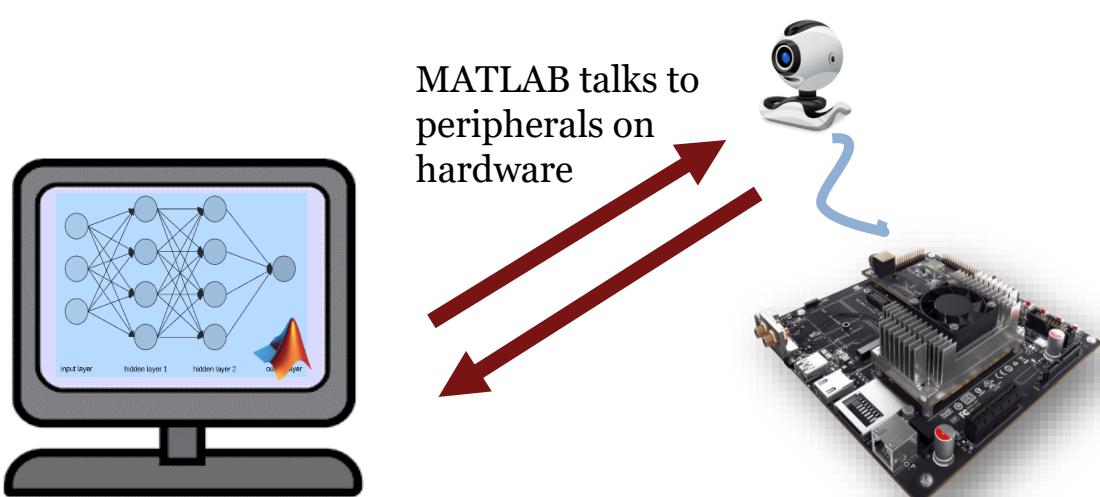
[Run Generated Code](#) [Stop](#) [Help](#)

[Run manually in MATLAB](#)



```
Editor - /mathworks/devel/sandbox/jshankar/GTC2019/demofolder/demo_files/lane_and_vehicleDetection.m
laneyolo.m x lane_and_vehicleDetection.m x +
1 function lane_and_vehicleDetection
2
3 videoFileReader = VideoReader('caltech_washington1.avi');
4 depVideoPlayer = vision.DeployableVideoPlayer('Name', 'simulation');
5 fps = 0;
6 while hasFrame(videoFileReader)
7     % grab frame from video
8     I = readFrame(videoFileReader);
9
10    % Run the detector on the input test image
11    tic;
12    sim_frame = lane_yolo_mex(I);
13    mltime = toc;
14
15    % Calculate fps
16    cur_fps = 1/mltime;
17    fps = .1*cur_fps + .9*fps;
18
19    % Display output on the target monitor
20    depVideoPlayer(sim_frame);
21    pause(.05);
22 end
23
24 release(depVideoPlayer);
25
26 end
```

