

MATLAB EXPO 2019

ディープラーニングアプリケーション
の組み込みGPU/CPU実装

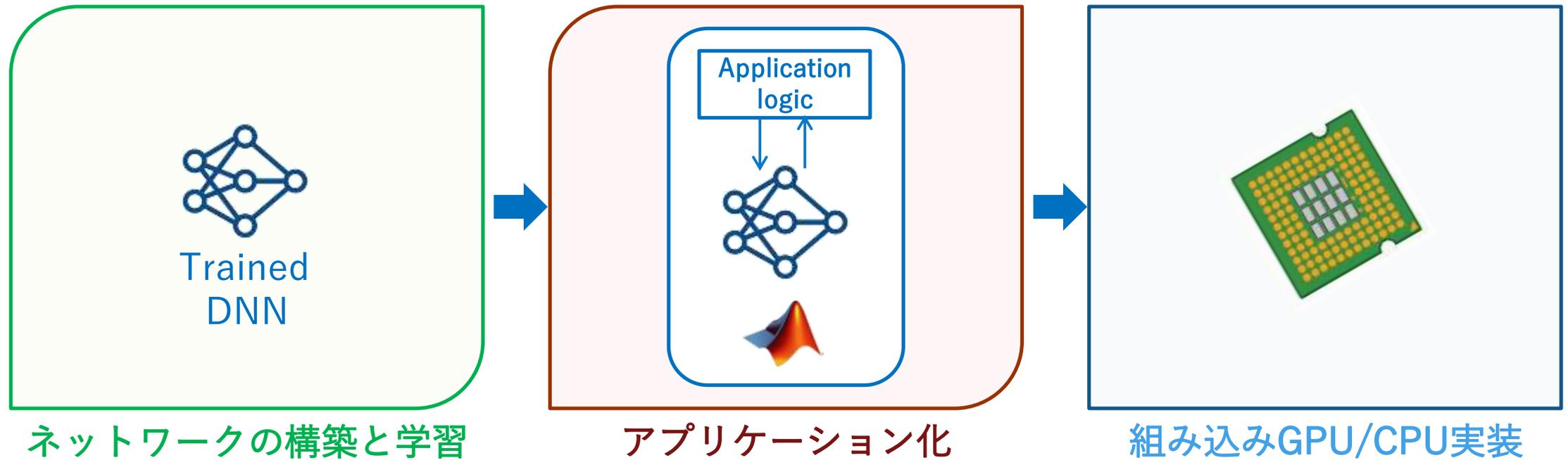
アプリケーション エンジニアリング部
町田 和也



アジェンダ

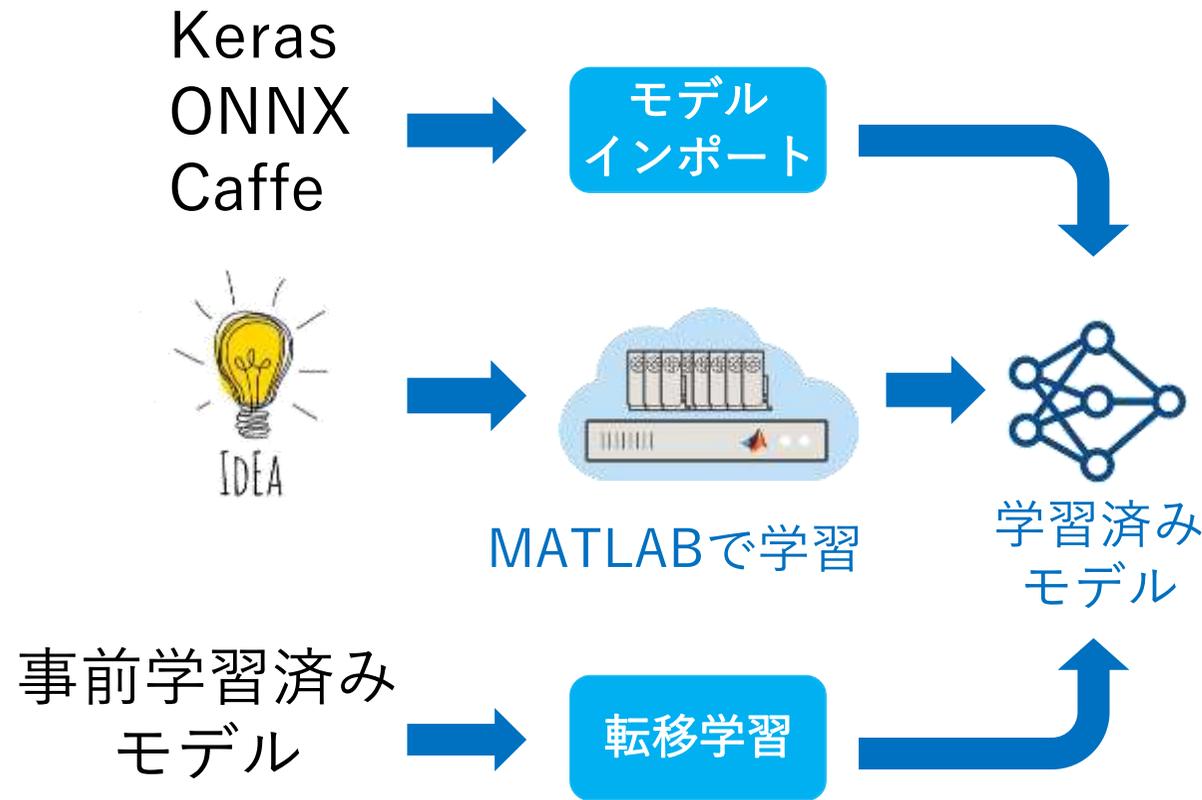
- MATLAB Coder/GPU Coderの概要
- ディープニューラルネットワークの組み込み実装ワークフロー
- パフォーマンスに関して
- まとめ

ディープラーニングワークフローのおさらい



ディープラーニングワークフローのおさらい

ネットワークの構築と学習



■ MATLABによる構築

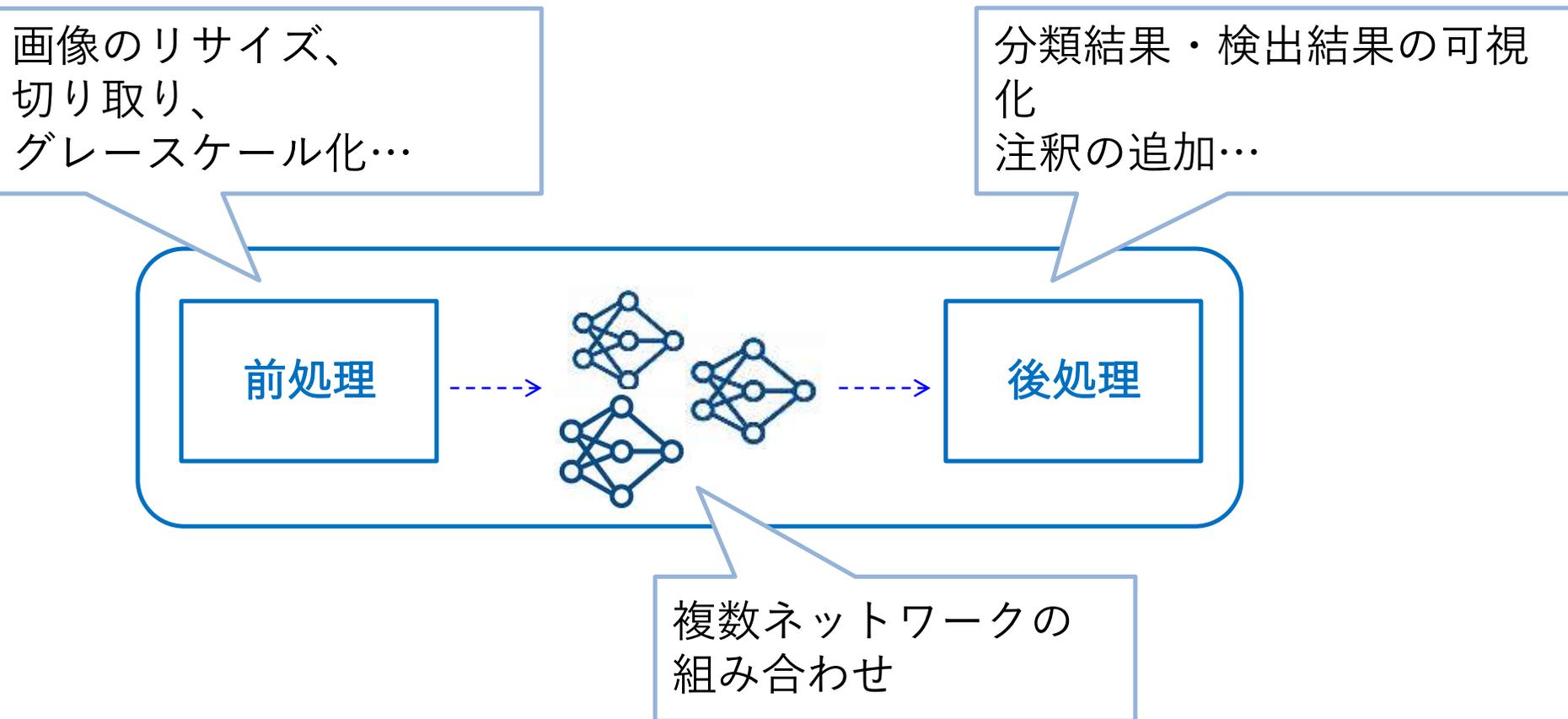
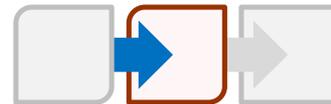
- 多量のデータの取り扱い
- ラベリングの自動化
- 様々なモデルへのアクセス

■ MATLABによる学習

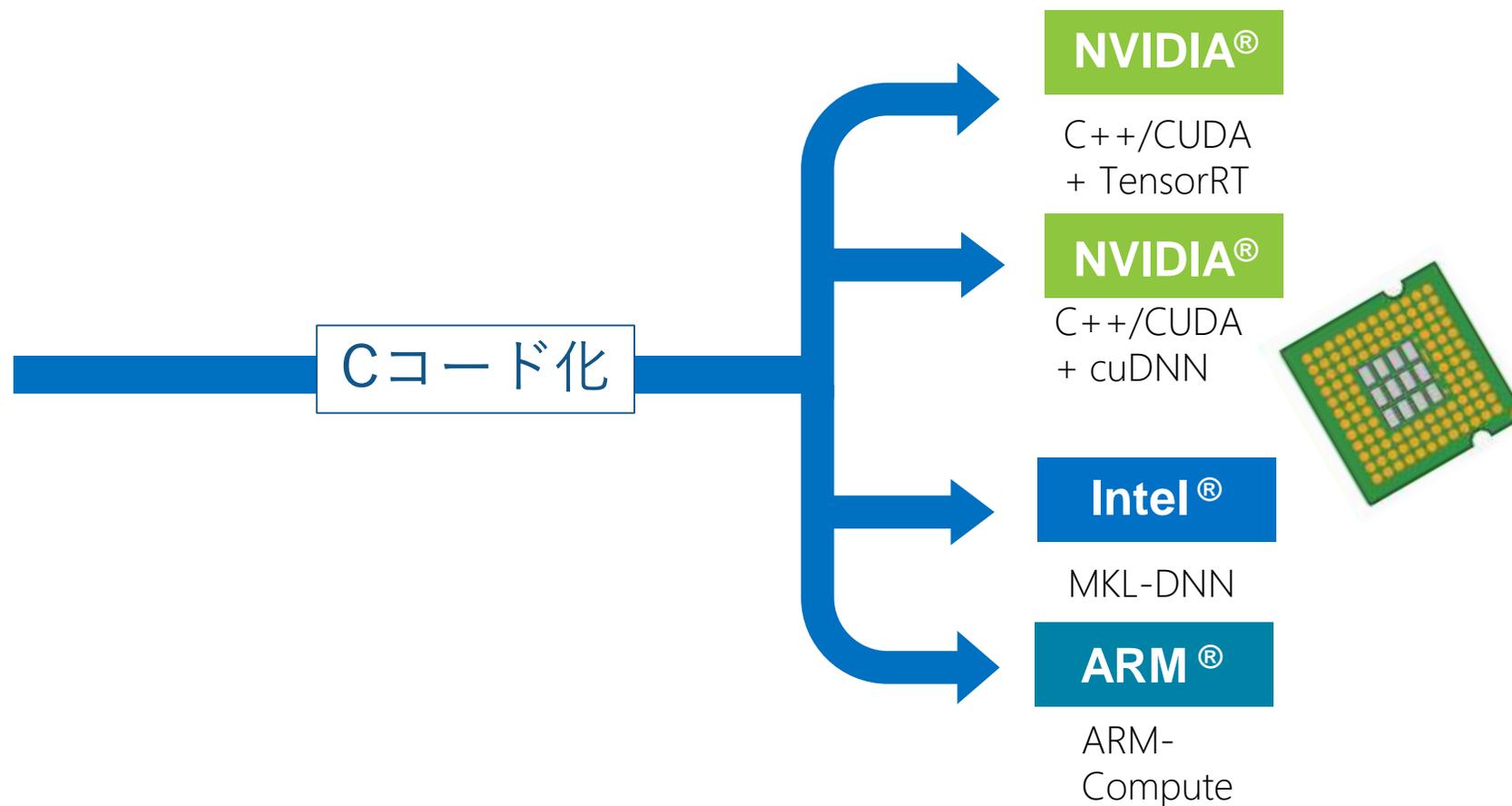
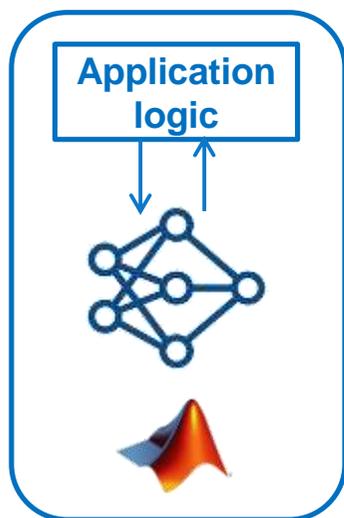
- GPUによる学習の高速化
- クラスタによる分散処理

