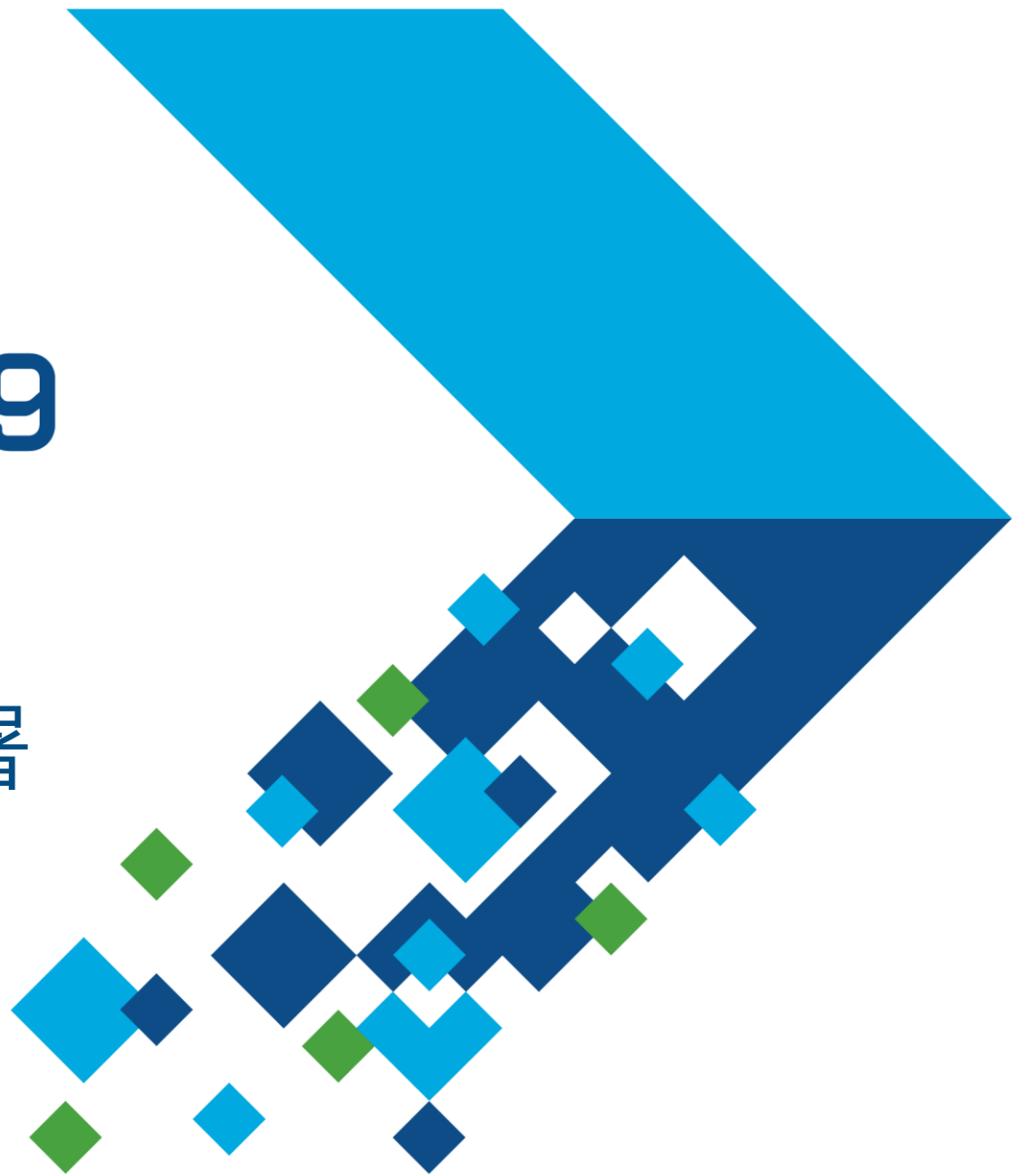


MATLAB EXPO 2019

在嵌入式 GPU 和 CPU 上部署
深度神经网络

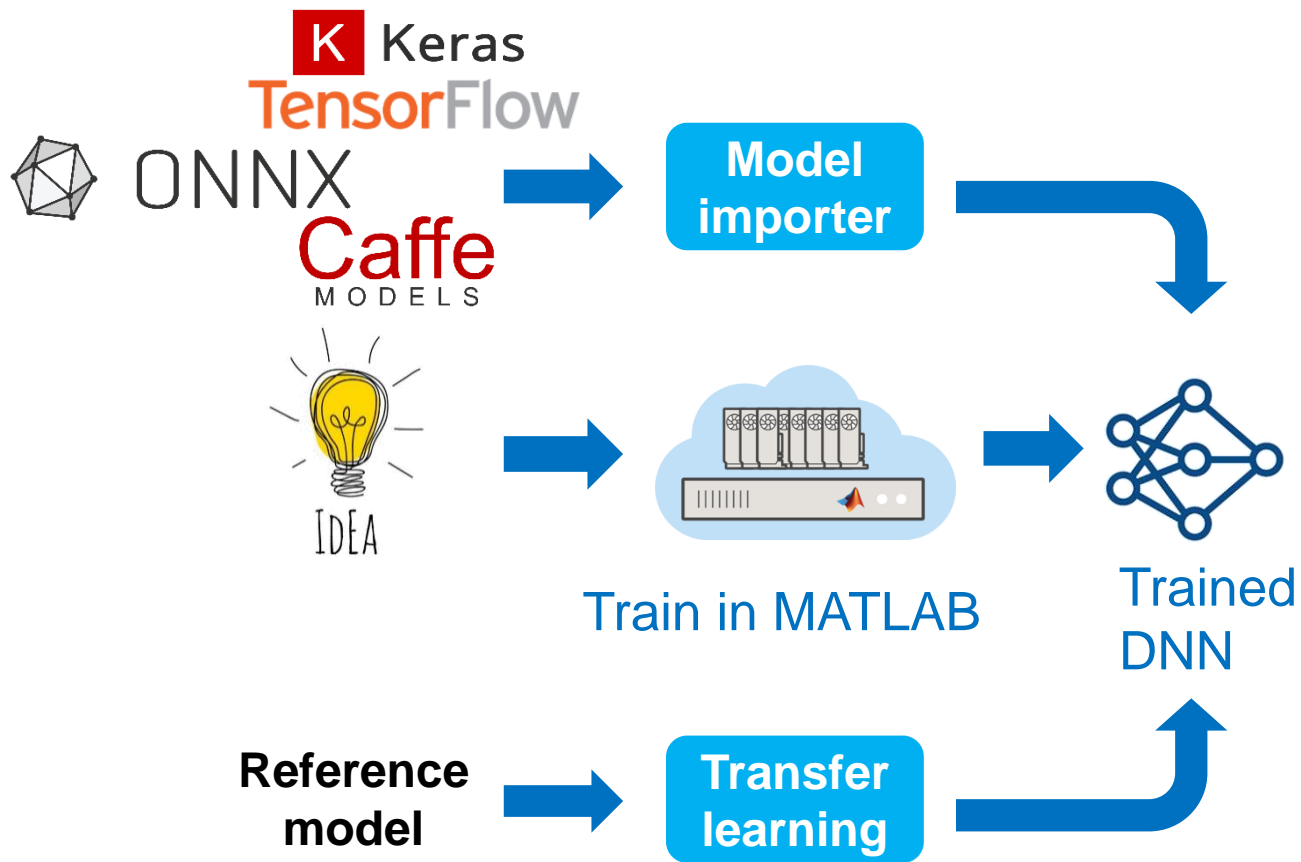
刘海伟



MATLAB 深度学习流程

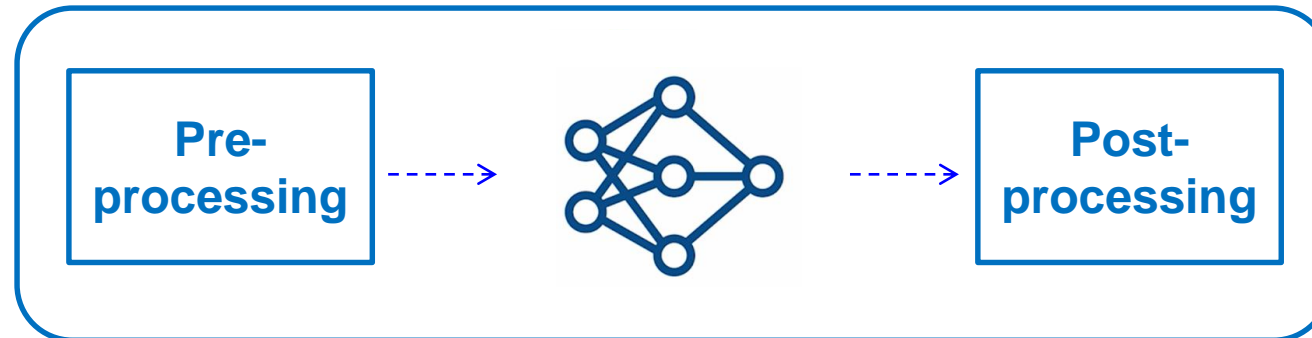
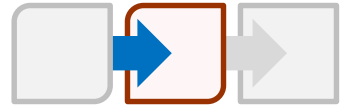


深度神经网络设计与训练

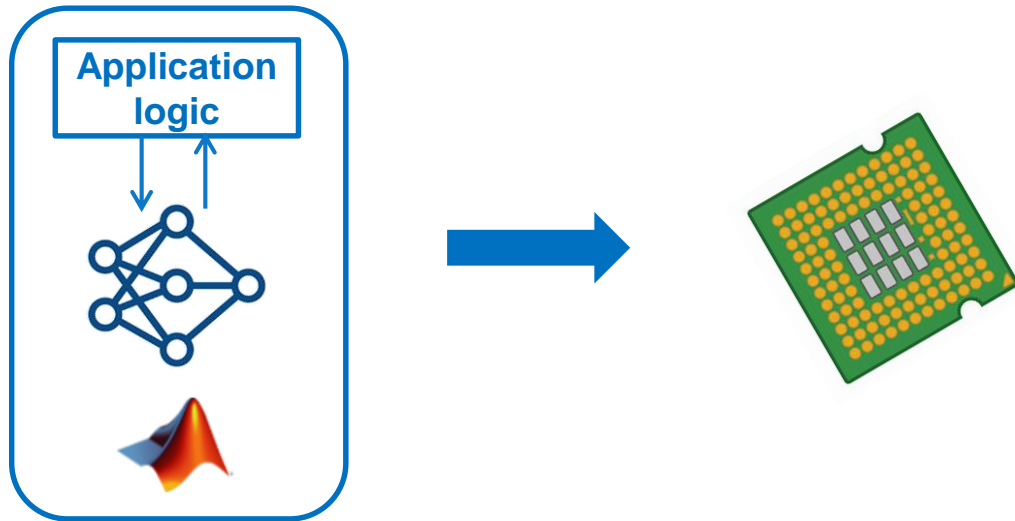


- 在 **MATLAB** 中设计
 - **管理**大数据集
 - **自动**数据标签
 - 模型**易于访问**
- 在 **MATLAB** 中训练
 - 使用 GPU **加速**
 - **扩展**到集群

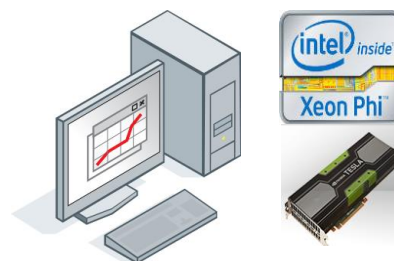
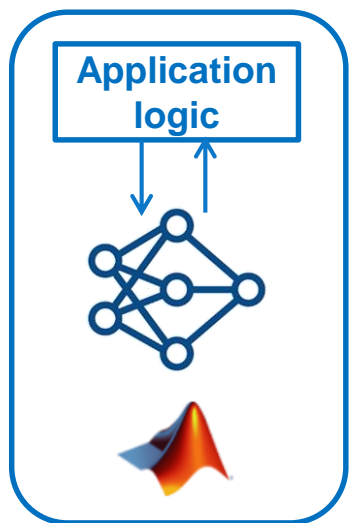
应用设计