Access to target peripherals from MATLAB



Verification with processor-in-loop mode

```
function sobelEdgeDetection()
%#codegen

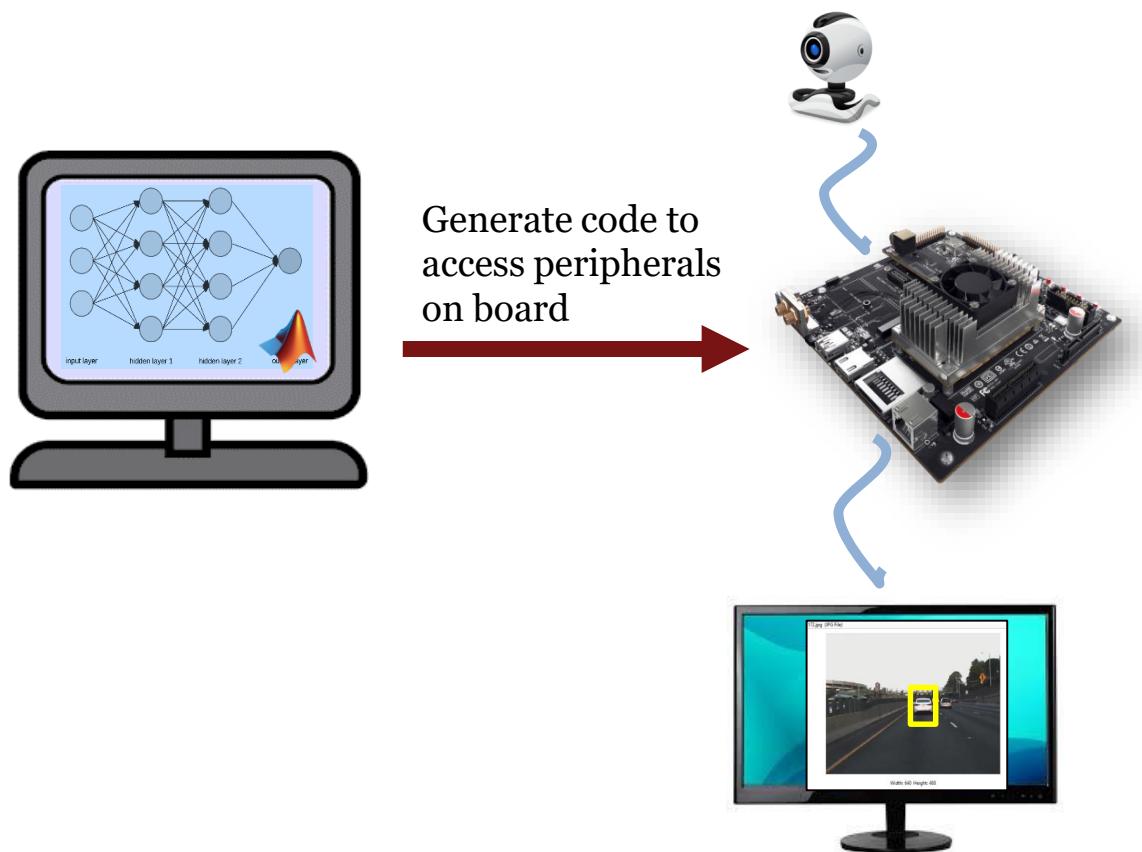
% access webcam on Jetson
hwobj = jetson;
webcamIndex = 1;
w = webcam(hwobj,webcamIndex);
d = imageDisplay(hwobj);

for k = 1:1800
    % Capture the image
    img = snapshot(w);
    % run deployed application on Jetson
    edgeImg = sobel_edge_pil(img, kernel);

    % Display edge image.
    image(d, edgeImg);
end
```

Streams from target peripherals to MATLAB

Generate code for peripheral access



Standalone deployment mode

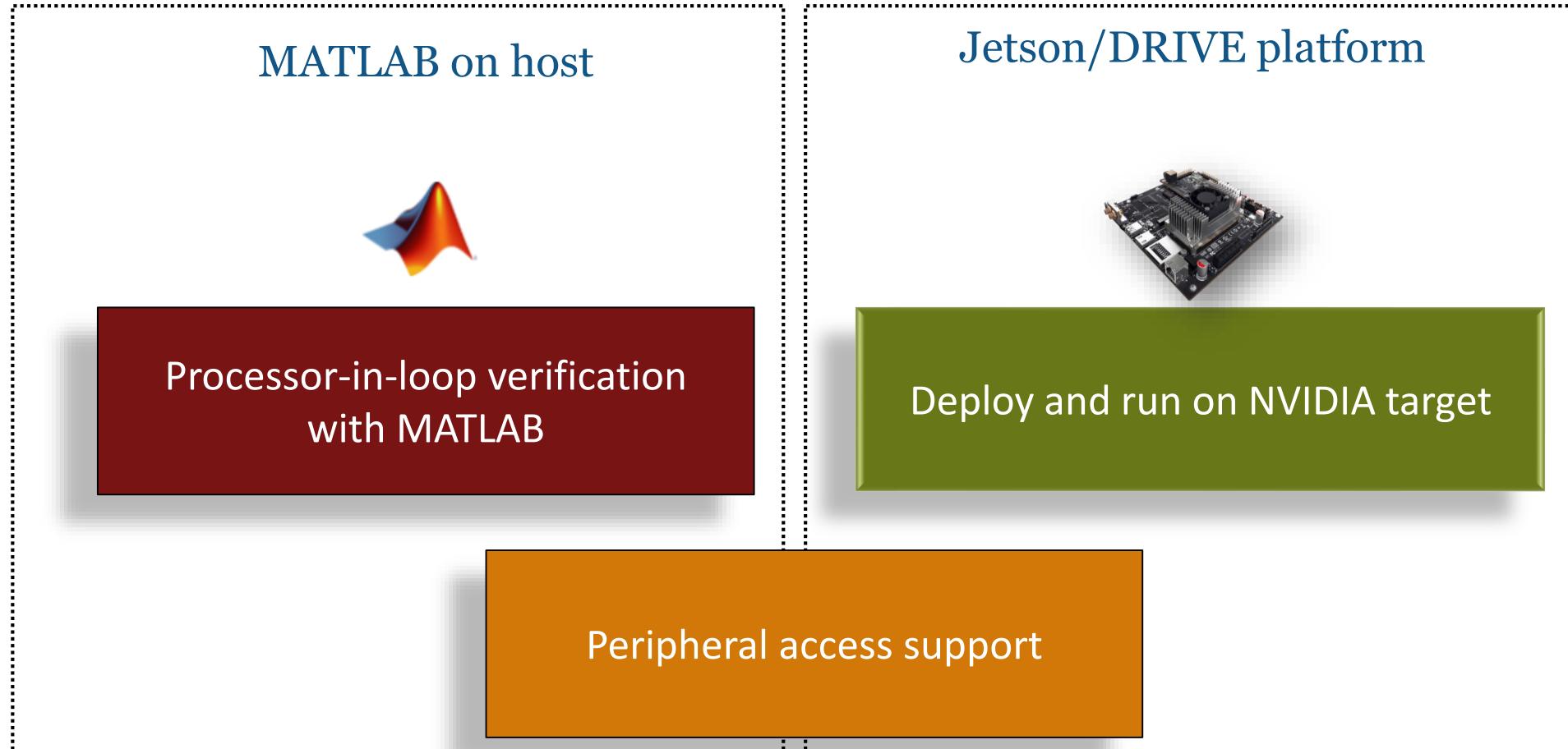
```
function sobelEdgeDetection()
%#codegen

% To enable code generation for hardware interfaces
hwobj = jetson;
w = webcam(hwobj,1);
d = imageDisplay(hwobj);

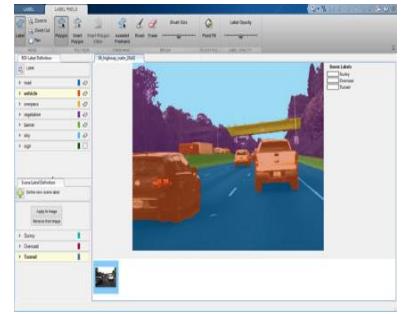
% Main loop
kern = [1 2 1; 0 0 0; -1 -2 -1];
for k = 1:1800
    % Capture the image from the webcam on hardware.
    img = snapshot(w);
    % finding horizon
    h = conv2(img(:,:,1));
    v = conv2(img(:,:,2));
    % Finding magnitude
    e = sqrt(h.*h + v.*v);
    % Threshold the edge
    edgeImg = uint8((e > 100) * 255);
    % Display edge image.
    image(d,edgeImg);
end
```