ディープラーニングワークフローのおさらい

アプリケーション化



ディープラーニングワークフローのおさらい マルチプラットフォーム環境への実装



Cコード化のハードル



Vectorized MATLAB

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

Cコード化



CUDA コード

```
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);
}
```

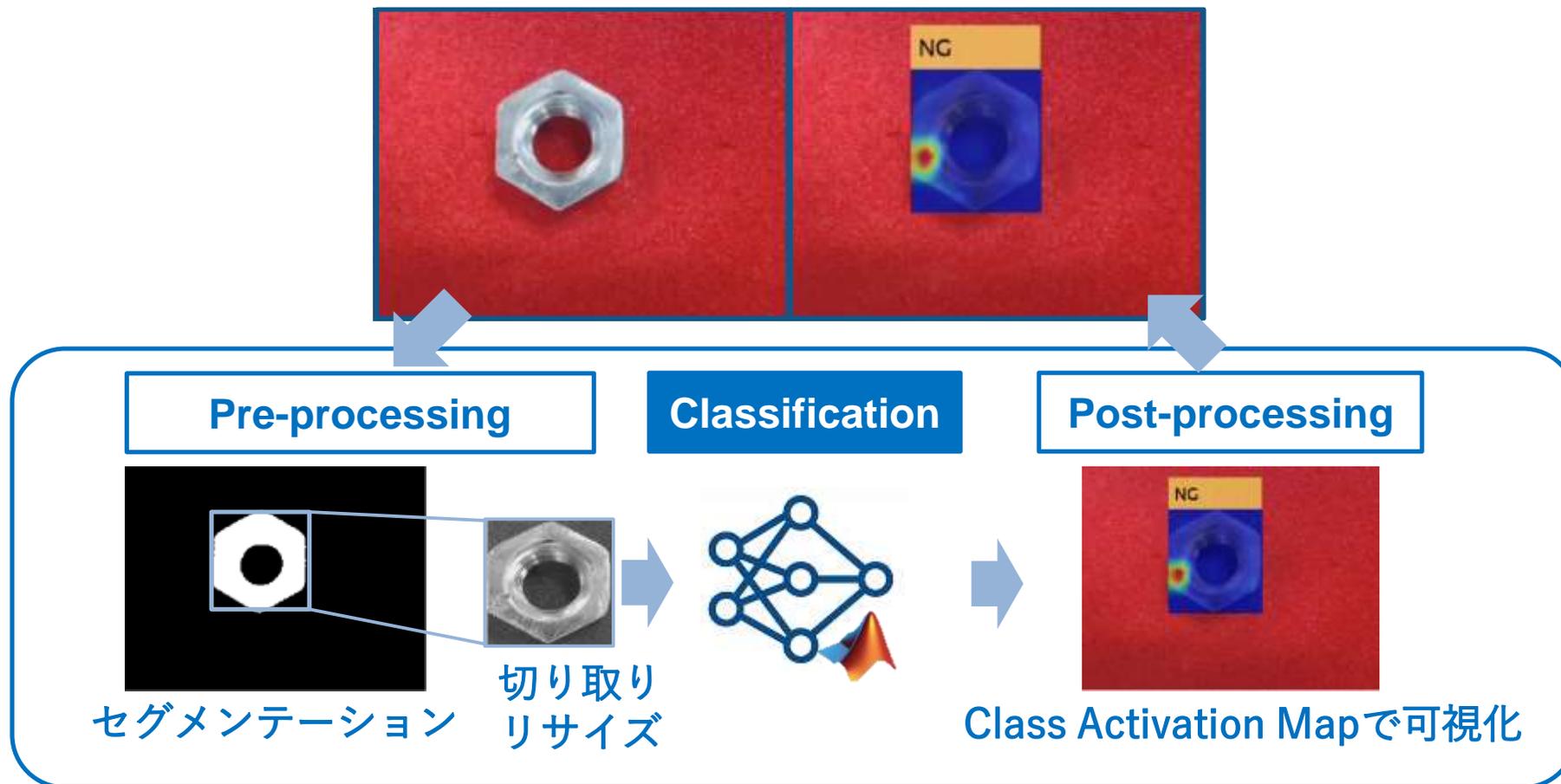
手書きによるバグ混入
等価性の維持が困難

GPUは非常に強力なハードウェアですが、…
プログラミングは専門の知識が必要となります。

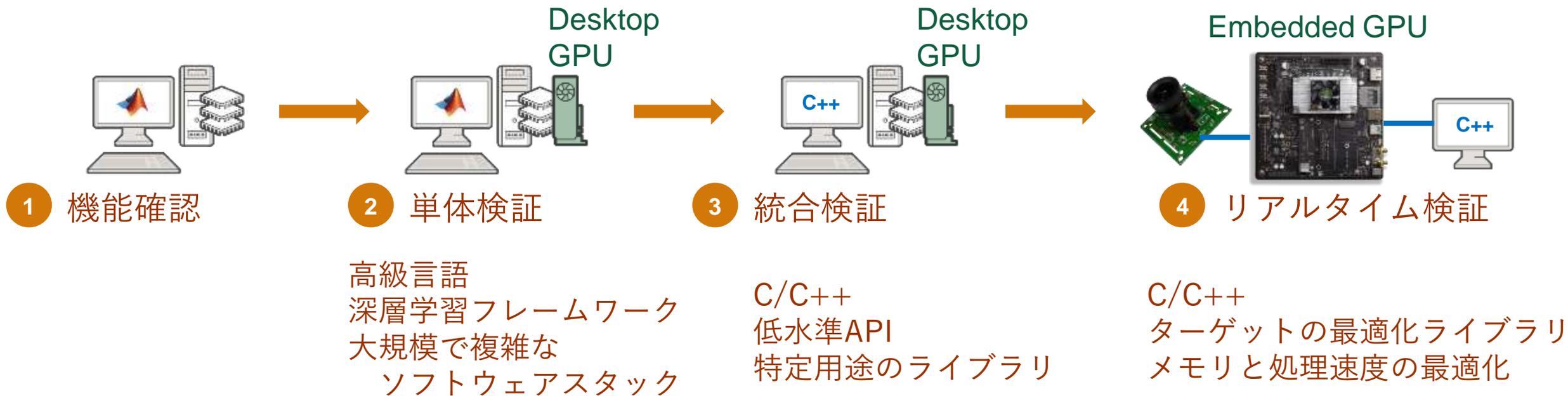
Cコード化のハードル



推論以外の処理もコード化する必要があります



組み込み実装の課題

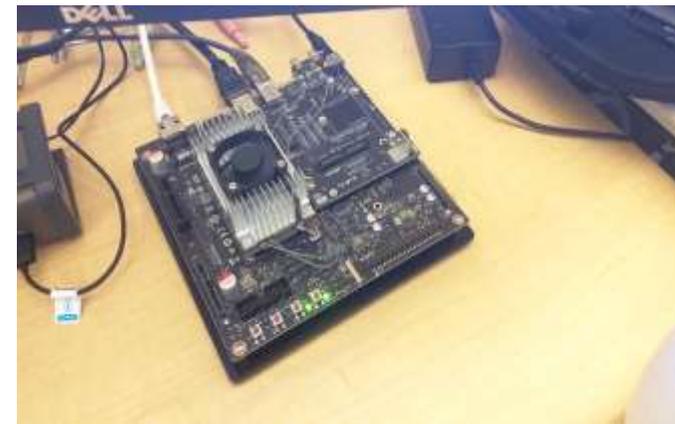
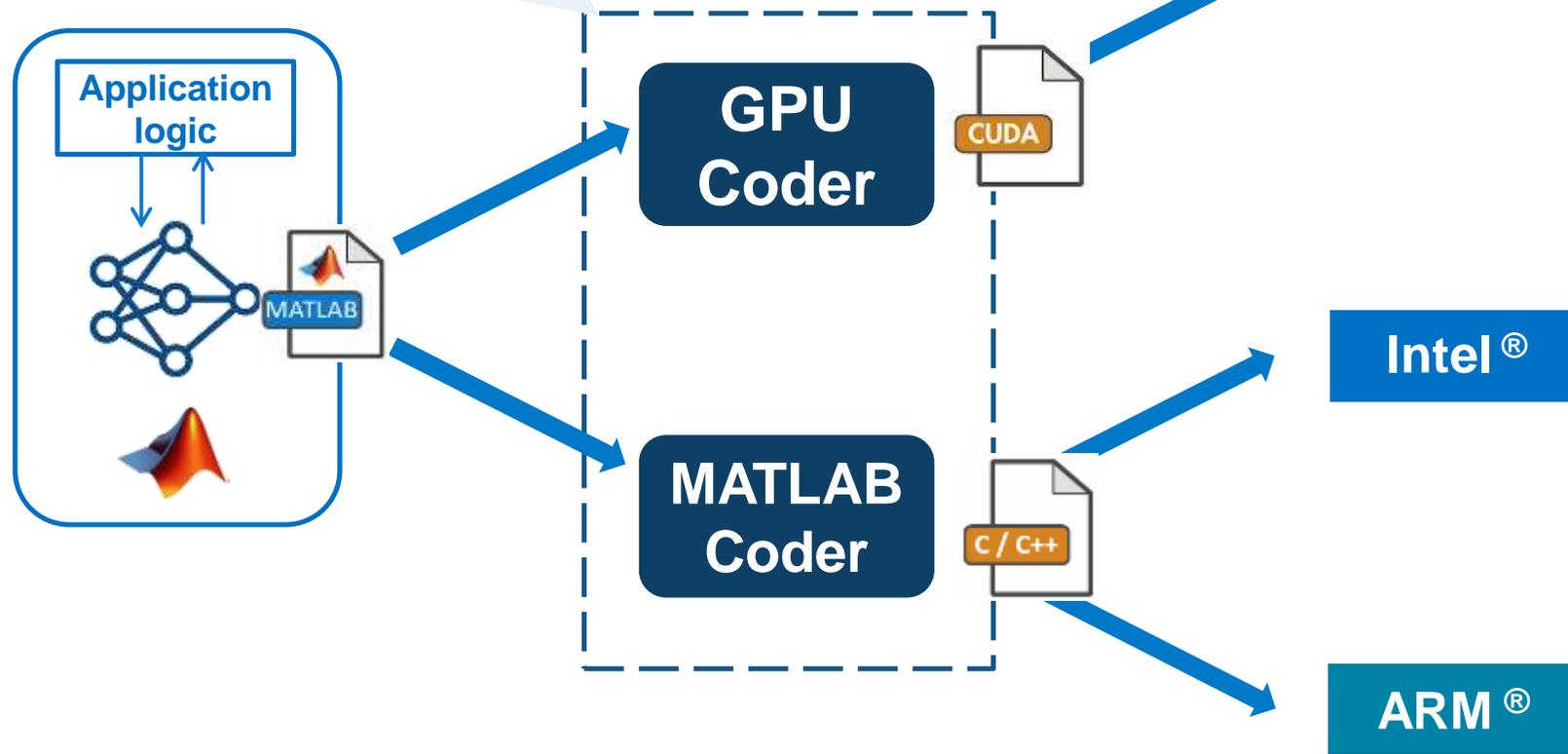