深度学习多平台部署



深度学习多平台部署



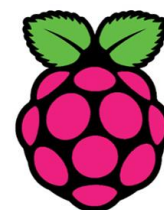
Desktop



Data Center



NVIDIA Jetson



Raspberry pi



Mobile



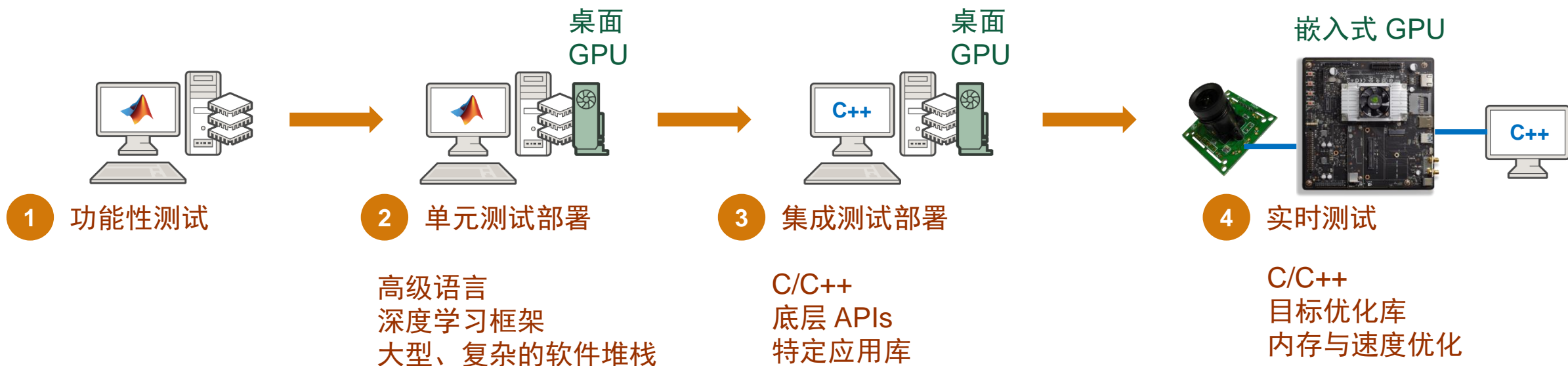
Beaglebone



嵌入式

从算法设计到嵌入式部署流程

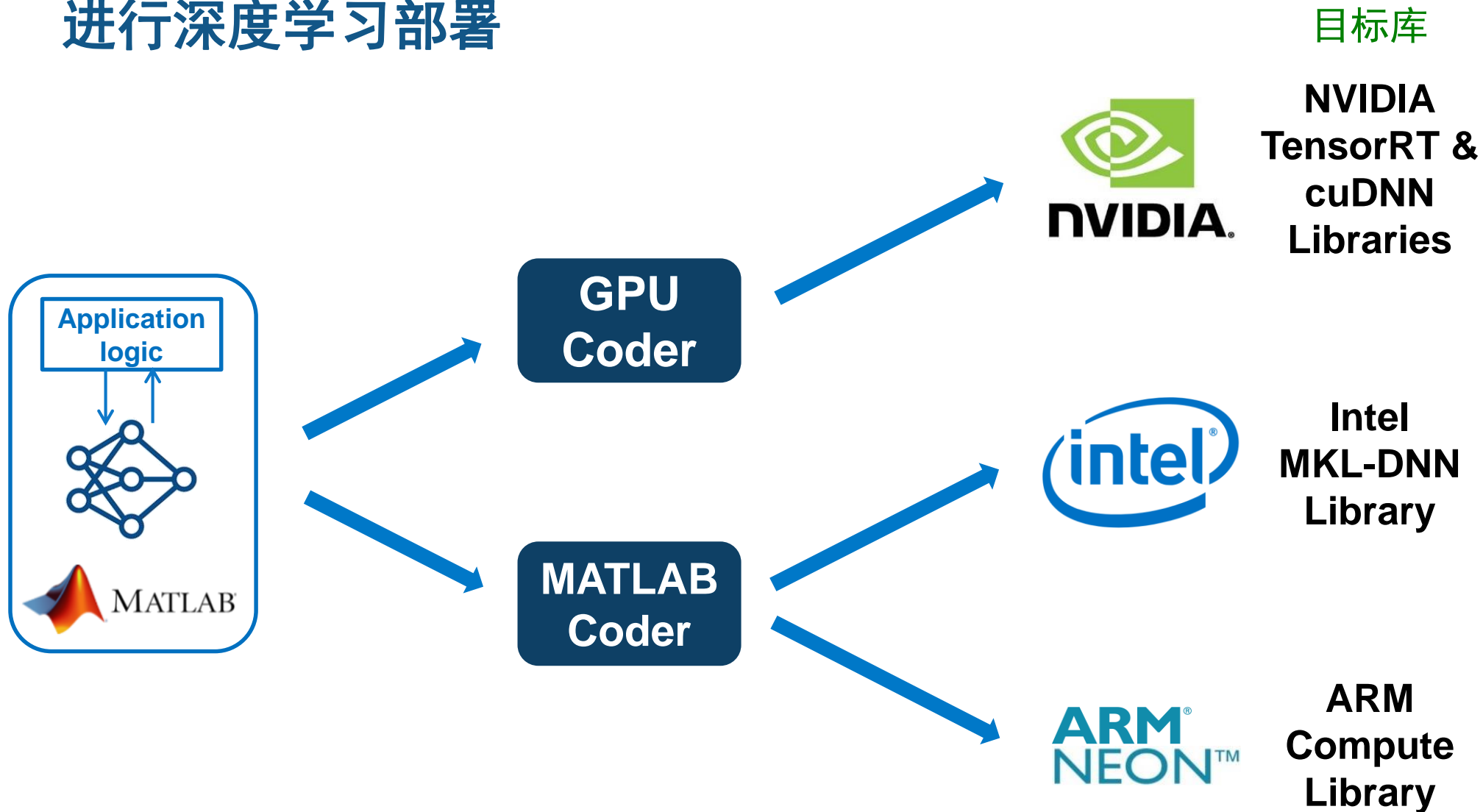
传统方法



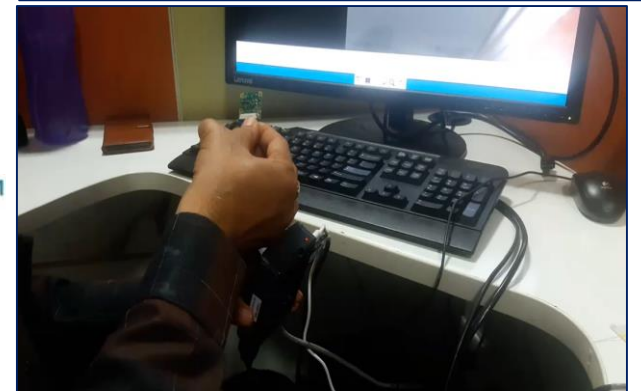
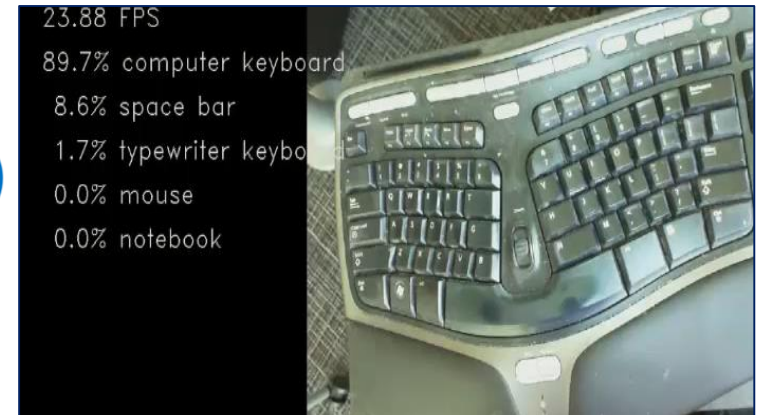
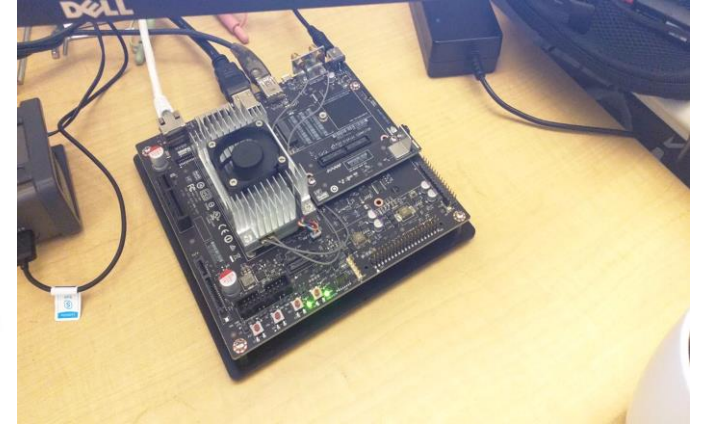
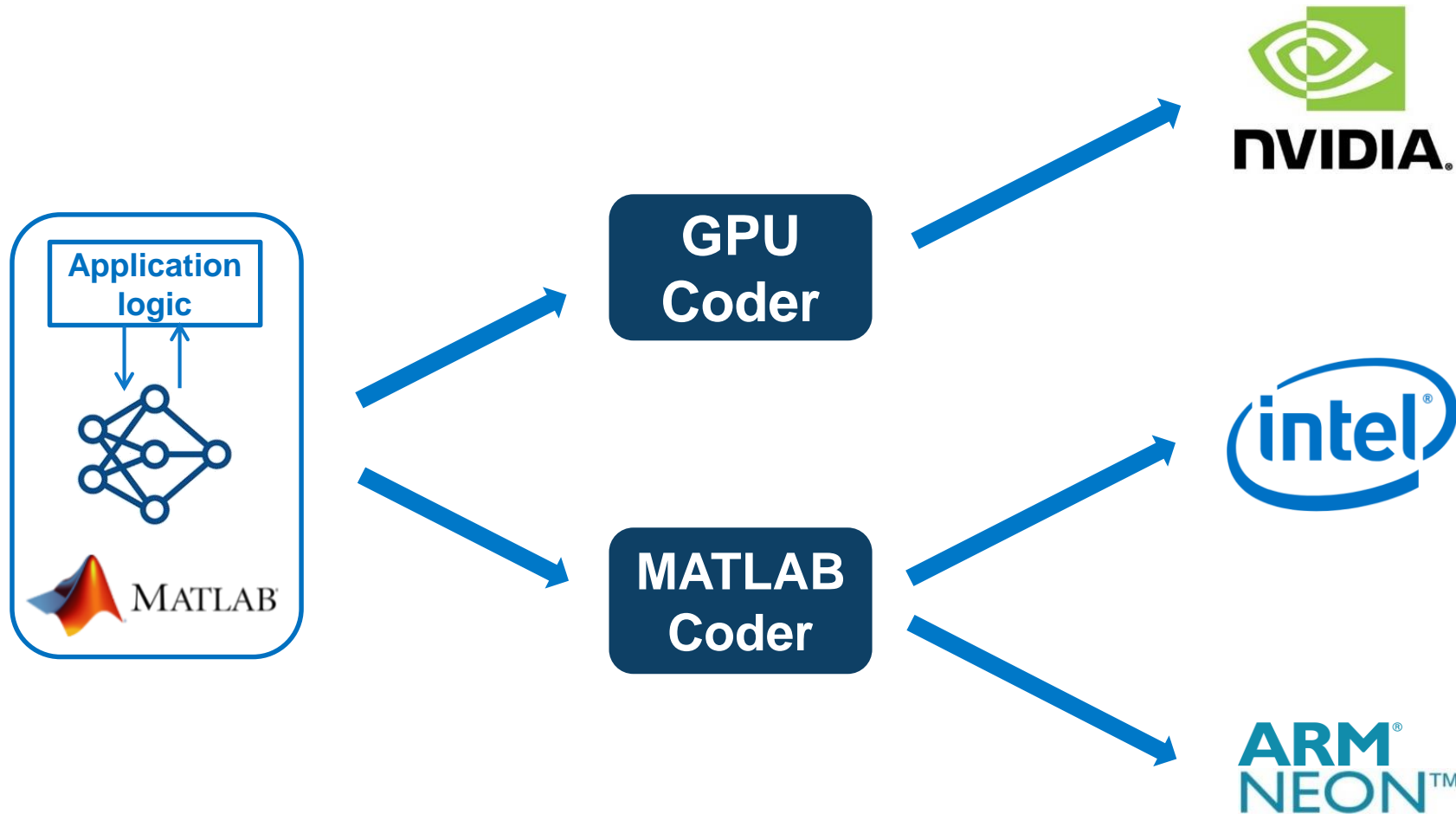
挑战