Generates interface code for target peripherals

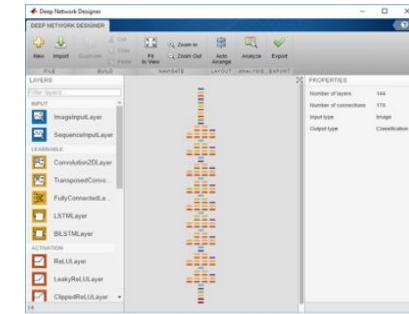
GPU Coder enables hardware prototyping and system integration



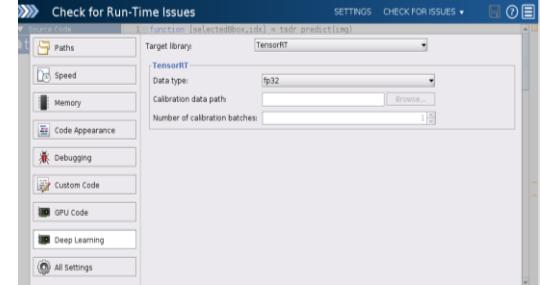
In Summary



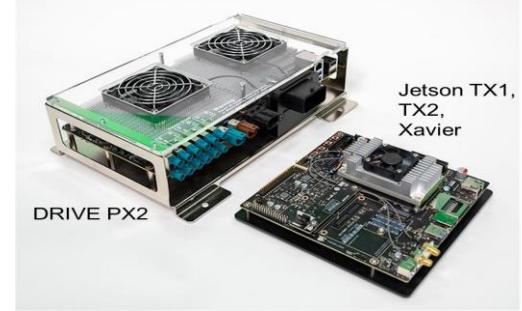
Ground Truth Labeling



Network Design and Training



CUDA and TensorRT Code Generation



Jetson Xavier and DRIVE Xavier Targeting

Network design

Platform Productivity: Workflow automation, ease of use

Framework Interoperability: ONNX, Keras-TensorFlow, Caffe

Deployment

Optimized CUDA and TensorRT code generation
Jetson Xavier and DRIVE Xavier targeting
Processor-in-loop(PIL) testing and system integration

Thank You