課題

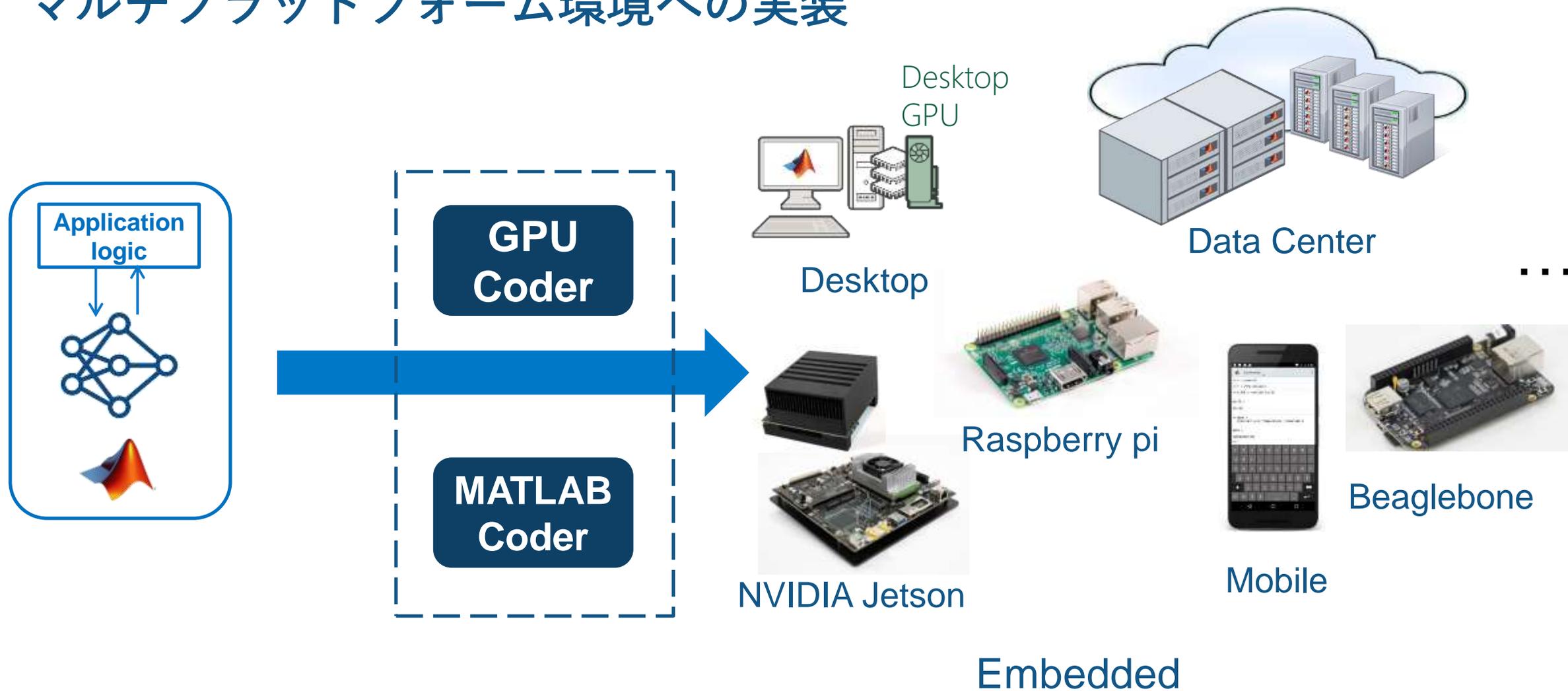
- 複数のライブラリとパッケージの統合
- 複数の実装の検証と維持
- アルゴリズムとベンダーロックインの検討

MATLAB Coder & GPU Coderによる実装ソリューション

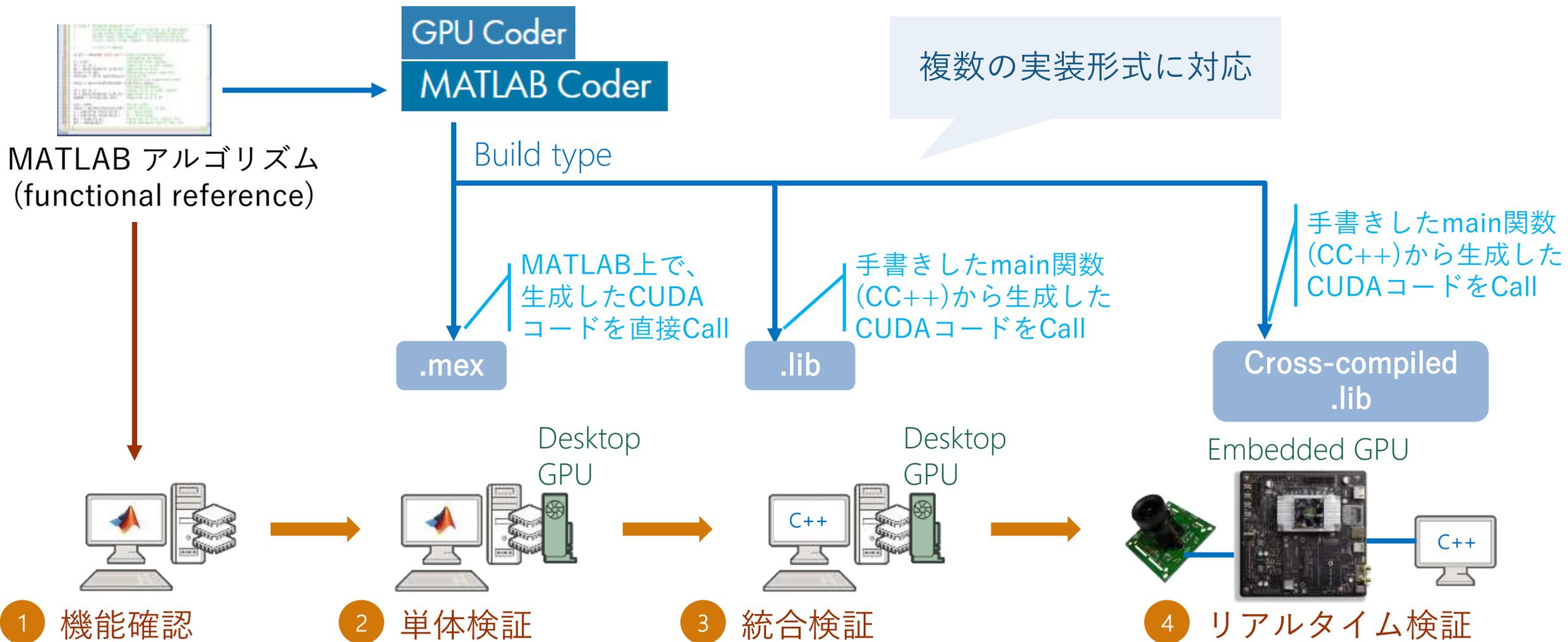
- 前処理・後処理を含めたコードの自動生成
- 生成されるコードの最適化
- 複数のターゲットへの実装



MATLAB Coder & GPU Coderによる実装ソリューション マルチプラットフォーム環境への実装



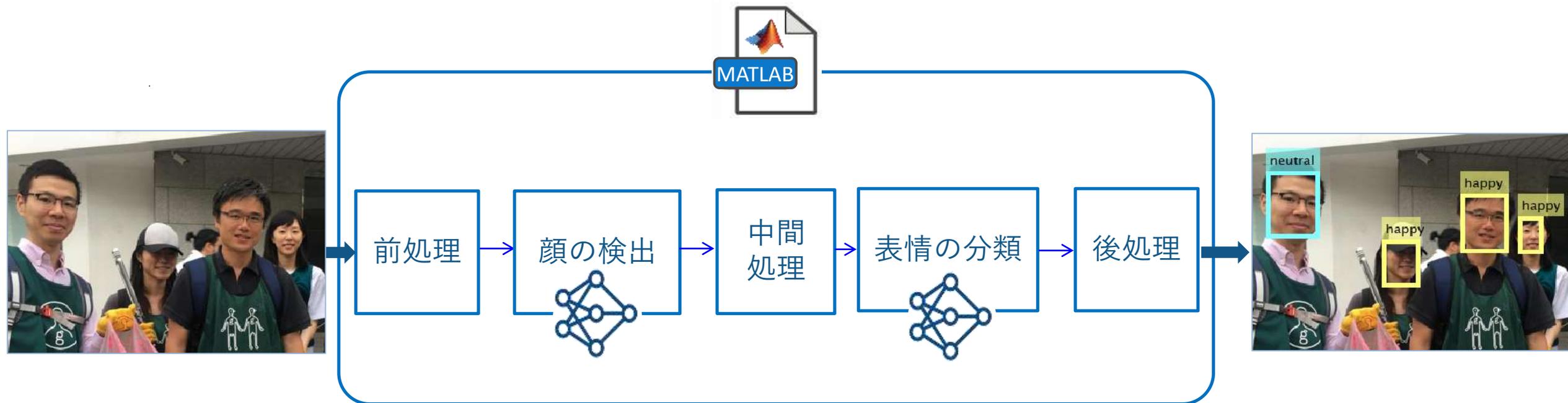
MATLAB Coder & GPU Coderによる組み込み実装



アジェンダ

- MATLAB Coder/GPU Coderの概要
- ディープニューラルネットワークの組み込み実装ワークフロー
- パフォーマンスに関して
- まとめ

【Demo】顔検出と表情分類アプリケーションのJetson実装

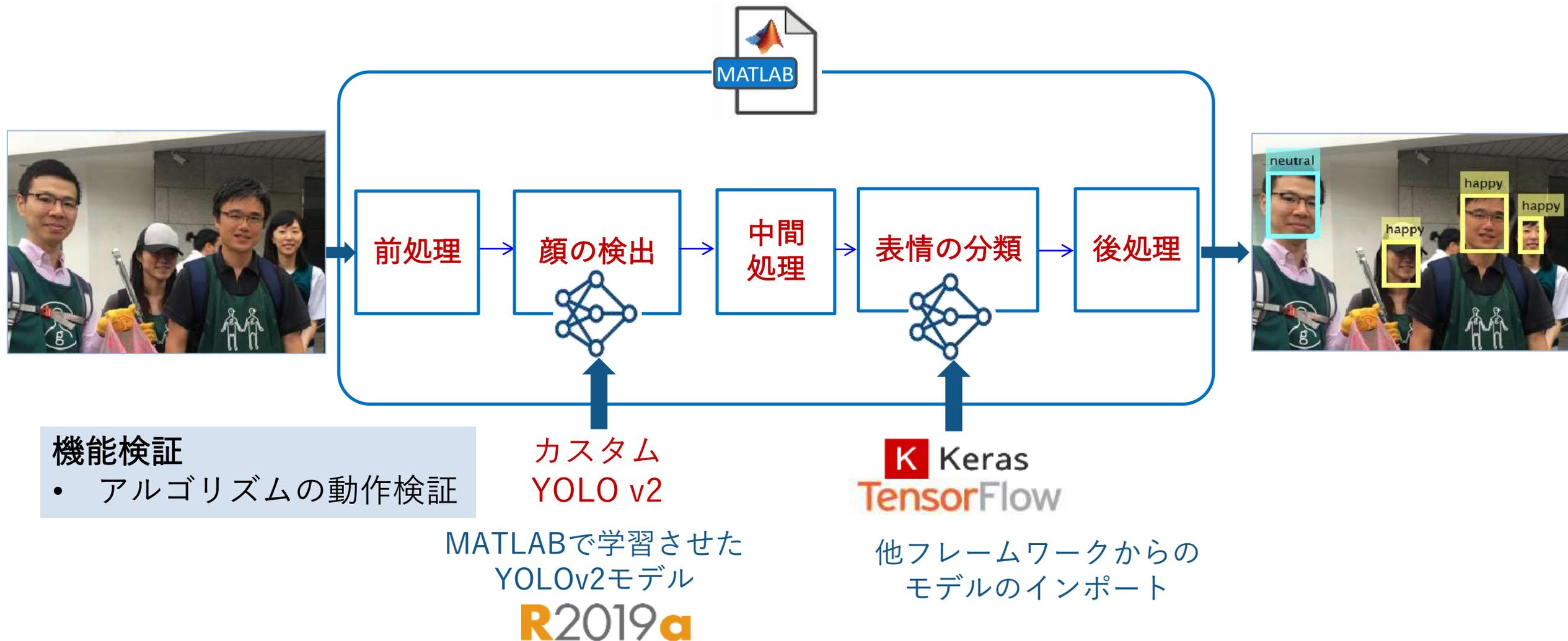


実装ワークフロー

- ① MATLABで検証
- ② コード生成→ホスト環境での検証
- ③ コード生成→組み込み環境での検証



① MATLABで検証



①MATLABで検証

顔の検出



R2019a

- YOLO v2による物体検出
 - YOLO(You Only Look Once)とは？
 - リアルタイム向けの物体検出モデル (Faster R-CNNの100倍程度高速)
 - 学習からCUDAコード生成まで対応
 - 自身のデータセットに合わせたカスタム YOLOv2モデルの学習から実装までサポート