- 集成多个库和包
- 验证和维护多个实现
- 算法和供应商锁定

解决方案：使用 MATLAB Coder 和 GPU Coder 进行深度学习部署

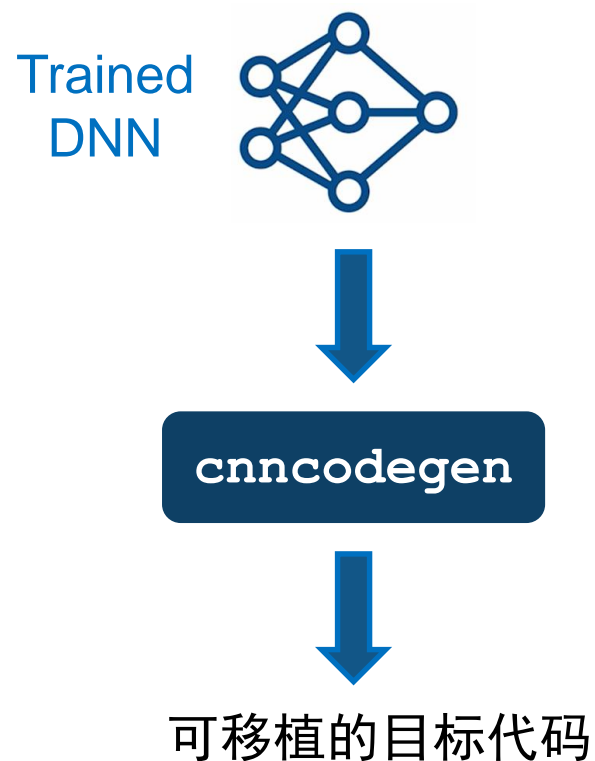


解决方案：使用 MATLAB Coder 和 GPU Coder 进行深度学习部署

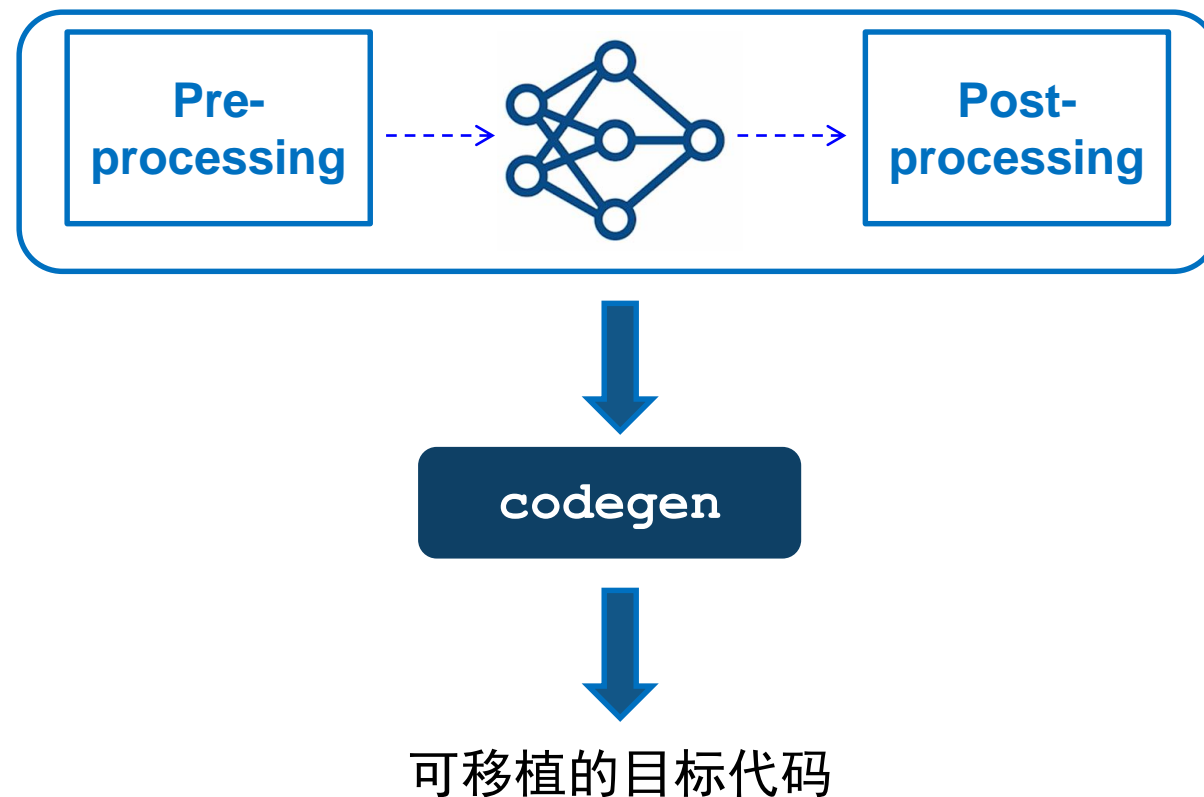


深入学习部署流程

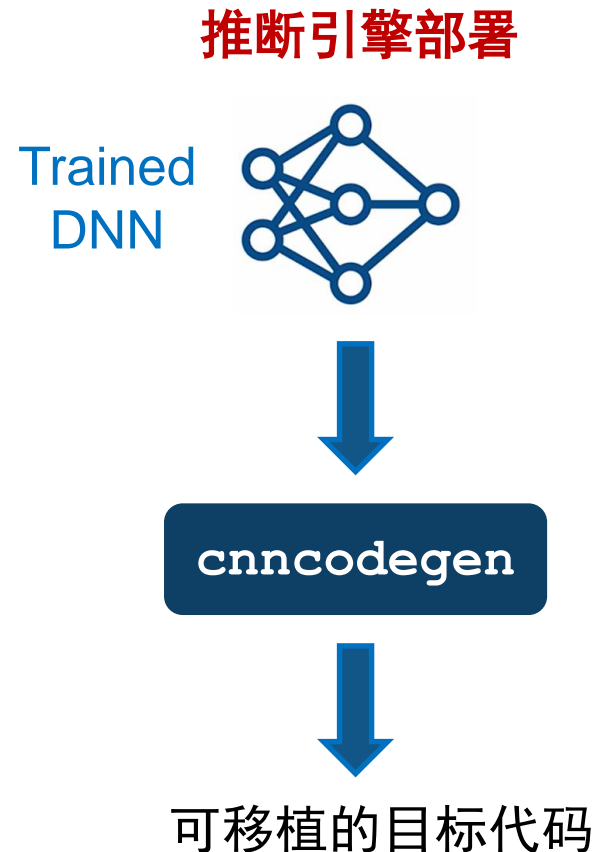
推断引擎部署



应用程序集成部署



推断引擎部署流程



推断引擎部署的步骤

1. 为训练好的模型生成代码

```
>> cnncodegen(net, 'targetlib', 'arm-compute')
```

2. 将生成的代码复制到目标器件上

3. 生成推断引擎

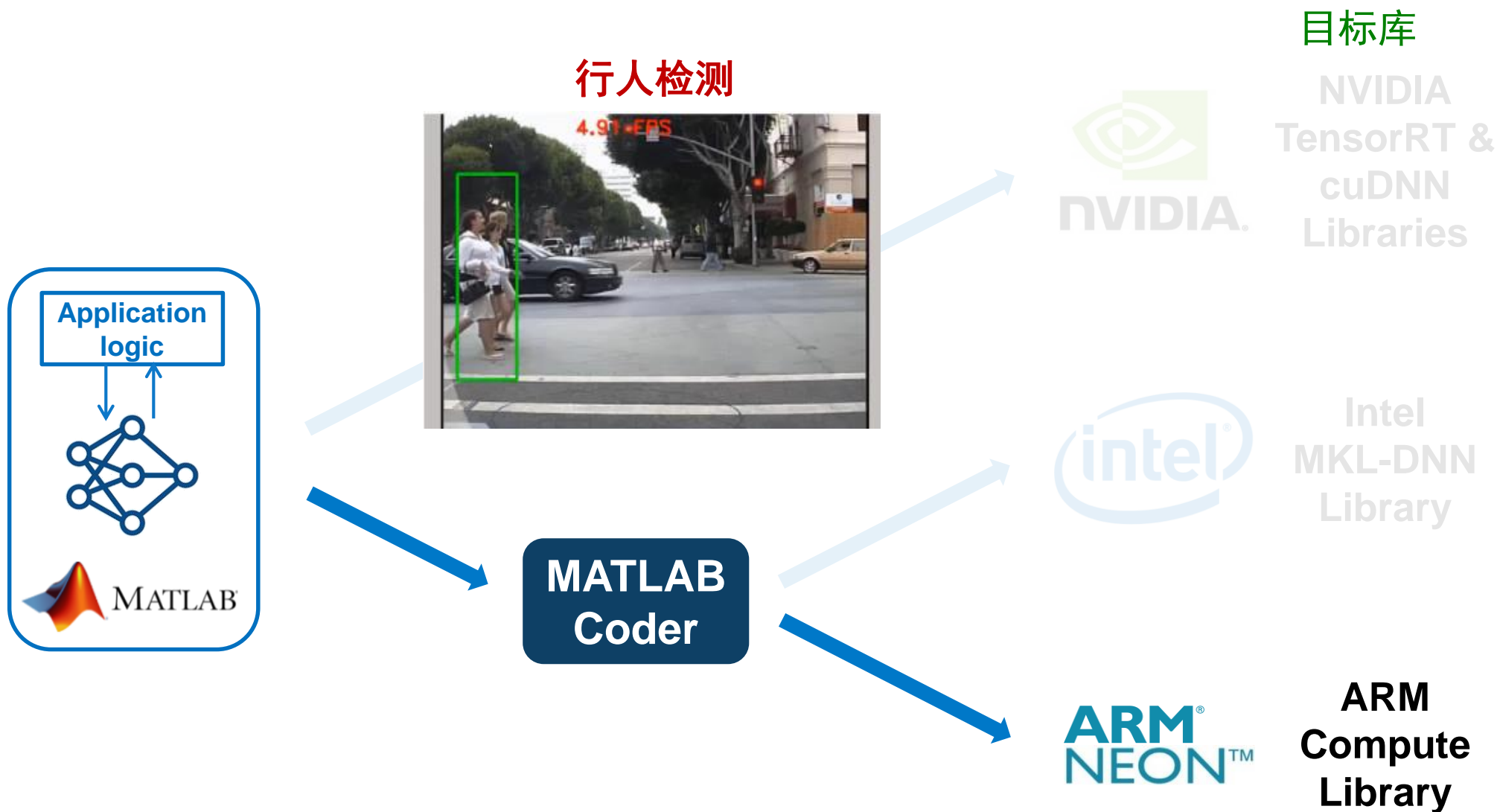
```