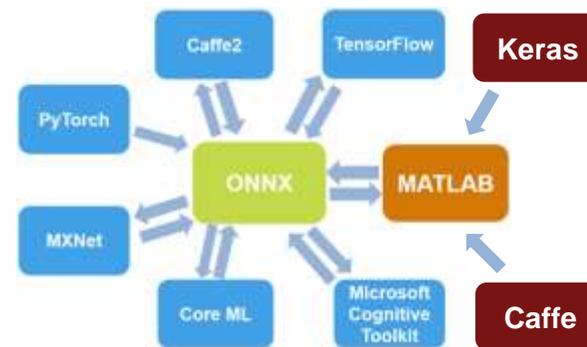


```
detector =
  trainYOLOv2ObjectDetector(trainingData,lgraph,options)
```

表情の分類

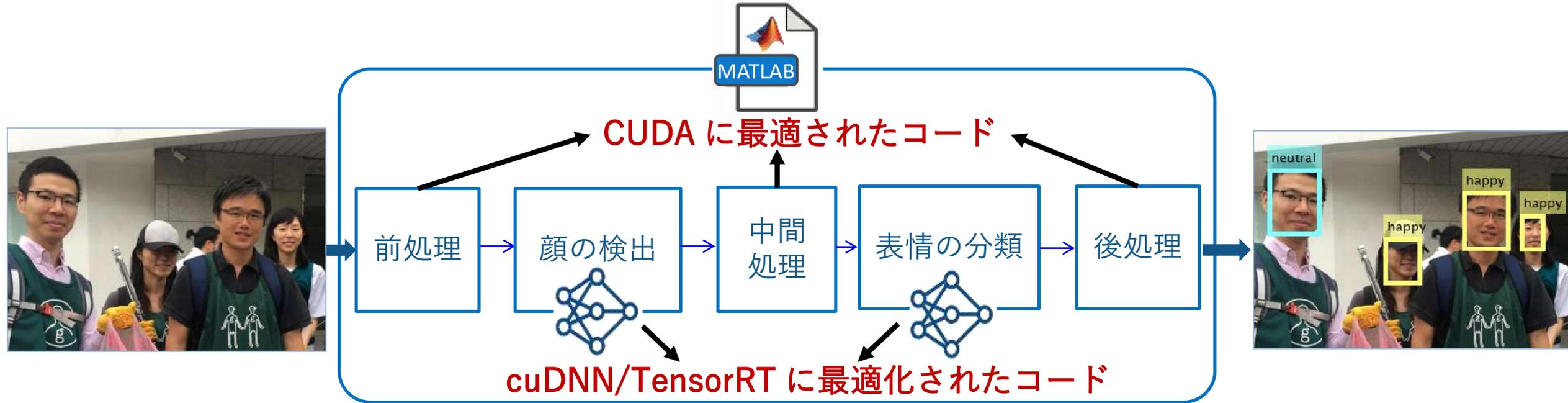


- 他フレームワークとの連携
 - Keras, Caffeモデルの取り込み
 - ONNXフォーマットを介したモデルのやり取り



```
net = importKerasNetwork(modelfile)
net = importCaffeNetwork(protofile,datafile)
net = importONNXNetwork(modelfile, ...
  'OutputLayerType',outputtype)
```

②コード生成→ホスト環境での検証



単体検証

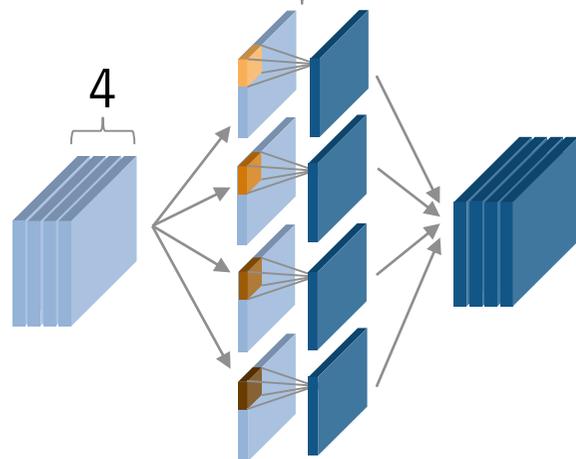
- コード生成可能か確認
→生成不可の場合は生成可能な形に修正

②コード生成→ホスト環境での検証 コード修正が必要な例

グループ化畳み込み層を通常の畳み込み層に変換

R2019aでは
コード生成未対応

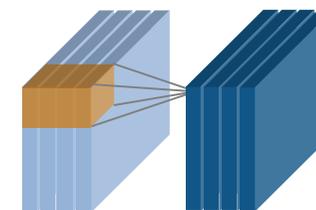
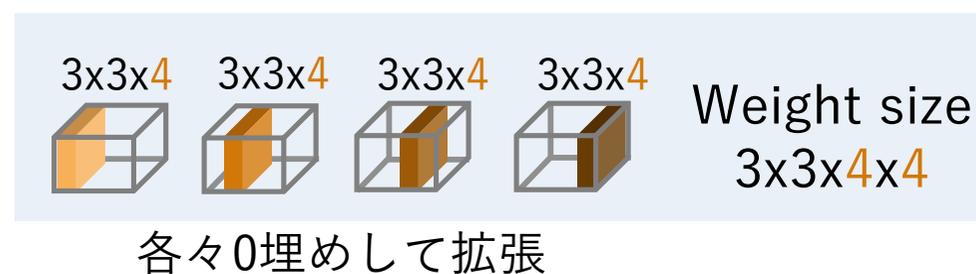
groupedConvolution2dLayer



例：
4チャンネルの入力を
1チャンネル毎に
グループ化した場合

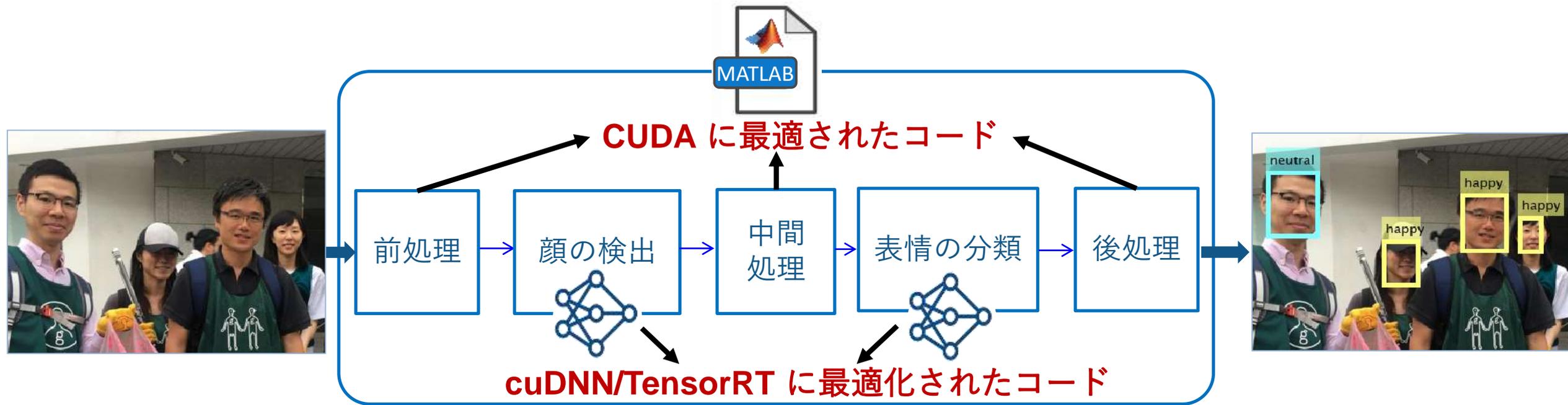
convolution2dLayer

コード生成
対応



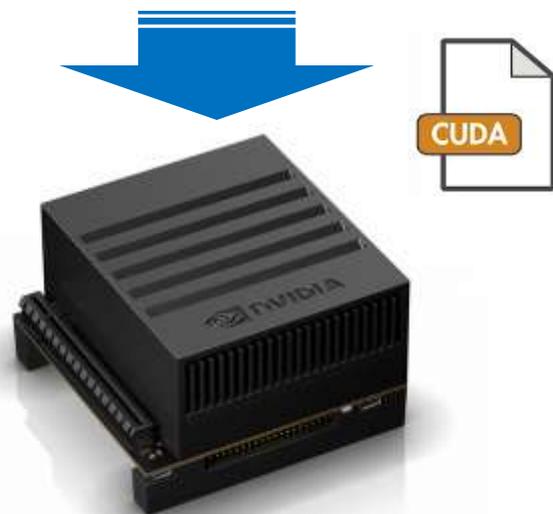
計算結果が等価になるように
通常の畳み込み層に重みとバイアスを移植

③コード生成→組み込み環境での検証

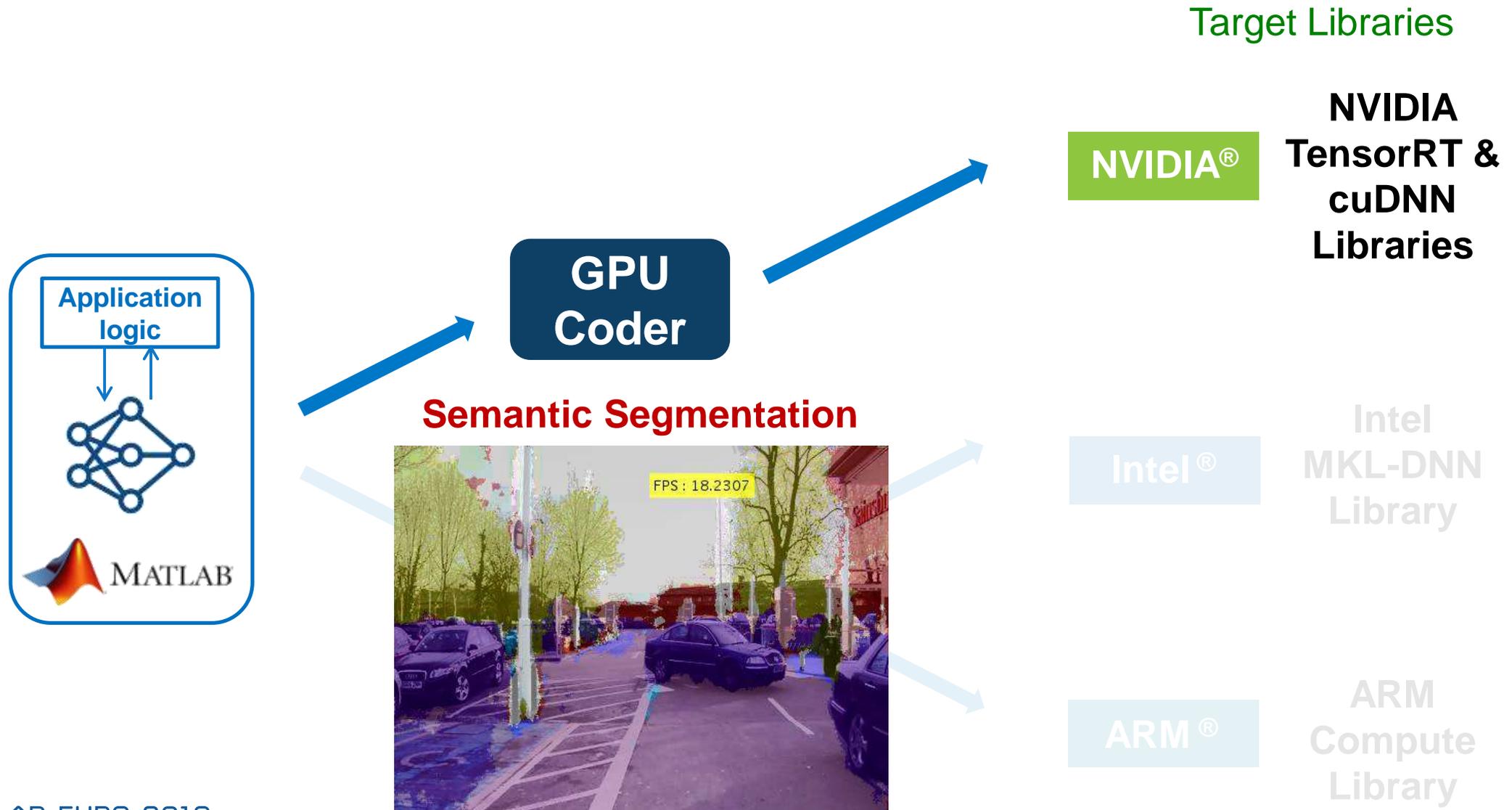


リアルタイム検証

- ターゲット環境でパフォーマンス確認

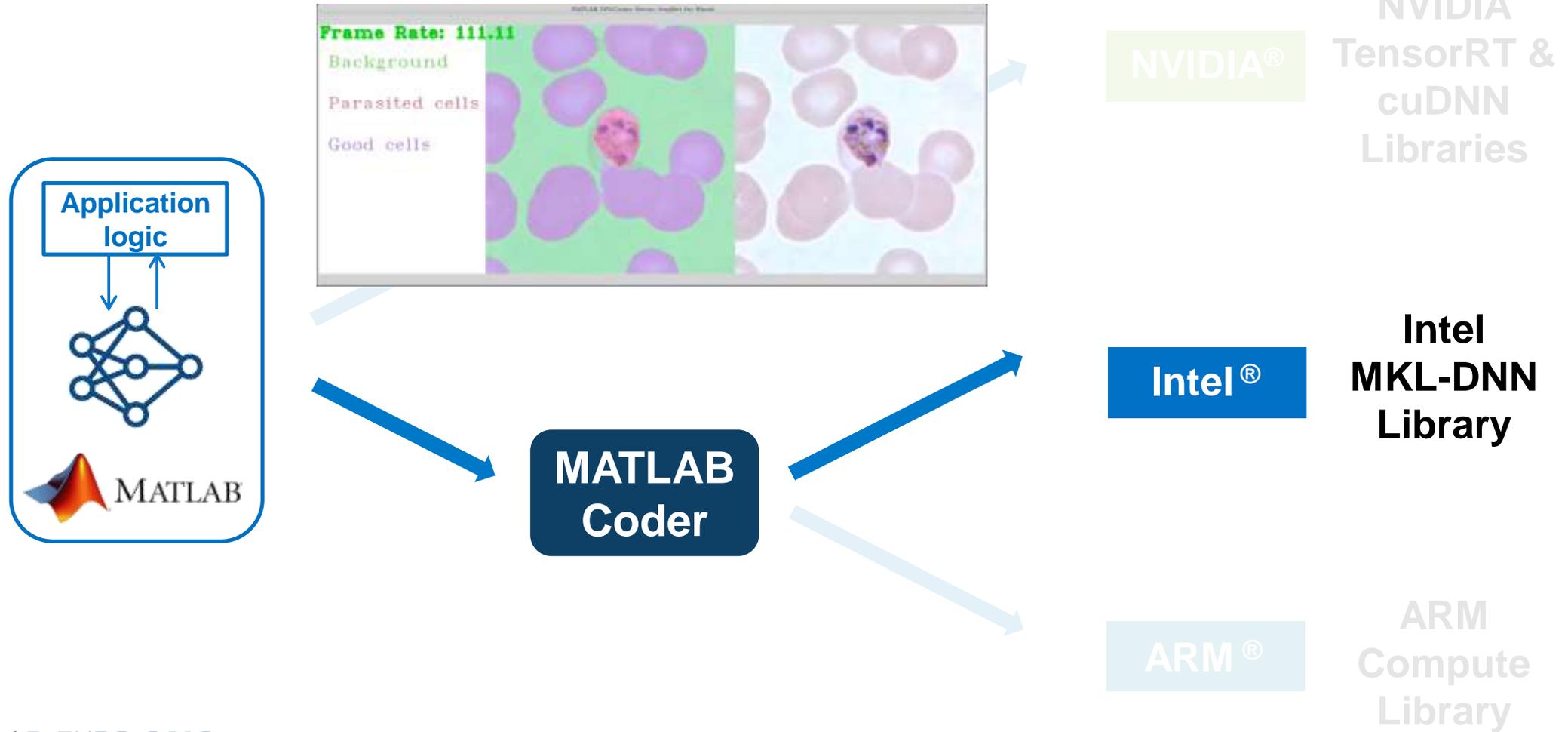


NVIDIA GPUでの実装例



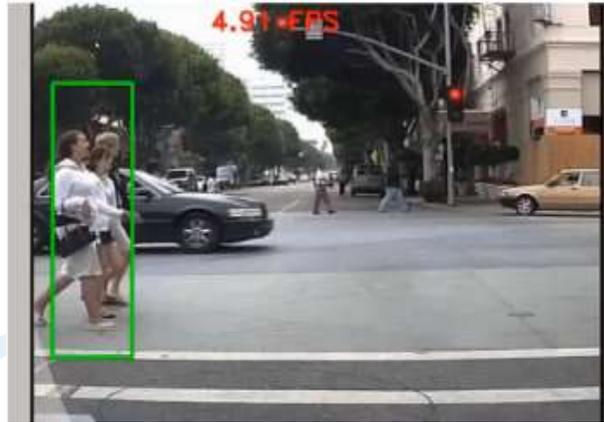
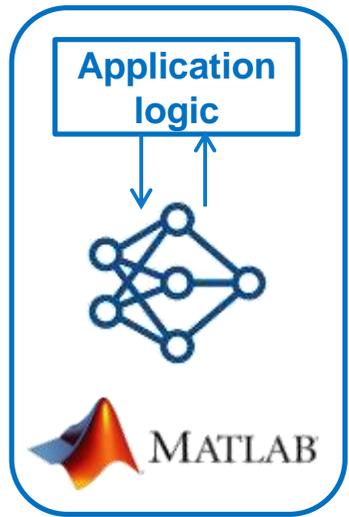
Intel CPUでの実装例