>> make -C ./codegen -f ...mk
```

4. 手写主函数，用以调用推断引擎

5. 生成并测试可执行 exe

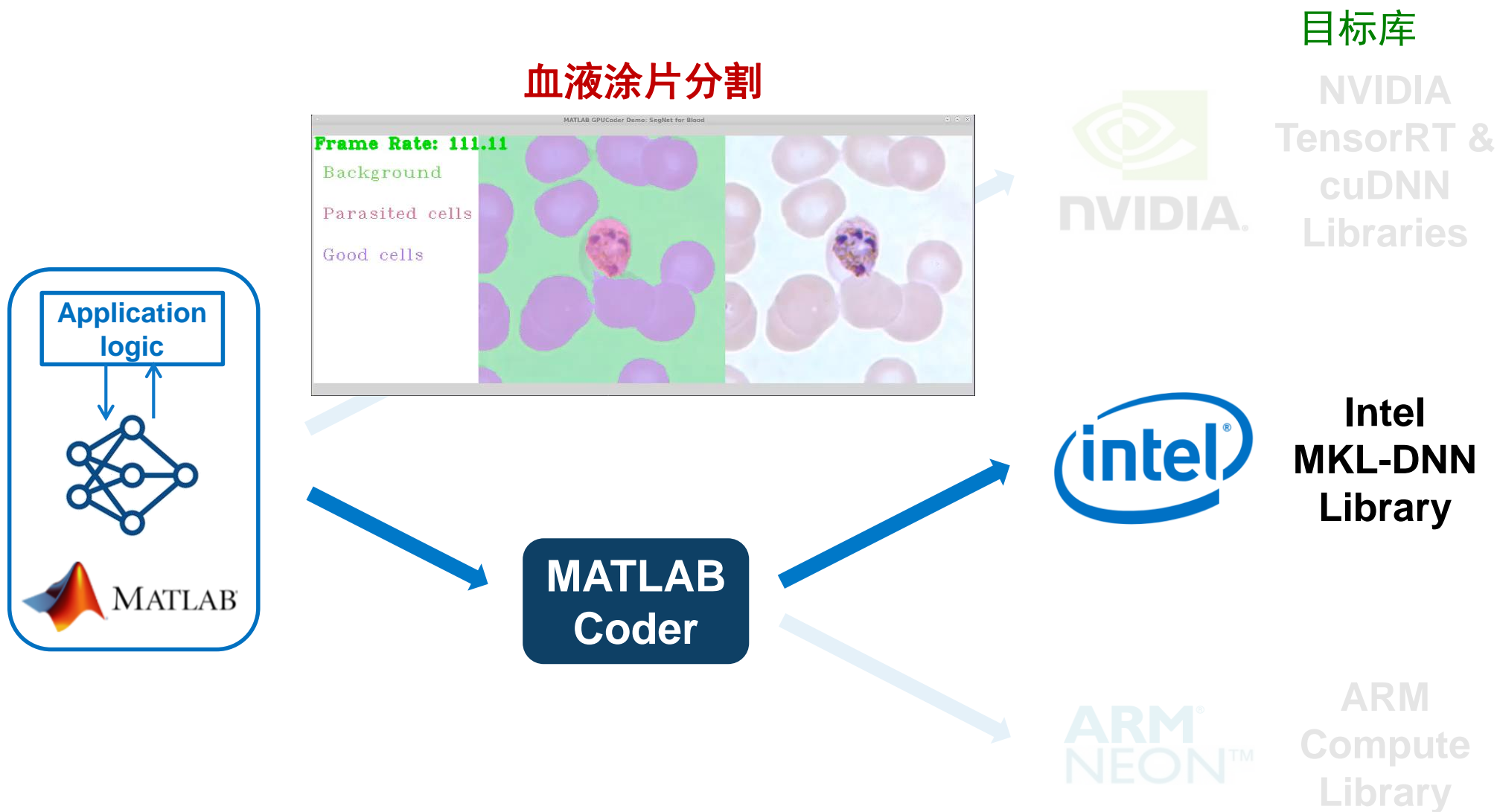
```
>> make -C ./ .....
```

深度学习推断部署

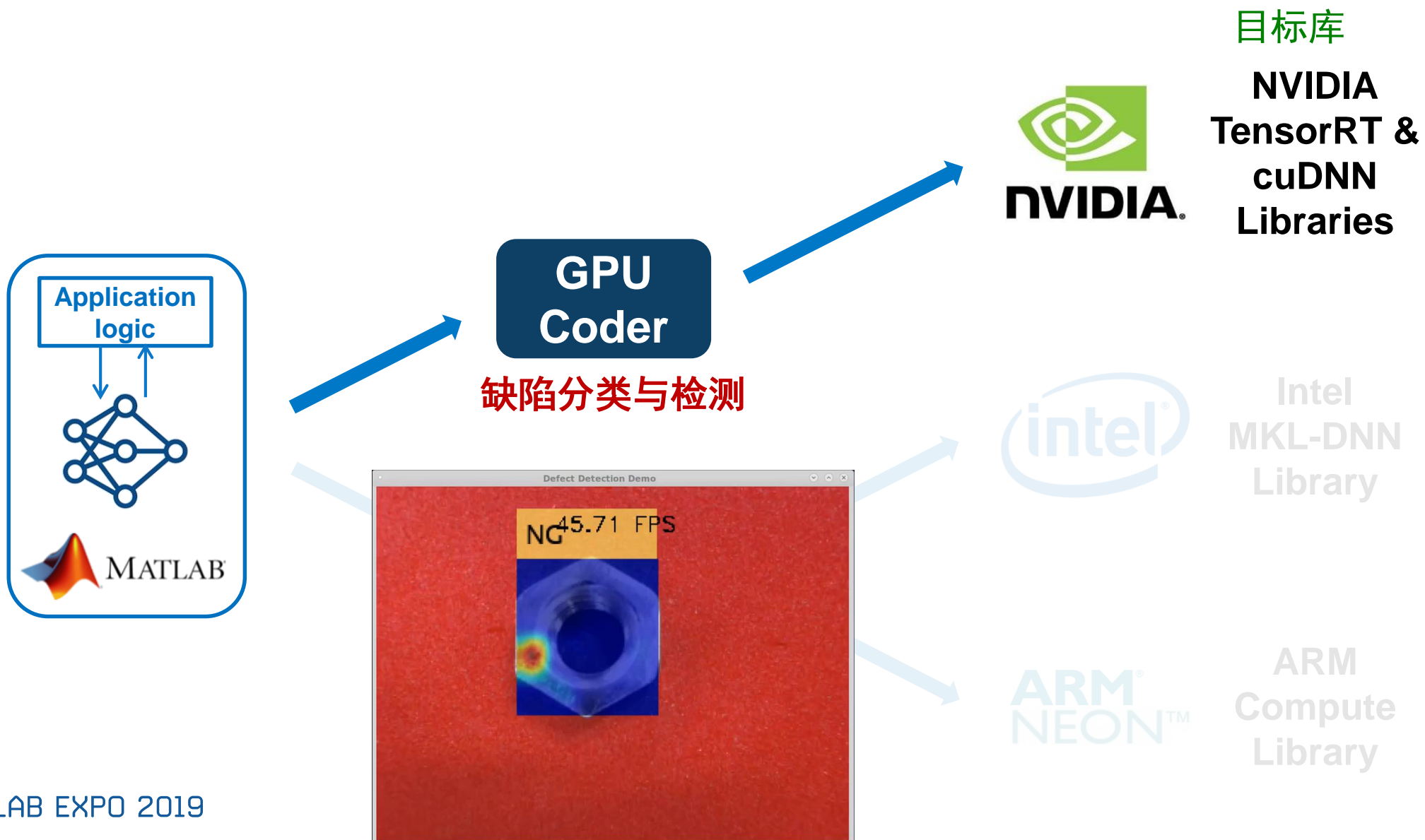


深度学习推断部署

血液涂片分割



深度学习推断部署

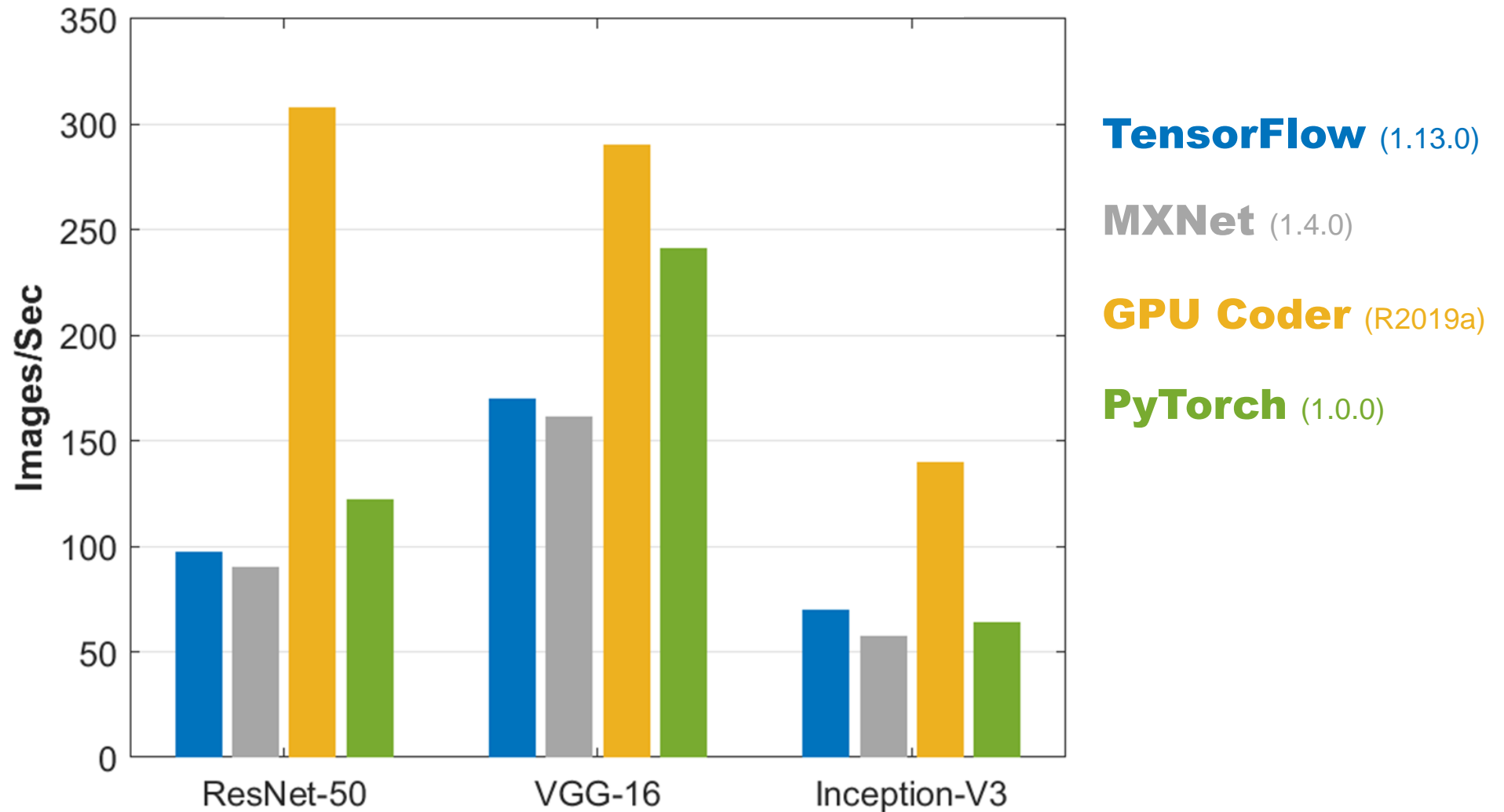




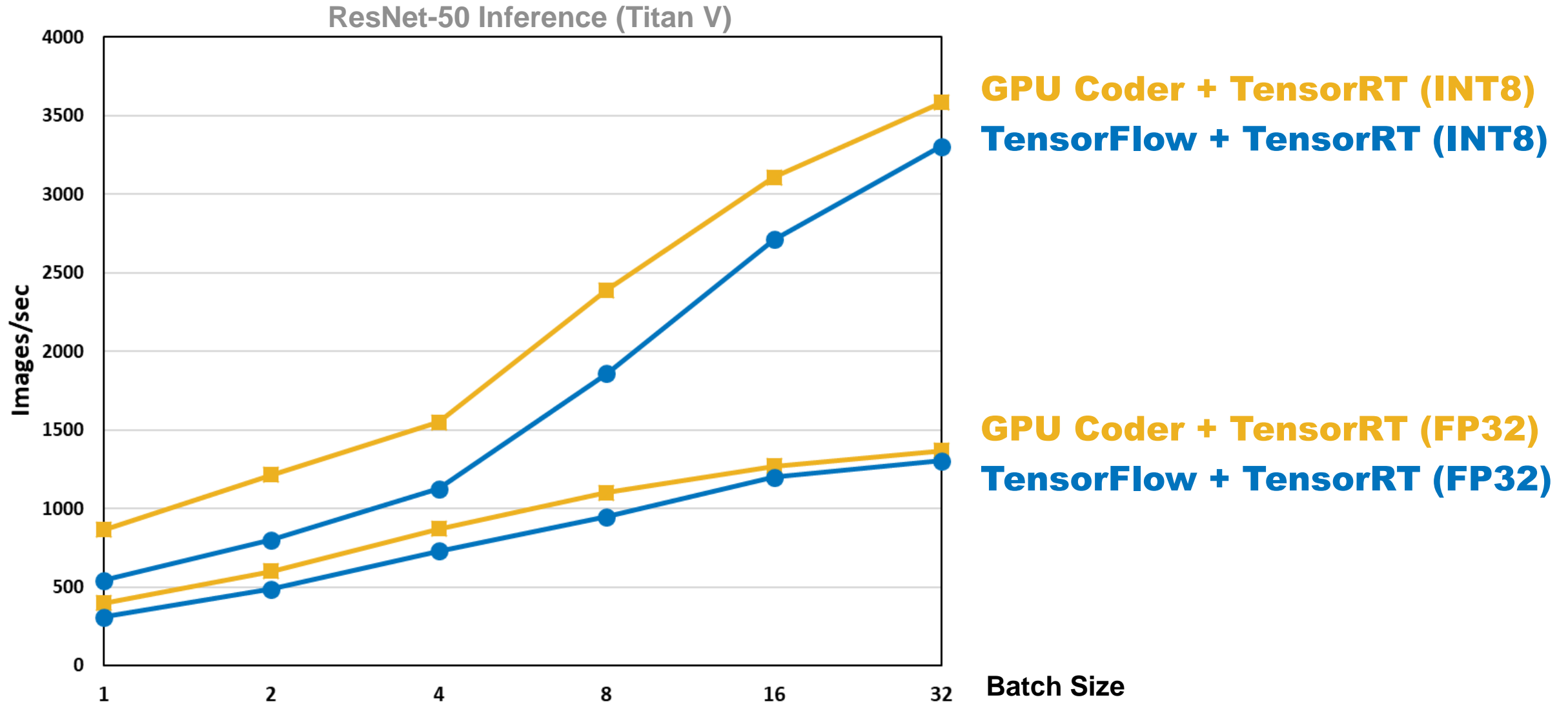
生成代码的性能

- Titan V GPU 上进行 CNN 推断(ResNet-50, VGG-16, Inception V3)
- Jetson TX2 上进行 CNN 推断(ResNet-50)
- Intel Xeon CPU 上进行 CNN推断(ResNet-50 , VGG-16, Inception V3)