Blood Smear Segmentation



ARMプロセッサでの実装例

Pedestrian Detection



**MATLAB
Coder**

Target Libraries

NVIDIA®

NVIDIA
TensorRT &
cuDNN
Libraries

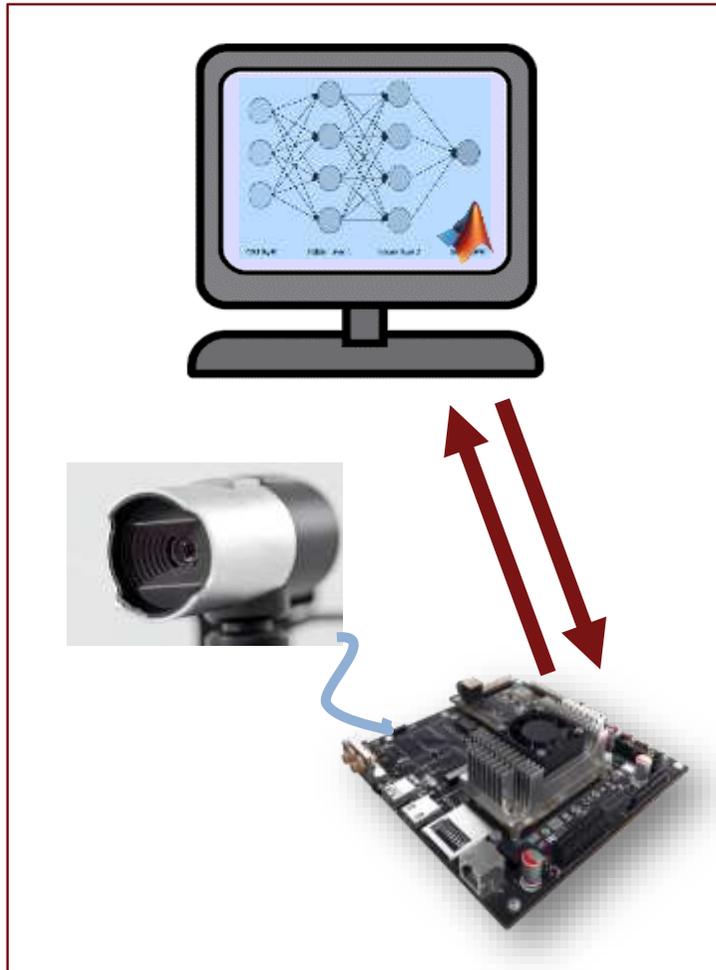
Intel®

Intel
MKL-DNN
Library

ARM®

**ARM
Compute
Library**

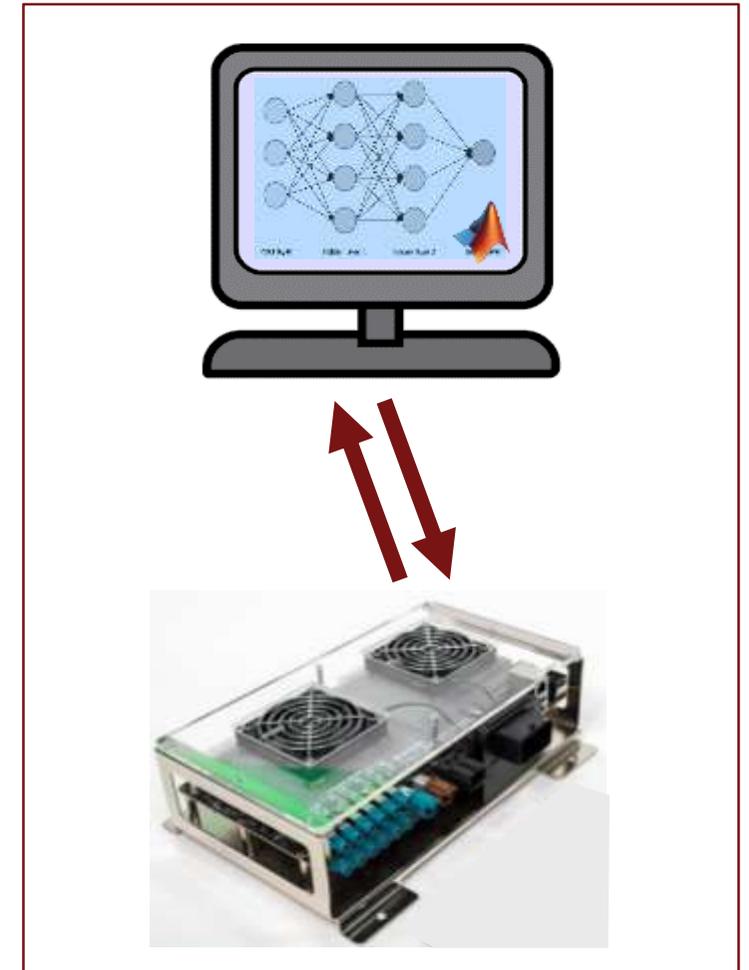
ハードウェアへのアクセス



MATLABから、
周辺機器にアクセス



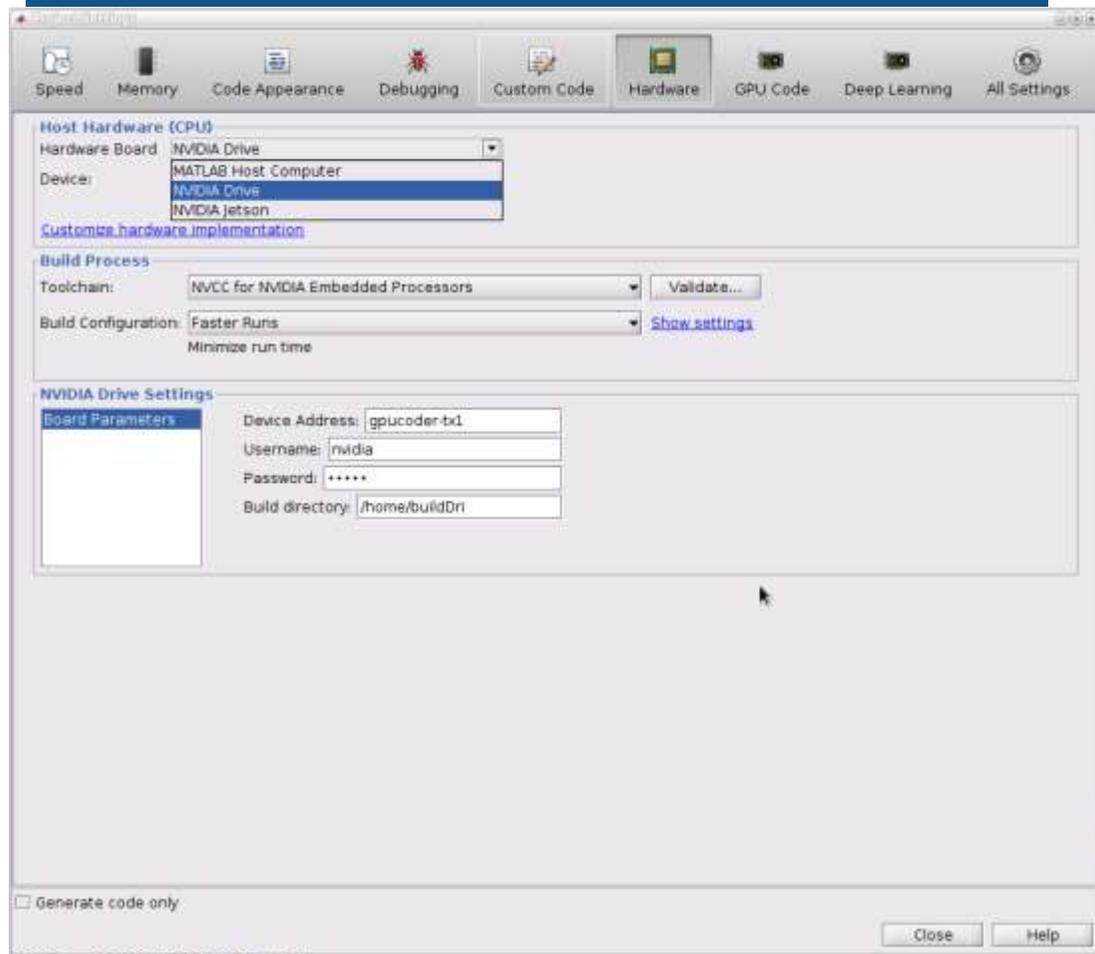
スタンドアロン
アプリケーションとして実装



Processor-in-Loop 検証

ターゲットへの実装方法

専用アプリによるコード生成と実装



コマンドによるコード生成と実装

```
%% Deploy and launch through NVIDIA HSP
```

```
%% setup hardware object
% create jetson/drive hardware object with IP or hostname of jetson/drive
%also pass credentials for login
hwObj = jetson('gpcoder-tx2-2','ubuntu','ubuntu');
hwObj.setupCodegenContext;
```

```
%% setup codegen config object
% create congen config and connect to hardware object.
cfg_hsp = coder.gpuConfig('exe');
cfg_hsp.Hardware = coder.hardware(hwObj.BoardPref);
buildDir = '~/buildDir';
cfg_hsp.Hardware.BuildDir = buildDir;
```

```
%% add user written main files for building executable
% and generate/build the code.
cfg_hsp.CustomSource = 'driver_files_alexnet/main.cu';
cfg_hsp.CustomInclude = 'driver_files_alexnet/';
```

```
codegen -config cfg_hsp -args {im, coder.Constant(cnnMatFile)} alexnet_test
```

```
%% copy input and run the executable
hwObj.putFile('input2.txt', buildDir);
hwObj.putFile('synsetWords.txt', buildDir);
```

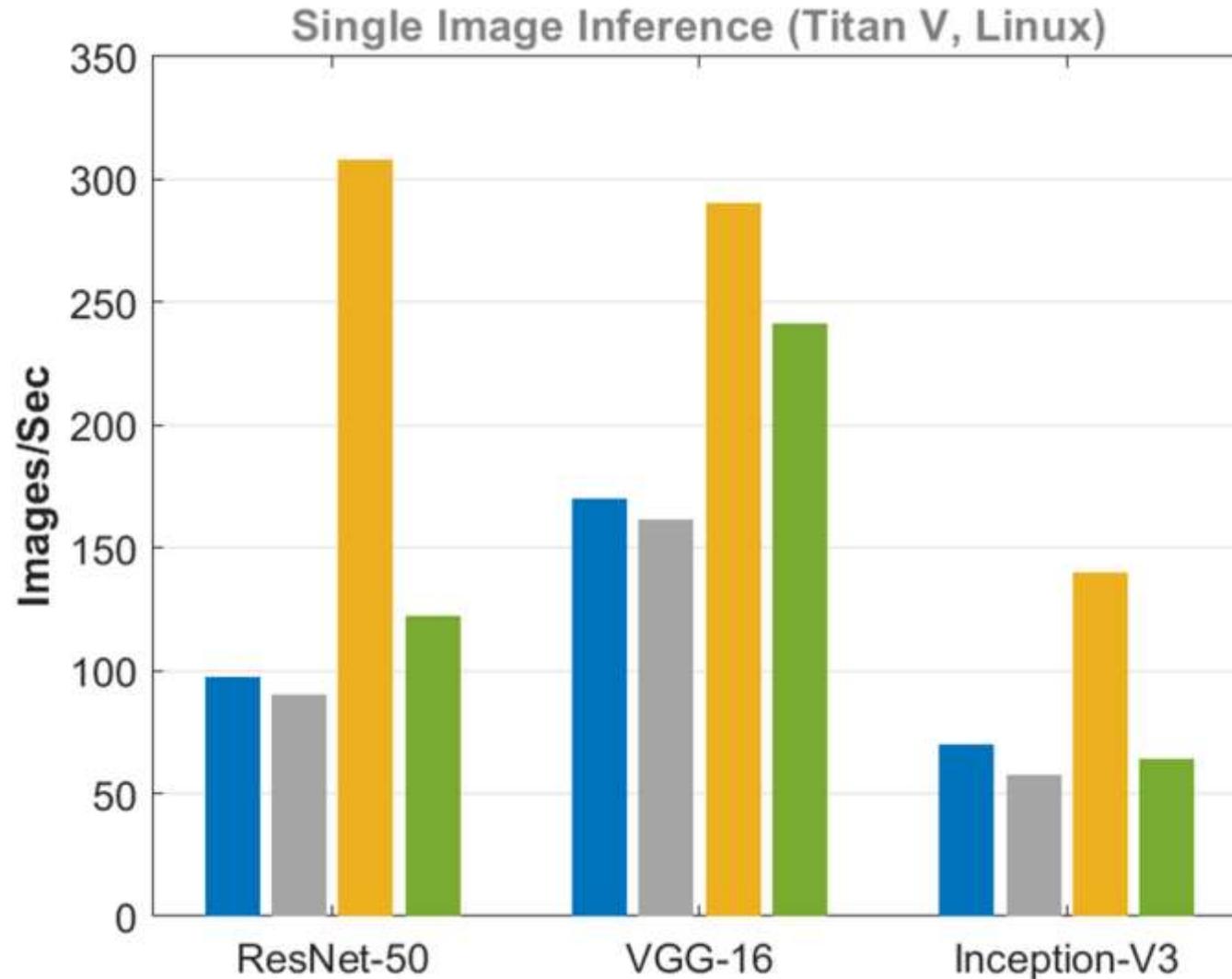
```
%execute on Jetson
hwObj.runExecutable([buildDir '/alexnet_test.elf'], 'input2.txt')
```

```
%% copy the output file back to host machine
hwObj.getFile([buildDir '/tOut.txt']);
```

アジェンダ

- MATLAB Coder/GPU Coderの概要
- ディープニューラルネットワークの組み込み実装ワークフロー
- パフォーマンスに関して
- まとめ

画像一枚に対する推論速度の比較



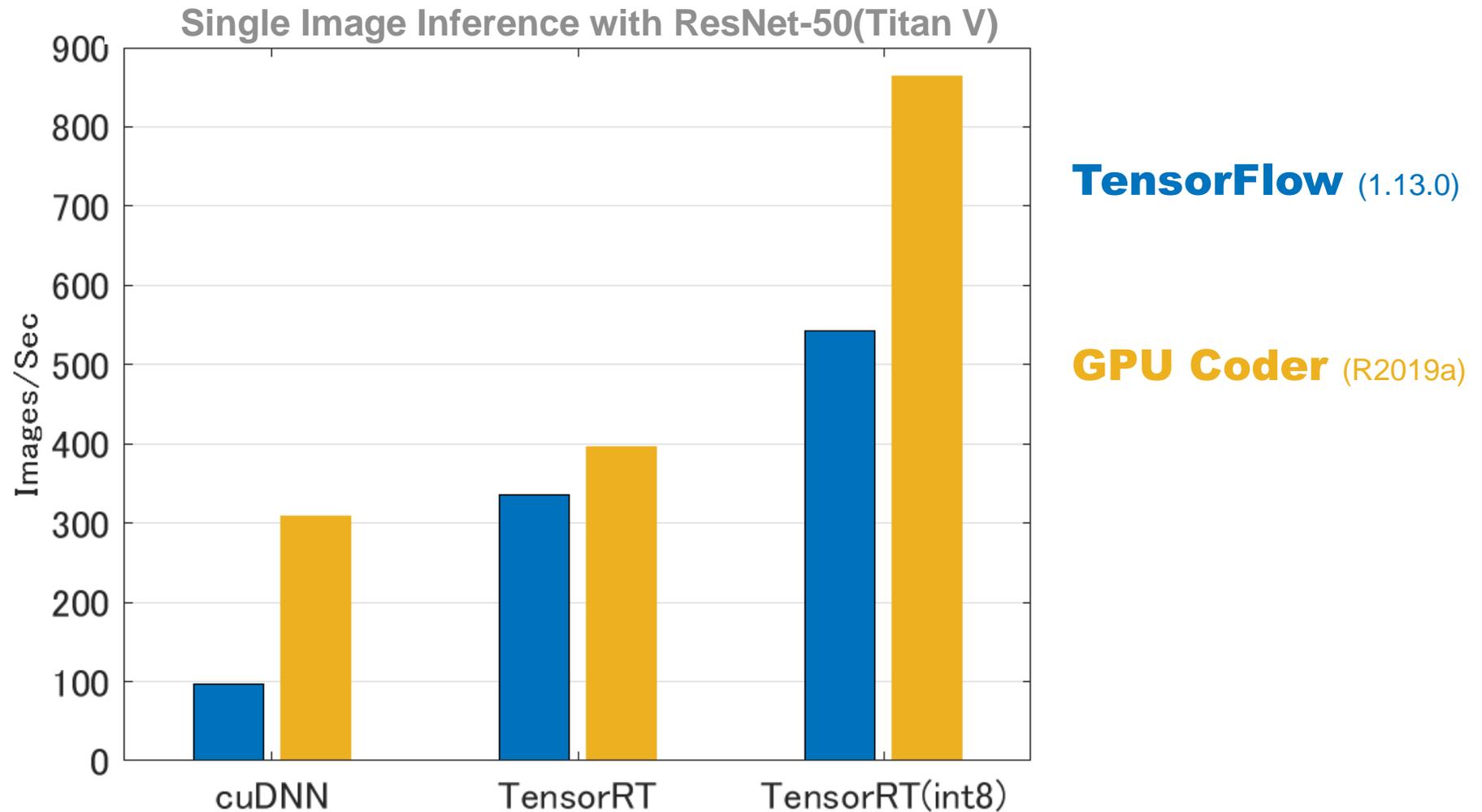
TensorFlow (1.13.0)

MXNet (1.4.0)

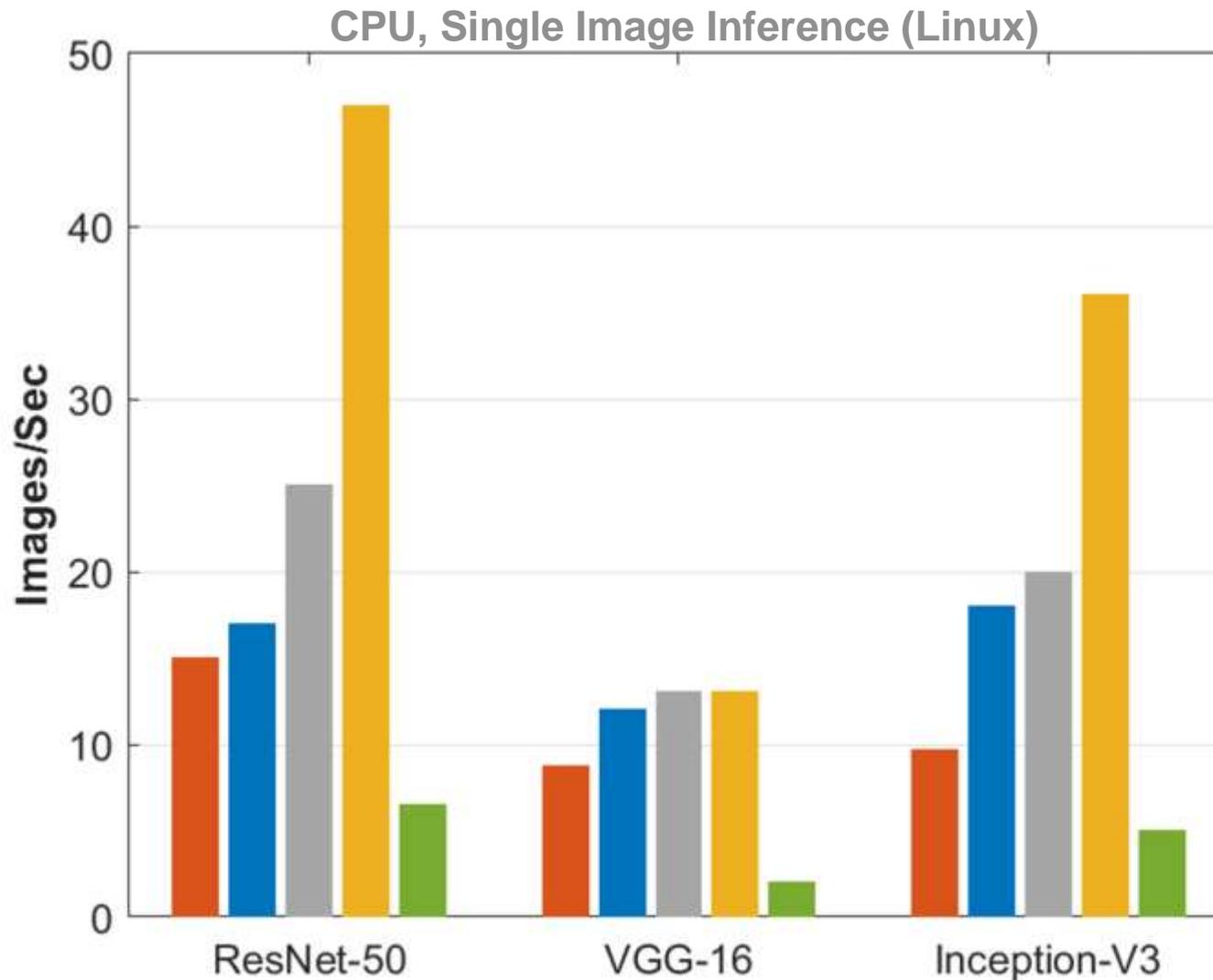
GPU Coder (R2019a)

PyTorch (1.0.0)

TensorRTによる高速化



CPUでの推論速度の比較



MATLAB

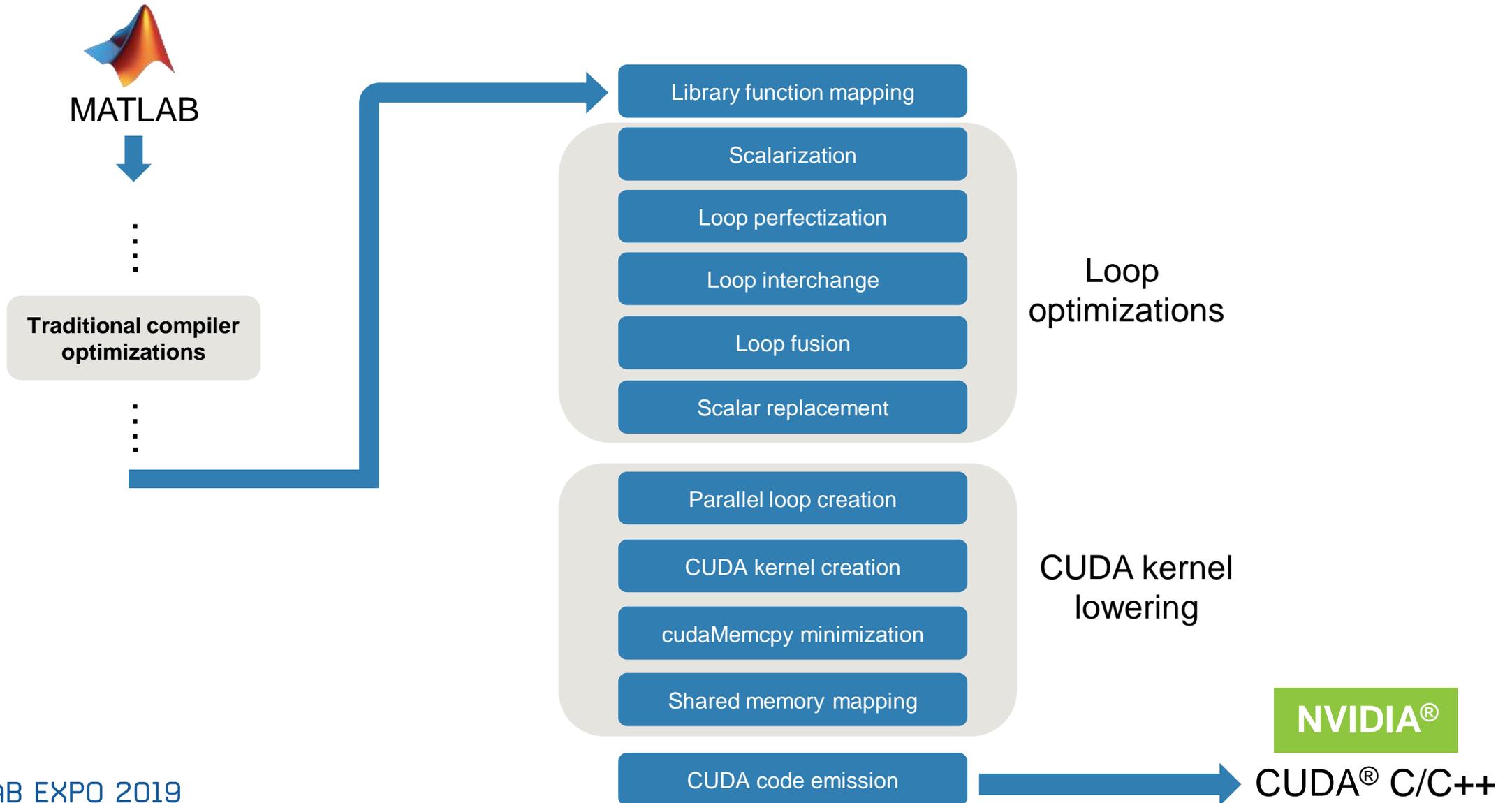
TensorFlow

MXNet

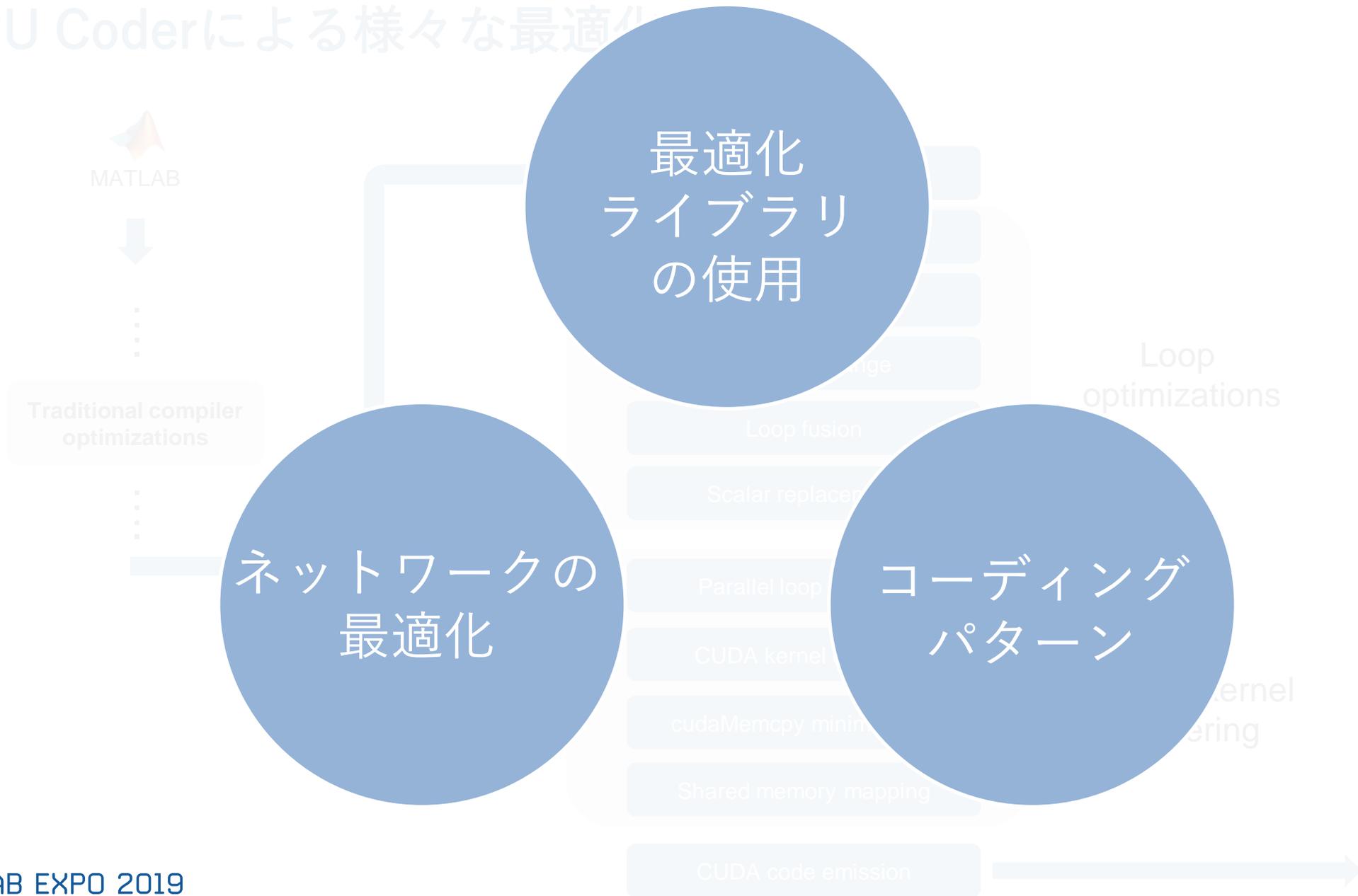
MATLAB Coder

PyTorch

GPU Coderによる様々な最適化



GPU Coderによる様々な最適化



最適化
ライブラリ
の使用

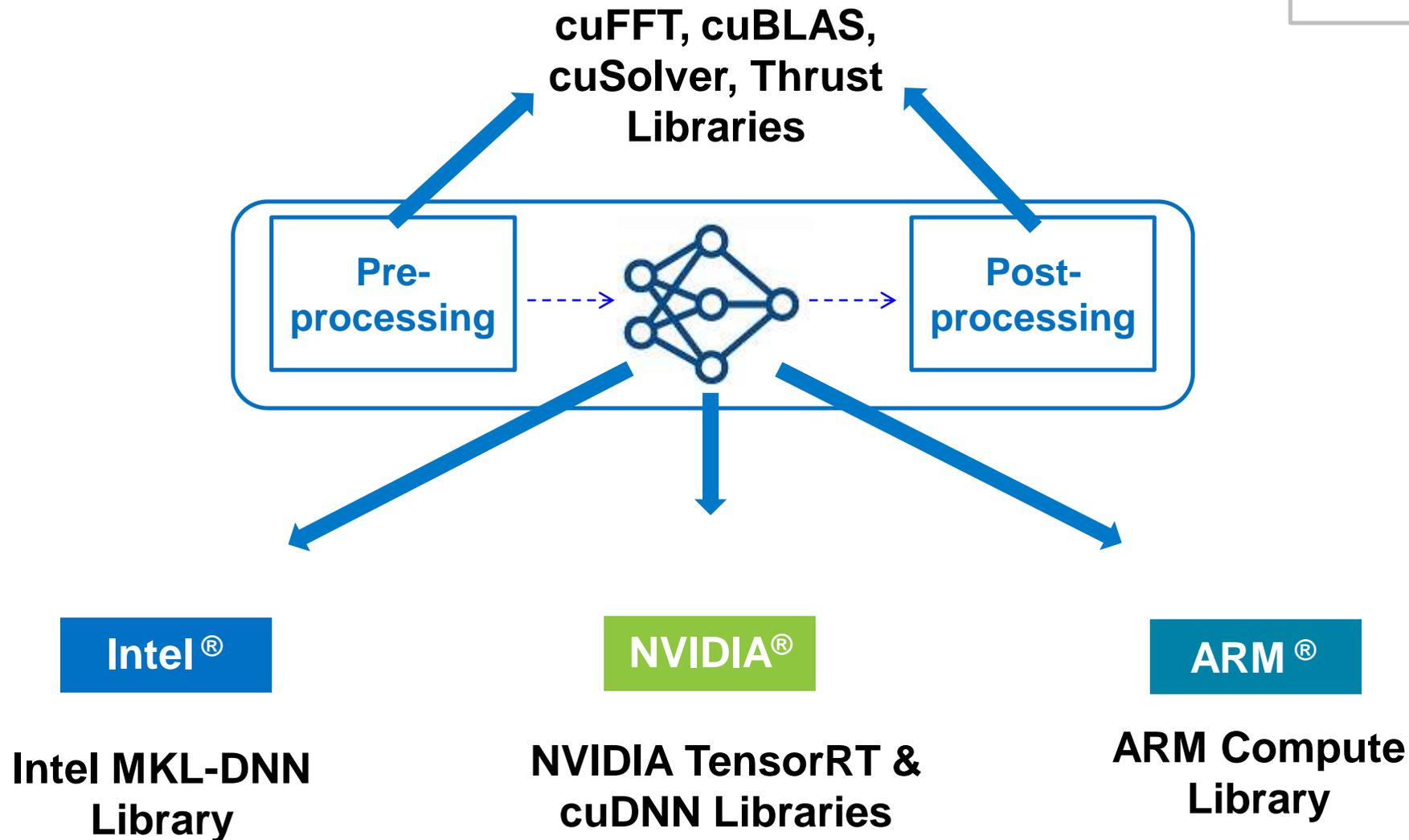
ネットワークの
最適化

コーディング
パターン

最適化ライブラリを用いたコード生成

Performance

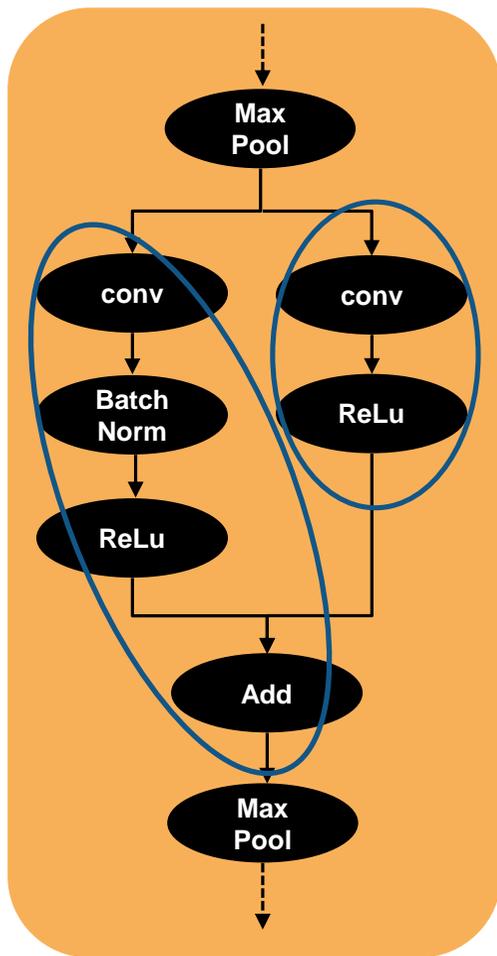
1. 最適化ライブラリの使用
2. ネットワークの最適化
3. コーディングパターン



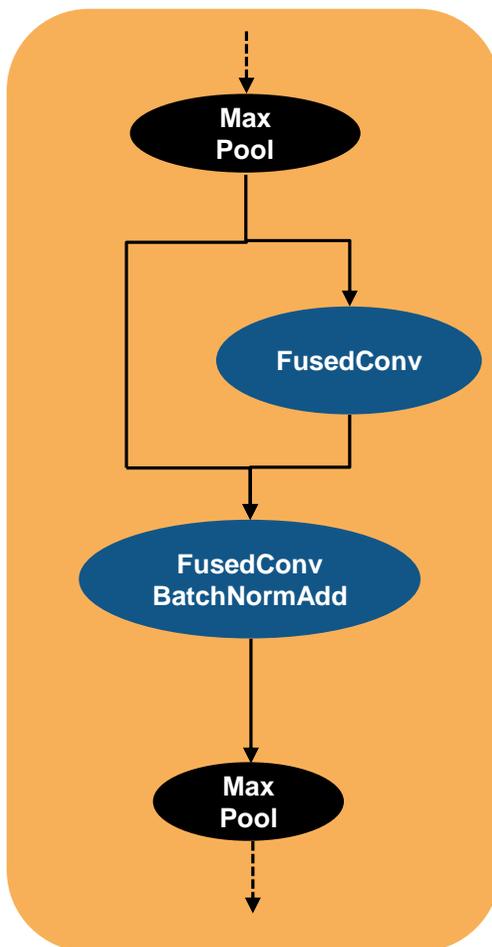
ディープニューラルネットワークの最適化

Performance

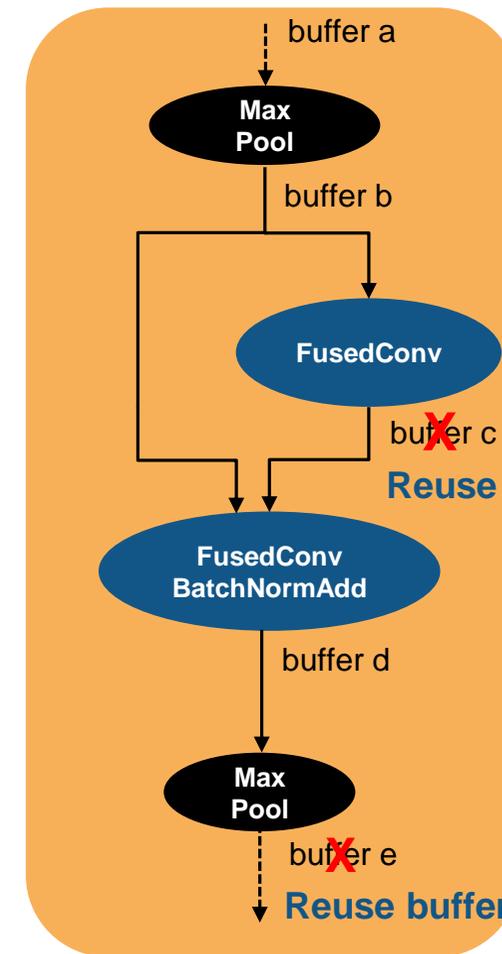
1. 最適化ライブラリの使用
2. ネットワークの最適化
3. コーディングパターン



Network



レイヤーフュージョン



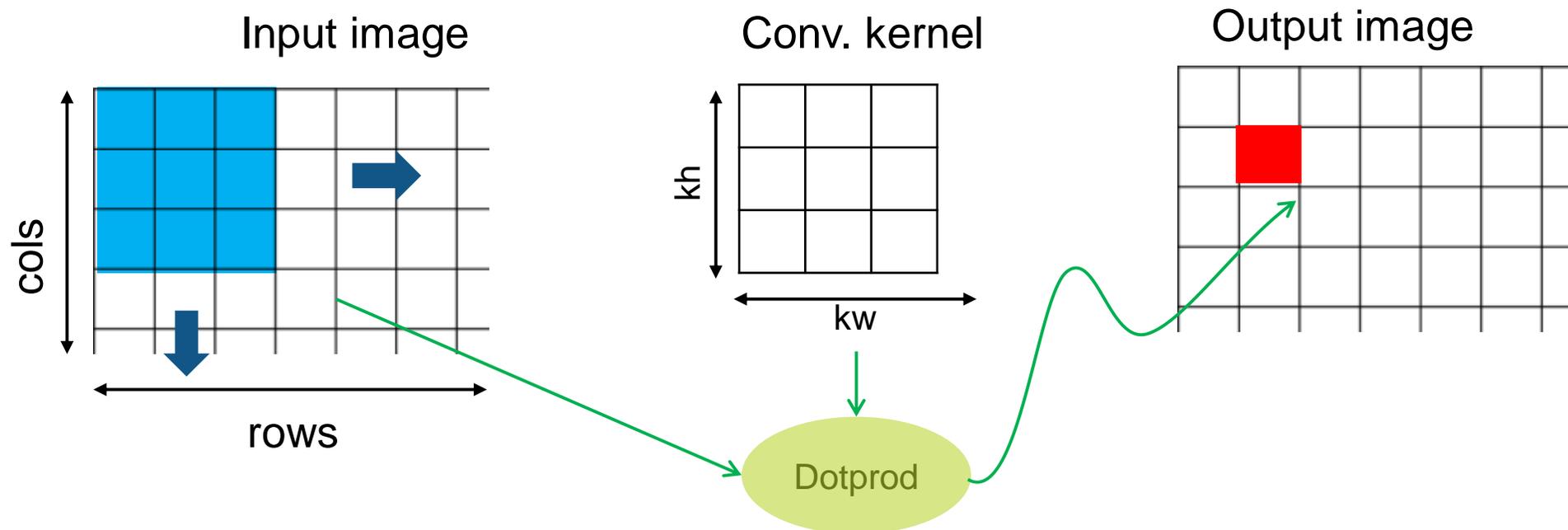
バッファの最小化

コーディングパターン：ステンシルカーネル

Performance

1. 最適化ライブラリの使用
2. ネットワークの最適化
3. **コーディングパターン**

- 画像処理系関数にはGPUによるステンシル演算が自動適用
 - 関数例： `imfilter`, `imerode`, `imdilate`, `conv2`, ...
- マニュアルでの指定には `gpuCoder.stencilKernel()`



コード生成でサポートされるレイヤはリリースごとに増加中

GPU Coder™

MATLAB Coder™

Layer Name	cuDNN	TensorRT	Intel MKL-DNN	ARM Compute
AdditionLayer	✓	✓	✓	✓
AveragePooling2DLayer	✓	✓	✓	✓
BatchNormalizationLayer	✓	✓	✓	✓
ClassificationOutputLayer	✓	✓	✓	✓
ClippedReLULayer	✓		✓	✓ (R2019a)
Convolution2DLayer	✓	✓	✓	✓
R2019a Crop2DLayer	✓ (R2019a)	✓ (R2019a)	✓ (R2019a)	
CrossChannelNormalizationLayer	✓	✓	✓	✓
DepthConcatenationLayer	✓	✓ (R2019a)	✓	✓
DropoutLayer	✓	✓	✓	✓
FullyConnectedLayer	✓	✓	✓	✓
ImageInputLayer	✓	✓	✓	✓
LeakyReLULayer	✓	✓	✓	✓ (R2019a)
MaxPooling2DLayer	✓	✓	✓	✓
MaxUnpooling2DLayer	✓	✓	✓ (R2019a)	
PixelClassificationLayer	✓	✓	✓	✓
ReLULayer	✓	✓	✓	✓
RegressionOutputLayer	✓	✓	✓	✓
SoftmaxLayer	✓	✓	✓	✓
TransposedConvolution2DLayer	✓	✓	✓	