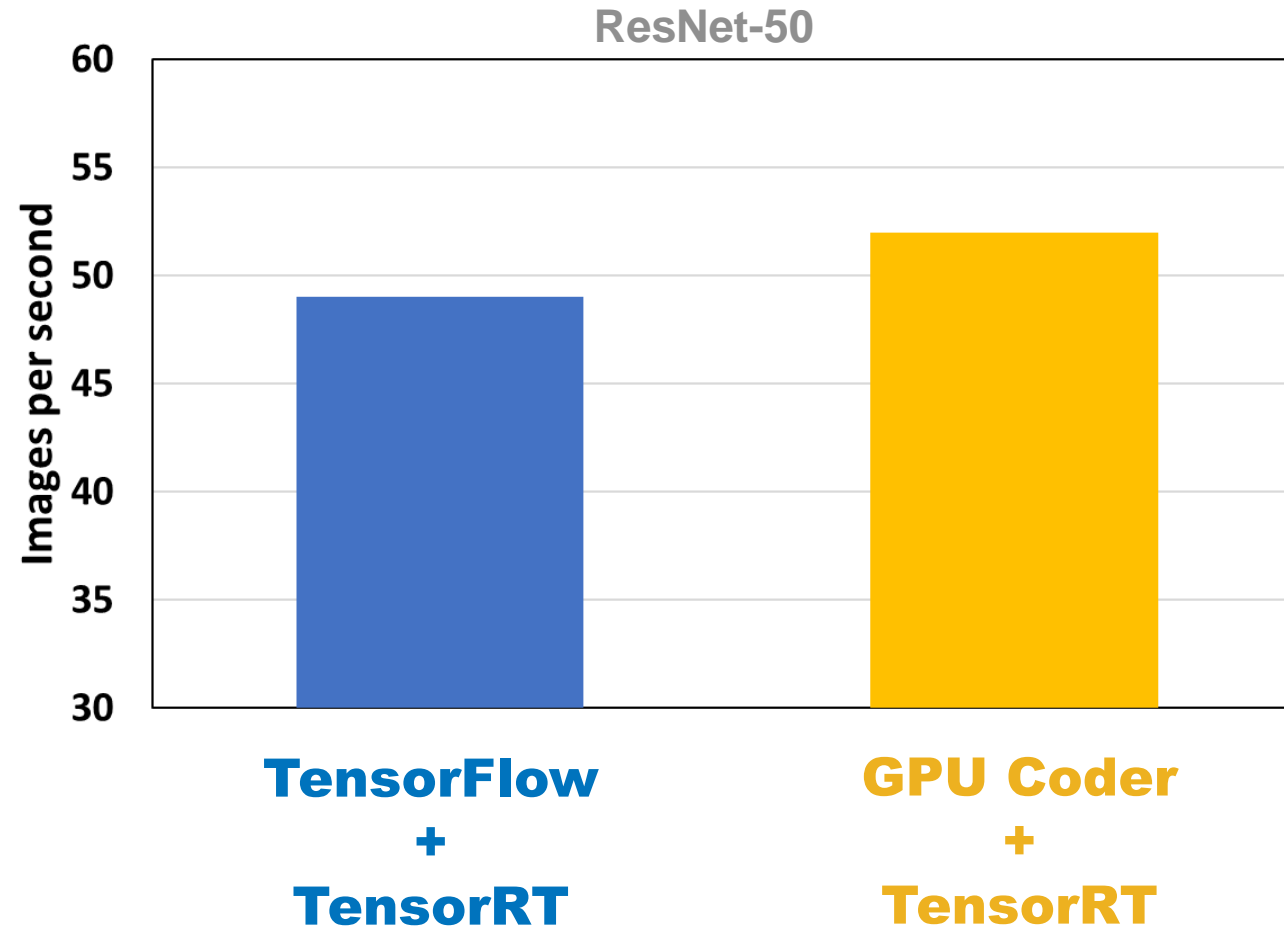
使用 Titan V cuDNN 进行单图像推断



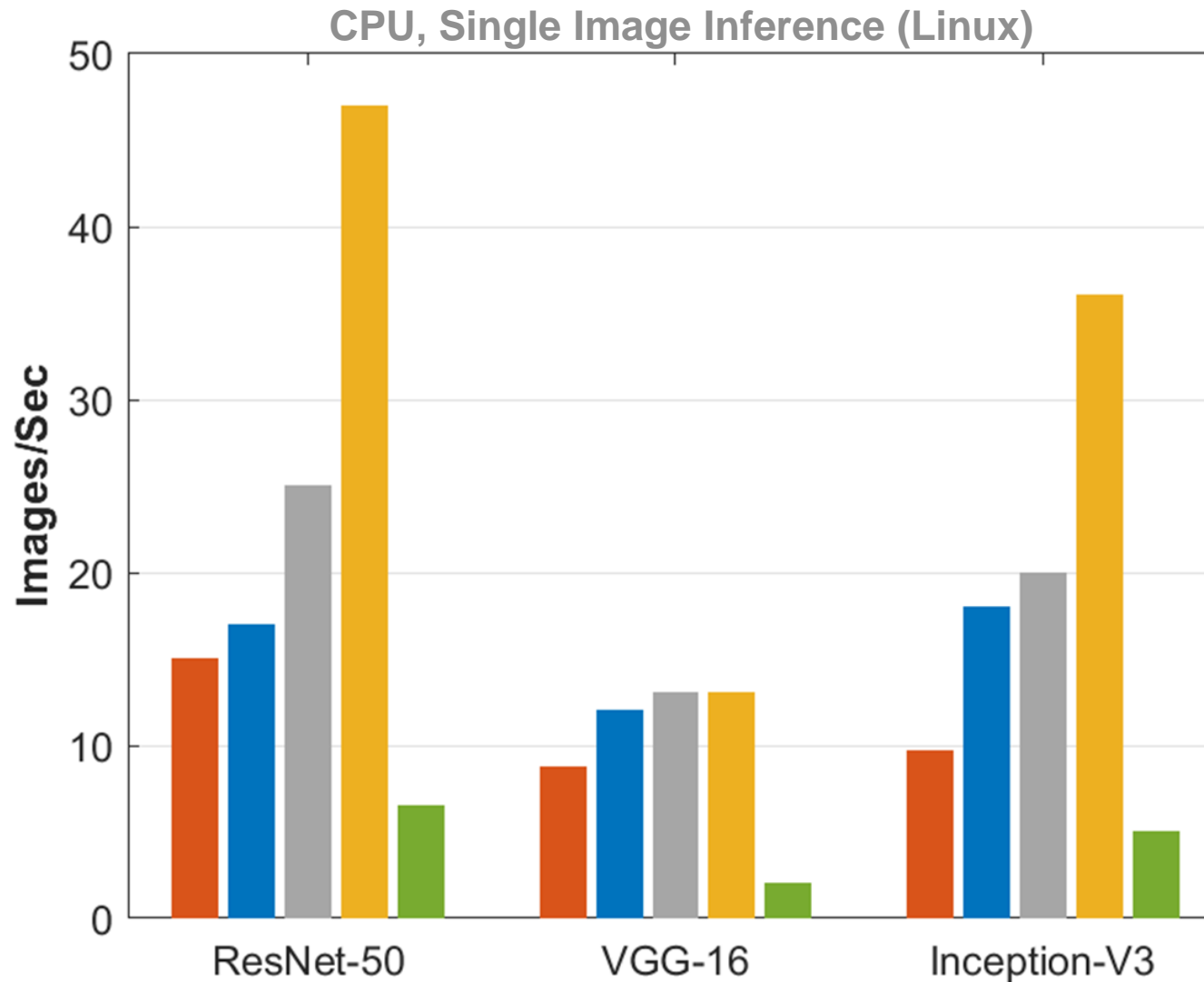
使用 TensorRT INT8 更强大的性能



Jetson TX2上的单图像推理



CPU 性能



MATLAB

TensorFlow

MXNet

MATLAB Coder

PyTorch

简要总结

DNN 库非常适合推断.....

MATLAB Coder 和 GPU Coder 生成的代码可以利用：



NVIDIA® CUDA 库，包括 TensorRT & cuDNN



用于深度神经网络的 Intel® 数学内核库 (MKL-DNN)



ARM® 移动平台计算库

简要总结

DNN 库非常适合推断.....

MATLAB Coder 和 GPU Coder 生成的代码可以利用：

**但是，应用程序需要的
不仅仅是推断**



NVIDIA® CUDA 库，包括 TensorRT & cuDNN

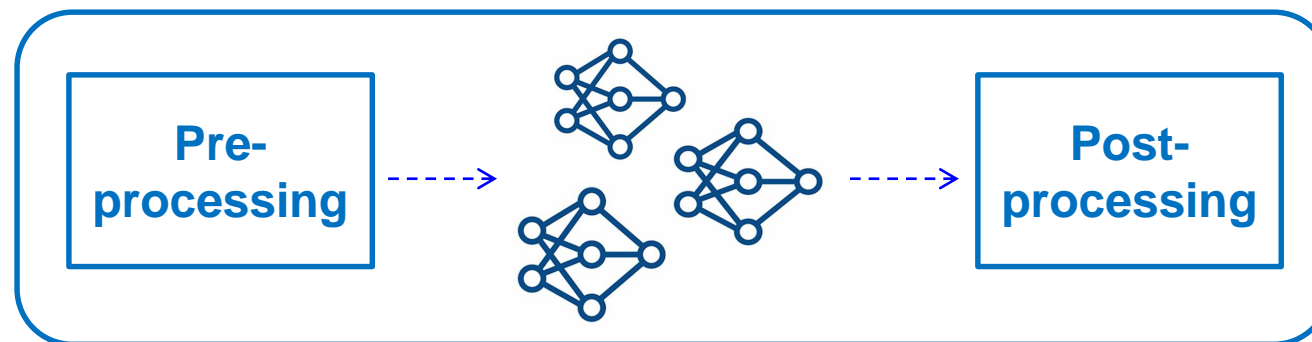


用于深度神经网络的 Intel® 数学内核库 (MKL-DNN)



ARM® 移动平台计算库

深度学习流程：应用程序集成部署



可移植的目标代码



使用 YOLO v2 进行车道线和目标检测



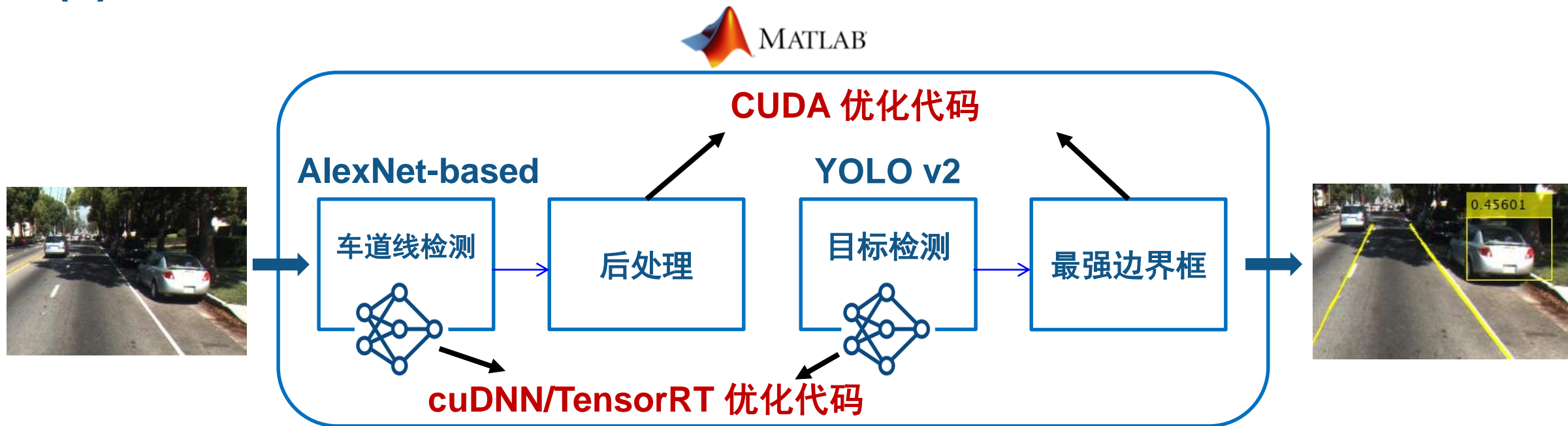
流程:

- 1) 在 CPU 上的 MATLAB 进行测试
- 2) 在桌面 GPU 上生成代码和测试
- 3) 在 Jetson AGX Xavier GPU 上生成代码和测试