R2019a YOLOv2OutputLayer	✓ (R2019a)	✓ (R2019a)	✓ (R2019a)	✓ (R2019a)
R2019a YOLOv2ReorgLayer	✓ (R2019a)	✓ (R2019a)		
R2019a YOLOv2TransformLayer	✓ (R2019a)	✓ (R2019a)		
(Keras Layer)FlattenCStyleLayer	✓	✓		
(Keras Layer)GlobalAveragePooling2dLayer	✓	✓	✓	✓
(Keras Layer)SigmoidLayer	✓	✓		
(Keras Layer)TanhLayer	✓	✓		
(Keras Layer)ZeroPadding2dLayer	✓	✓		

R2019a

✓: R2019aからコード生成可能なレイヤー

サポートされるネットワーク(GPU編)

CuDNN

Image
Classification



VGG, ResNet,
SqueezeNet, GoogLeNet

Encoder/Decoder
Semantic
Segmentation



DenoiseNet, SegNet
FCN (R2019a)

Object detection
ONNX/Keras imported
networks



YOLOv2 (R2019a)
Keras (R2019a)

TensorRT

VGG, ResNet
SqueezeNet, GoogLeNet
(R2019a)

DenoiseNet, SegNet
FCN (R2019a)

YOLOv2 (R2019a)
Keras (R2019a)

サポートされるネットワーク(CPU編)

Image
Classification



Encoder/Decoder
Semantic
Segmentation



Intel MKLDNN

VGG, ResNet,
SqueezeNet, GoogLeNet

DenoiseNet
Unet, SegNet
(R2019a)

ARM CPUs

VGG, ResNet
SqueezeNet, GoogLeNet

DenoiseNet

アジェンダ

- MATLAB Coder/GPU Coderの概要
- ディープニューラルネットワークの組み込み実装ワークフロー
- パフォーマンスに関して
- まとめ

まとめ

ネットワークの構築と学習

Keras
ONNX
Caffe



モデル
インポート



MATLABで学習



学習済み
モデル

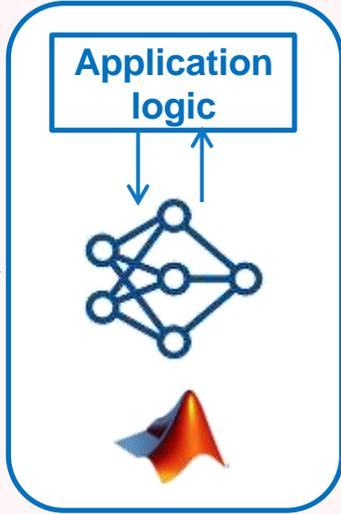
事前学習済み
モデル



転移学習



アプリケー ション化



組み込み GPU/CPU実装

Coders

NVIDIA®

TensorRT and
cuDNN Libraries

Intel®

MKL-DNN
Library

ARM®

ARM Compute
Library



© 2019 The MathWorks, Inc. MATLAB and Simulink are registered trademarks of The MathWorks, Inc. See www.mathworks.com/trademarks for a list of additional trademarks. Other product or brand names may be trademarks or registered trademarks of their respective holders.