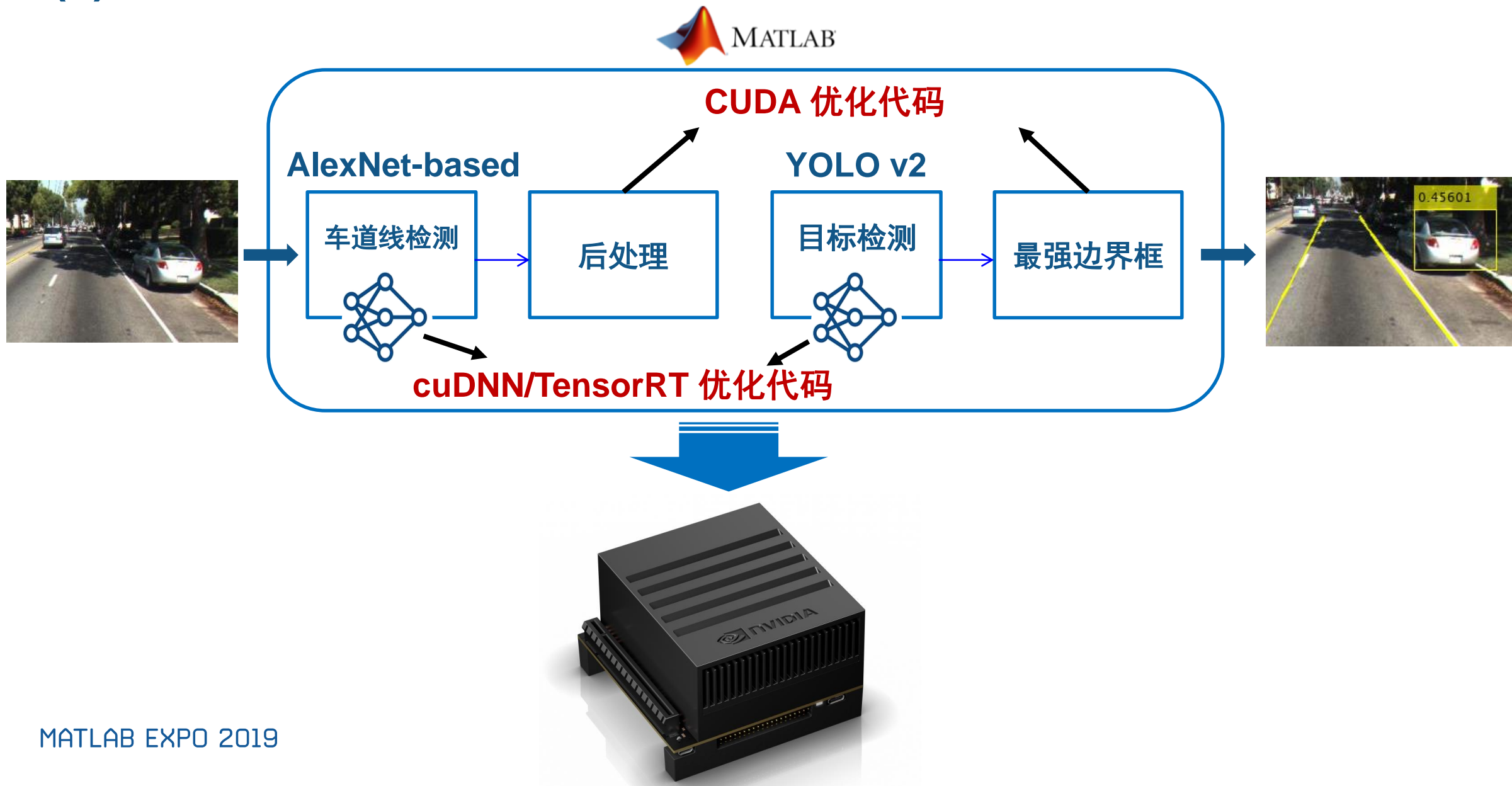
(1) 在 CPU 上用 MATLAB 进行测试



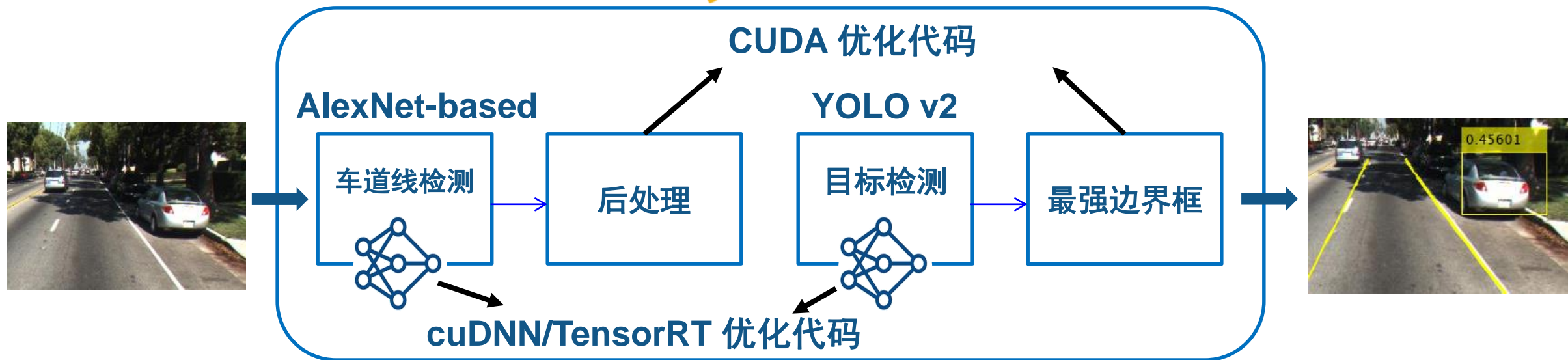
(2) 在桌面 GPU 上生成代码并测试



(3) 在Jetson AGX Xavier GPU 上生成代码并测试

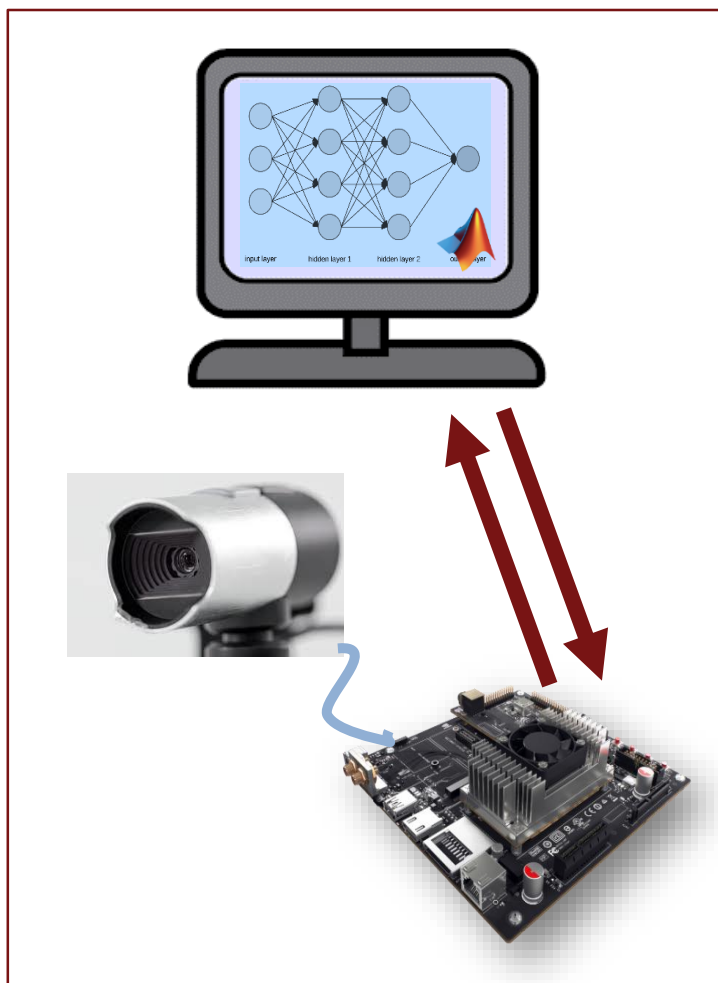


使用YOLO v2进行车道线和目标检测



- 1) CPU 上运行
- 2) 在桌面 GPU 上生成代码，运行速度提高7倍
- 3) 在Jetson AGX Xavier GPU 上生成代码和测试

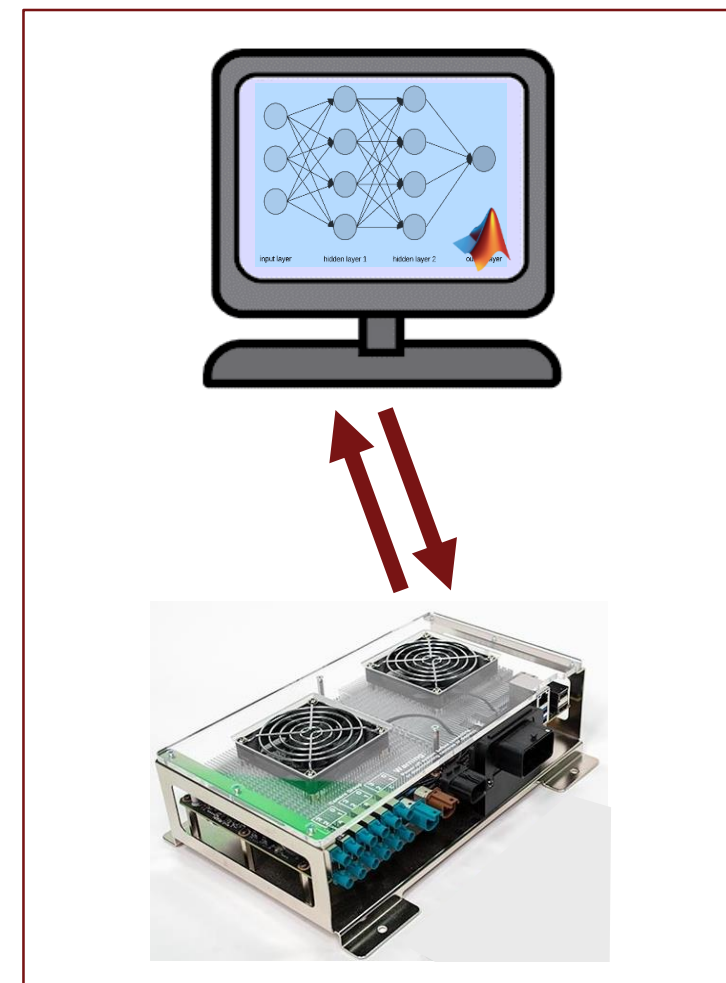
访问硬件



从 MATLAB 访问外设

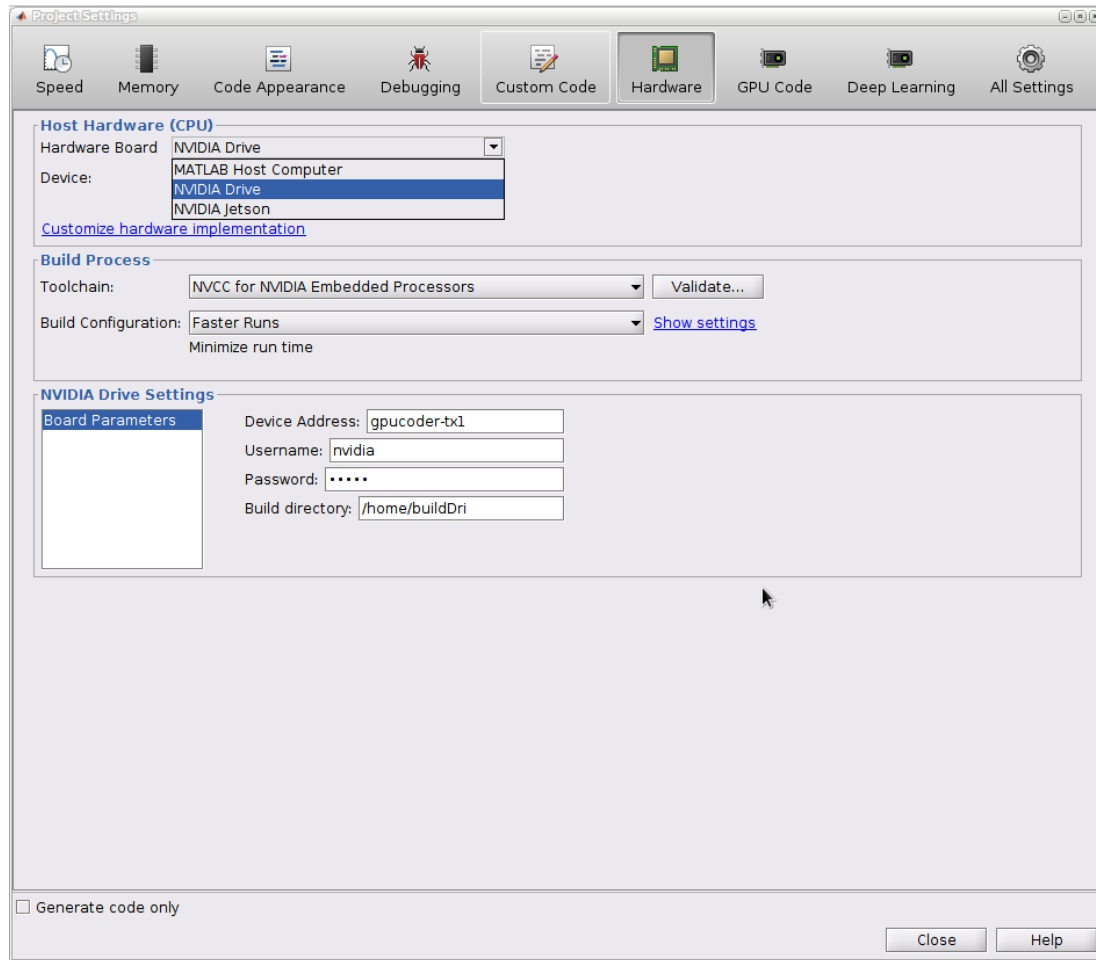


部署独立应用程序



处理器在环验证

通过 Apps 和命令行实现目标硬件部署



```
%% 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;
```

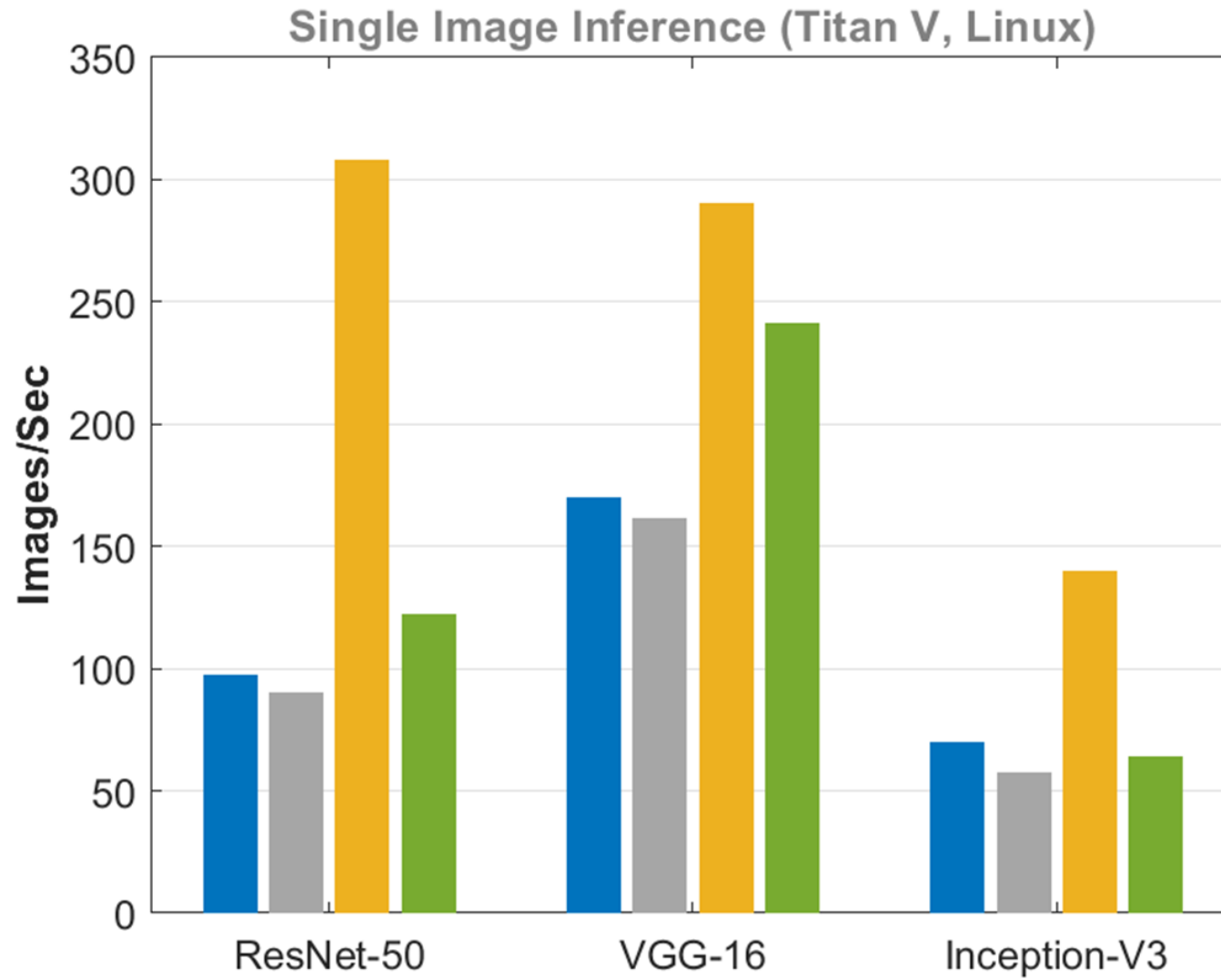
```
%% add user written main files for building executable
% and generate/build the code.
cfg_hsp.CustomSource = 'driver_files_alexnet/main.cu';
cfg_hsp.CustomInclude = 'driver_files_alexnet/';
```

```
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']);
```

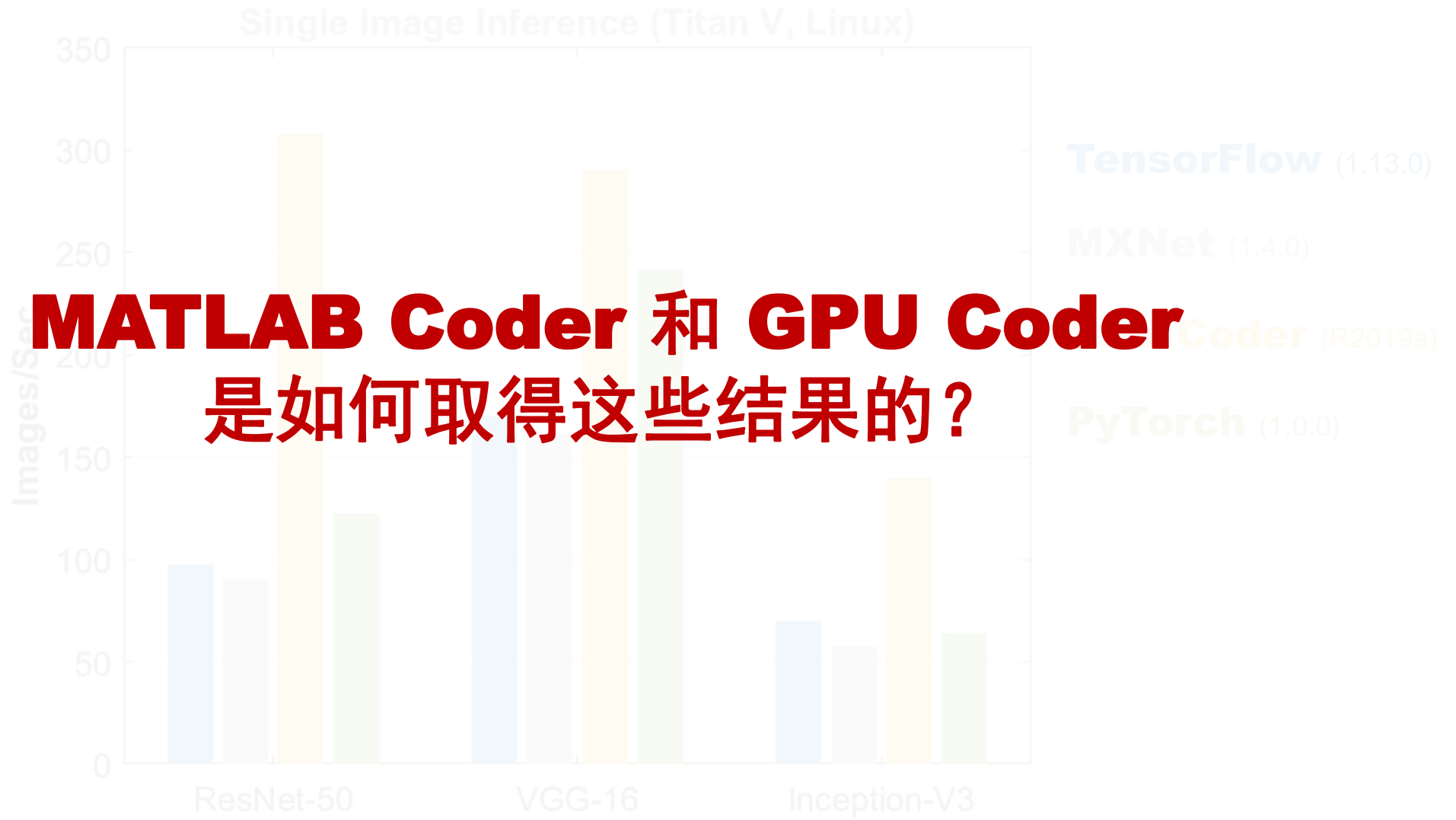


TensorFlow (1.13.0)

MXNet (1.4.0)

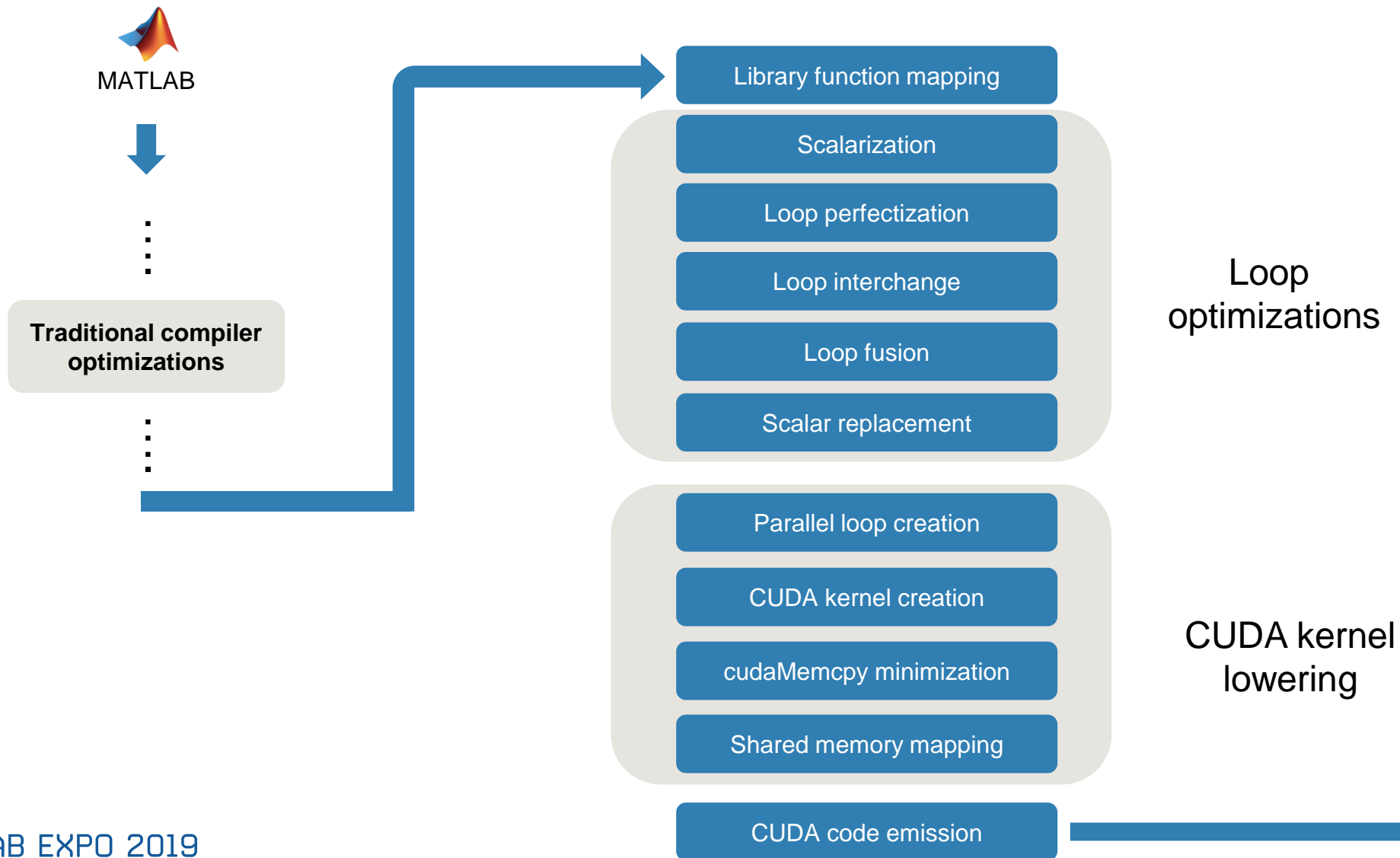
GPU Coder (R2019a)

PyTorch (1.0.0)

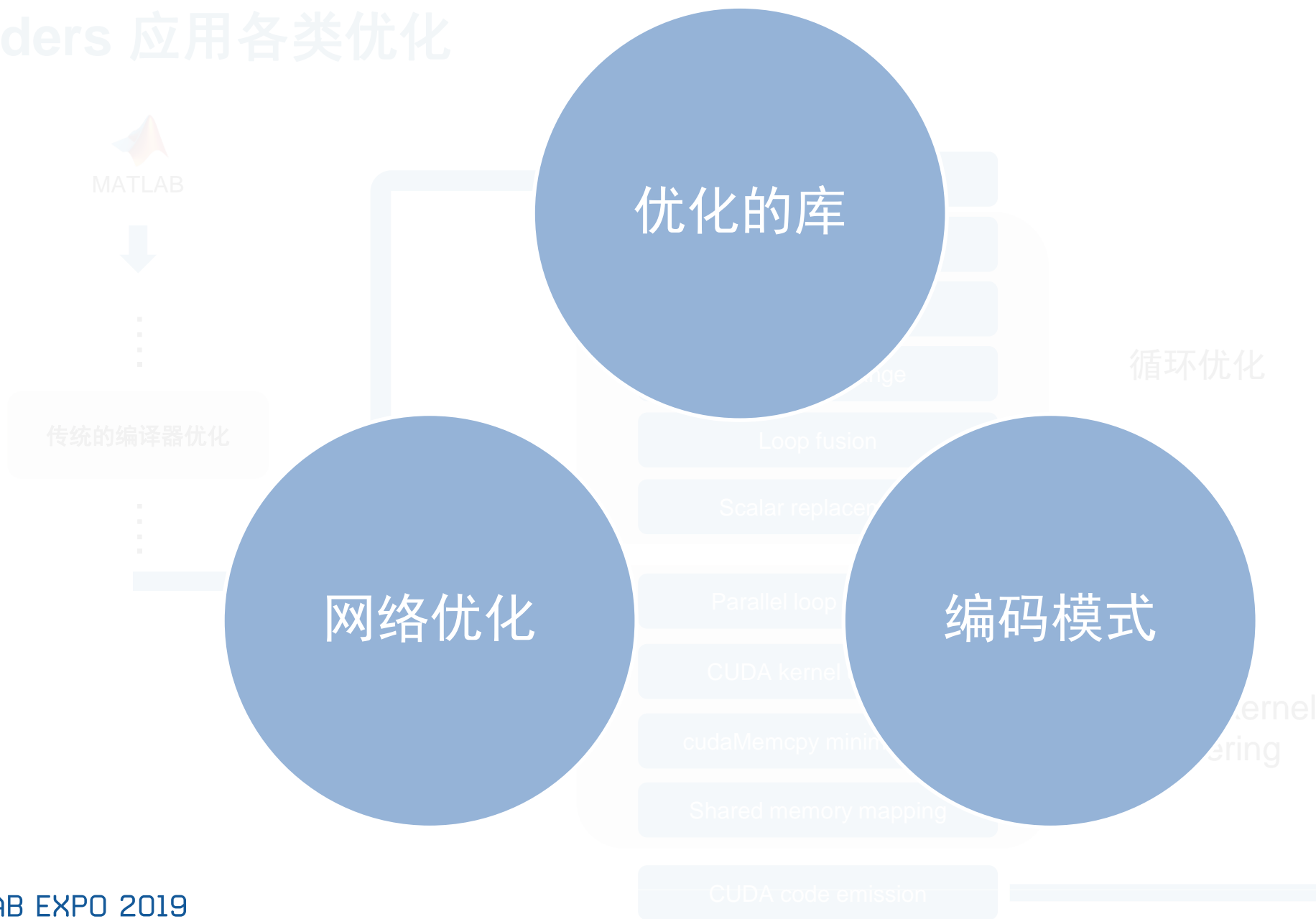


MATLAB Coder 和 GPU Coder
是如何取得这些结果的？

Coders 应用各类优化



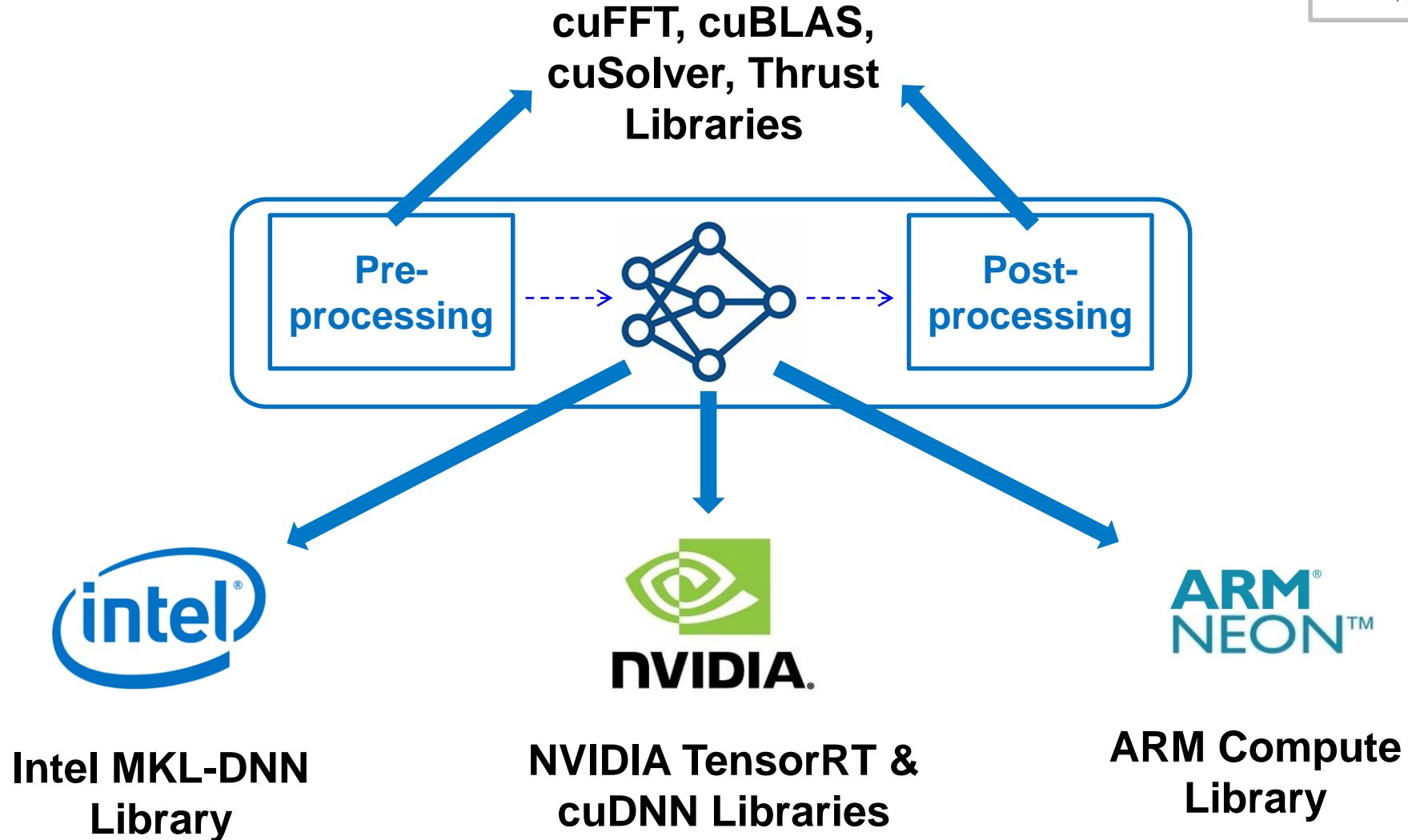
Coders 应用各类优化



生成的代码调用优化的库

性能

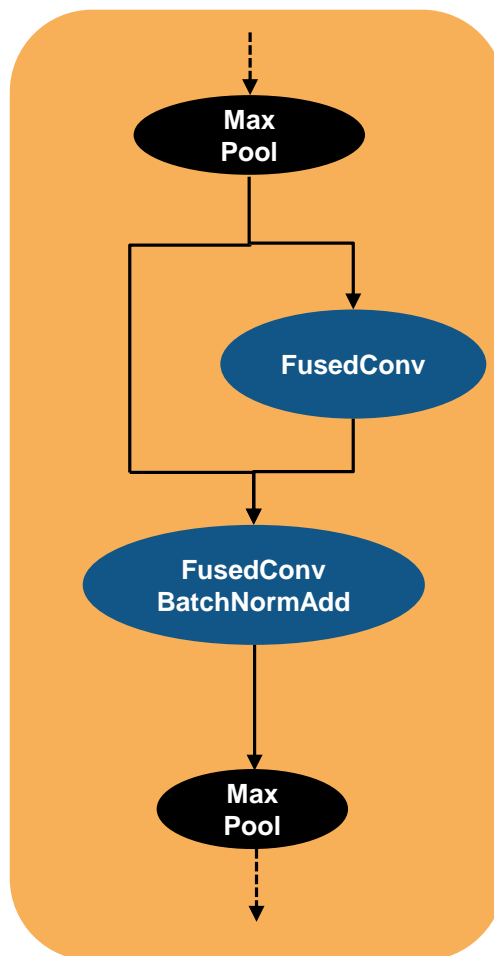
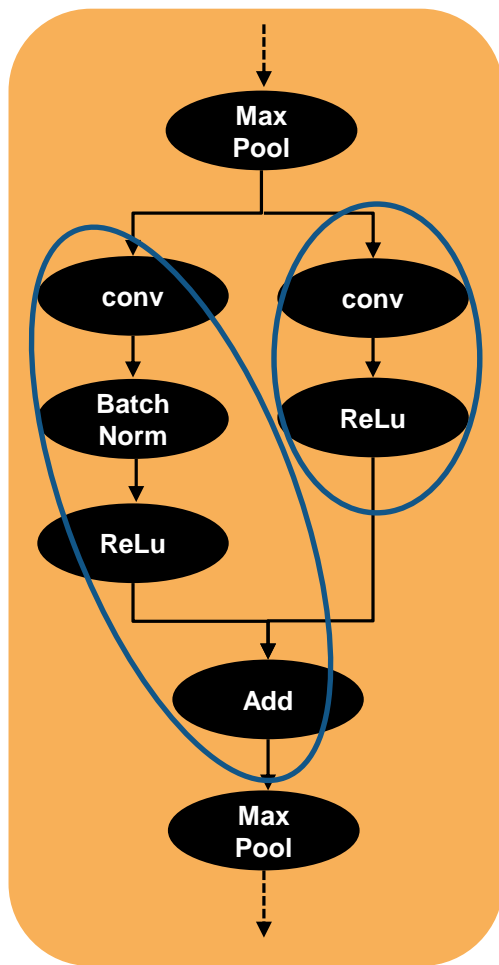
1. 优化的库
2. 网络优化
3. 编码模式



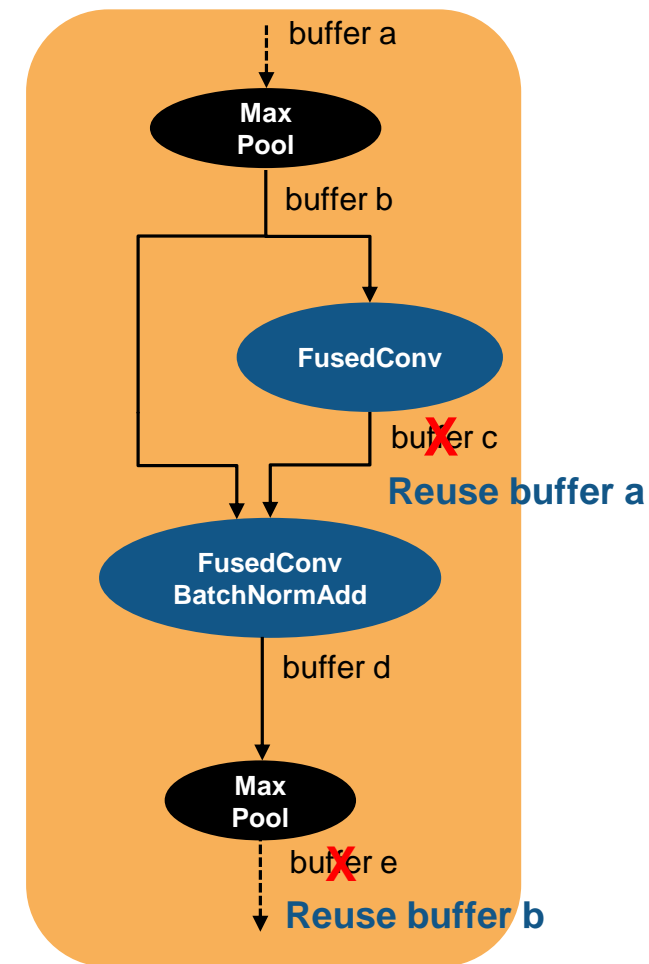
深度学习网络优化

性能

1. 优化的库
2. **网络优化**
3. 编码模式



层融合
优化计算



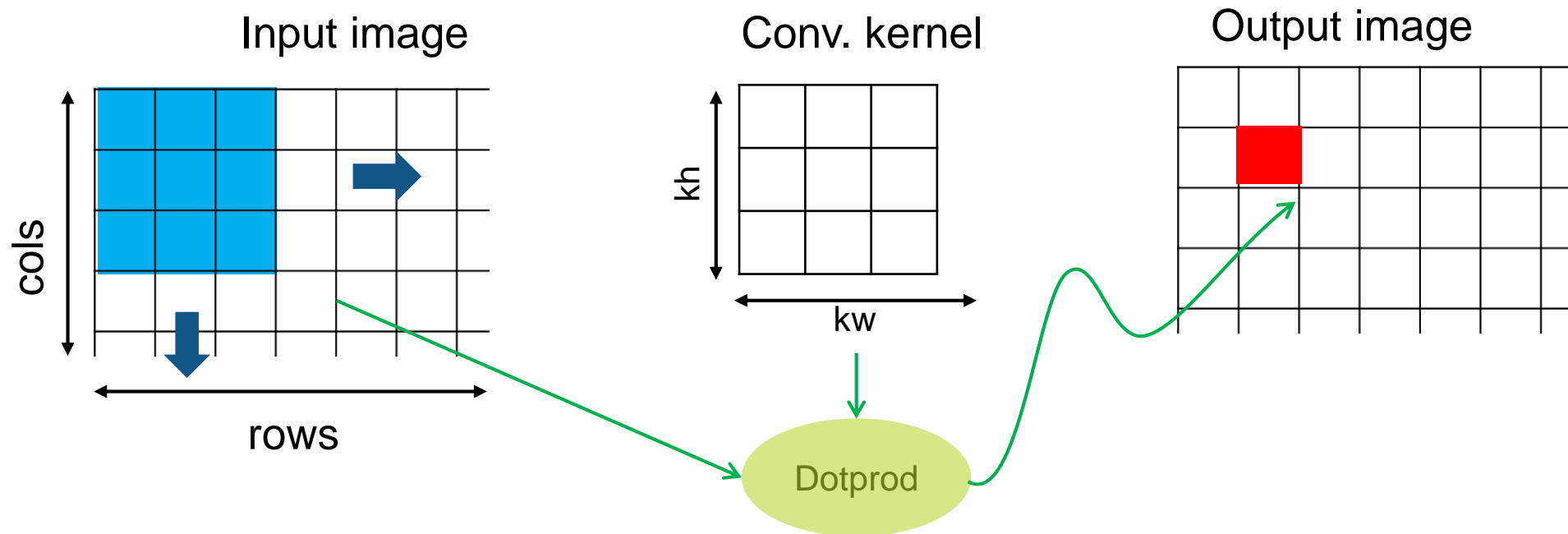
缓冲区最小化
优化内存

编码模式：Stencil 内核

性能

1. 优化的库
2. 网络优化
3. **编码模式**

- 自动应用于图像处理函数 (e.g. `imfilter`, `imerode`, `imdilate`, `conv2`, ...)
- 使用 `gpuCoder.stencilKernel()` 手动启用

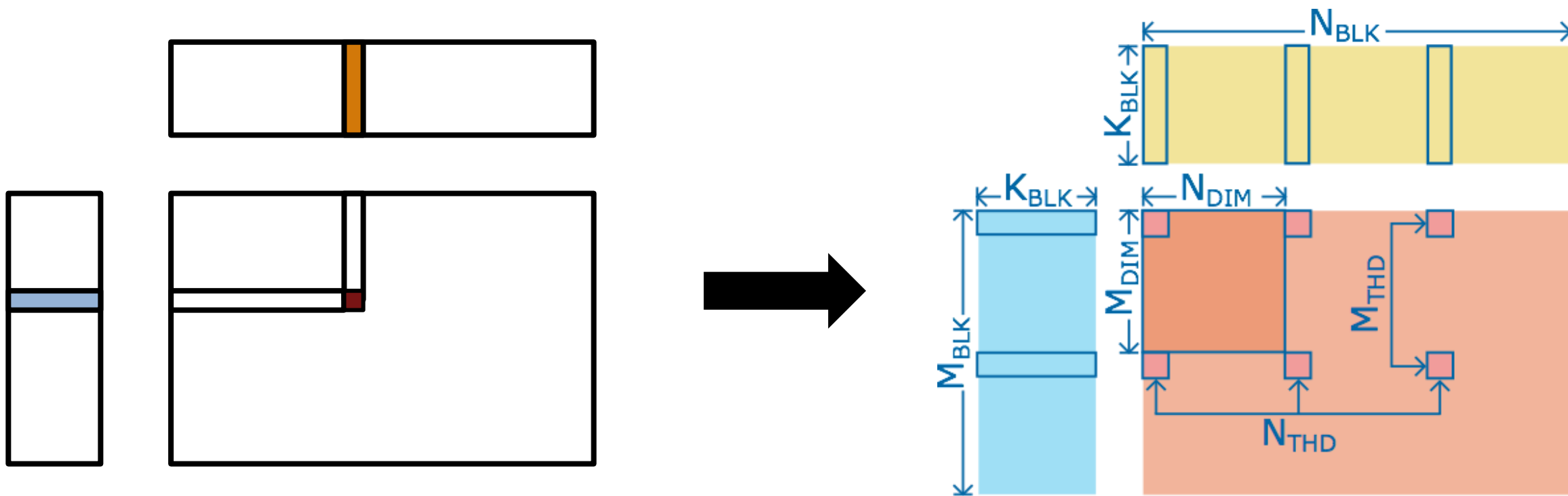


编码模式：Matrix-Matrix 内核

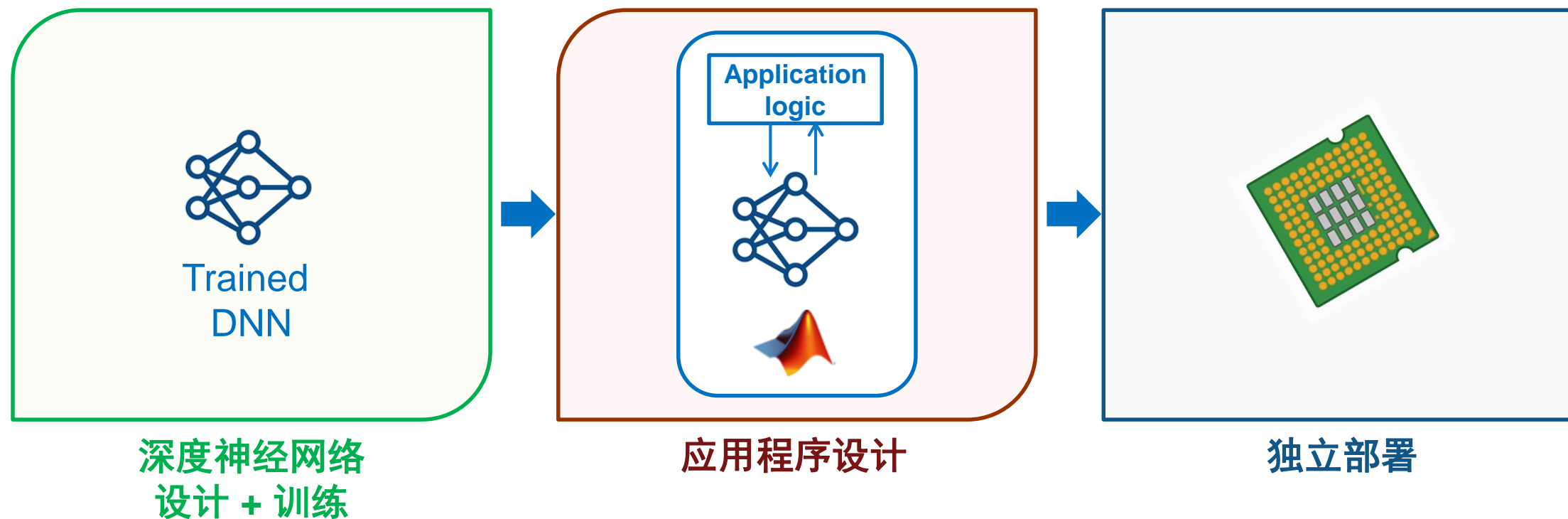
性能

1. 优化的库
2. 网络优化
3. **编码模式**

- 自动应用于很多 MATLAB 函数 (e.g. matchFeatures SAD, SSD, pdist, ...)
- 使用 `gpuCoder.matrixMatrixKernel()` 手动启动



MATLAB 深度学习流程



MATLAB 中的深度学习流程

