

2020  
embedded  
**VISION**  
summit<sup>®</sup>

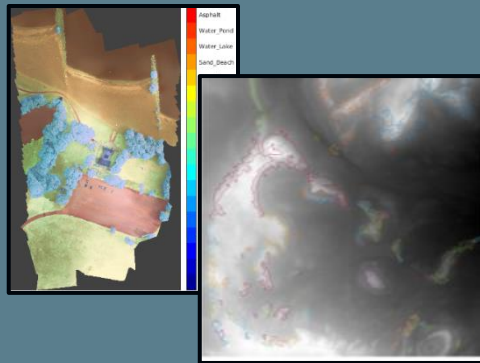
## Deploying Deep Learning Application on FPGAs with MATLAB

Jack Erickson  
Technical Marketing  
September 2020

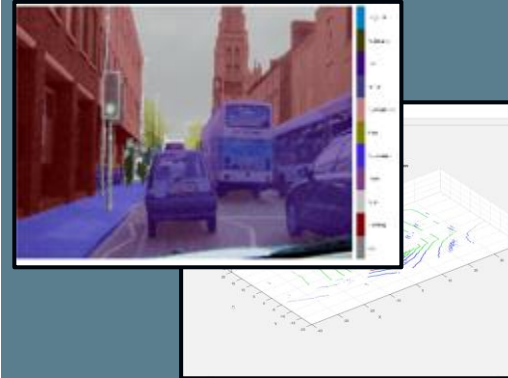


# Deep Learning Deployment on Embedded Devices

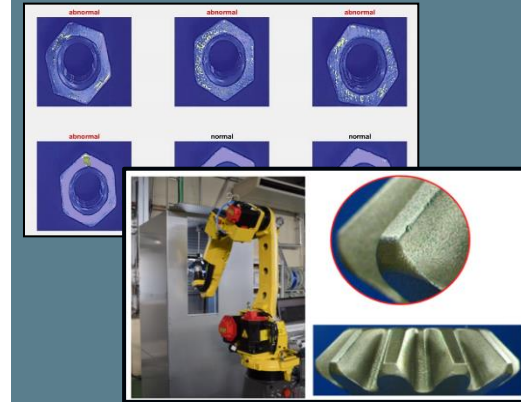
## Airborne Image Analysis



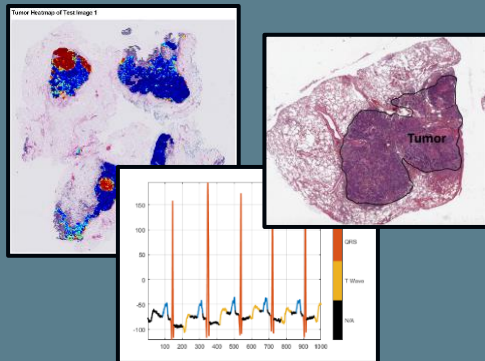
## Autonomous Driving



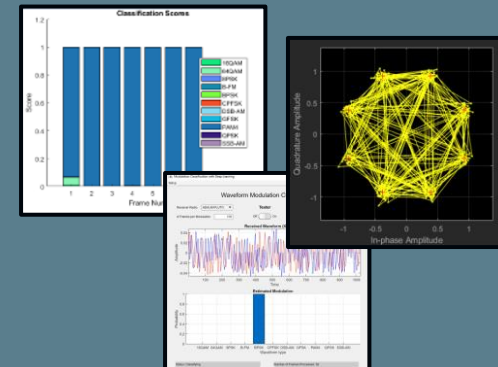
## Industrial Inspection



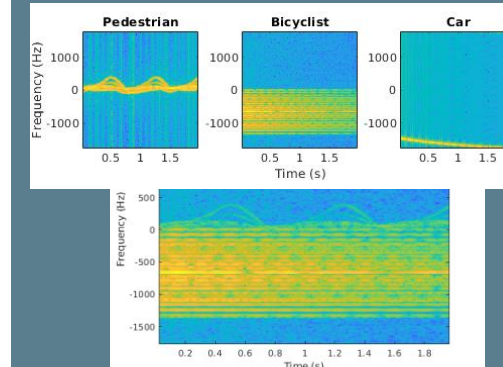
## Medical Image Analysis



## Wireless Modulation Classification



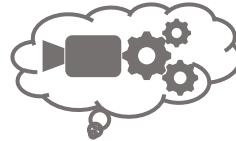
## Radar Signature Classification



# System Requirements Drive Network Design



Deep Learning  
Practitioner



Systems  
Engineer

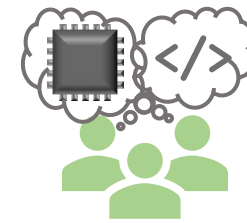
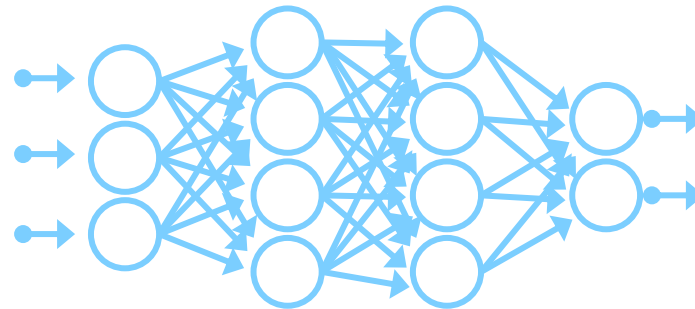
Camera specs

Accuracy

Latency

Cost

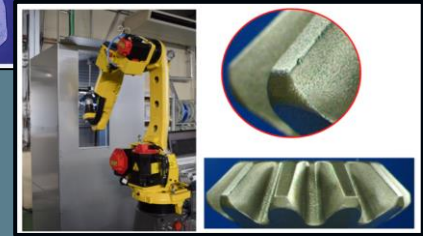
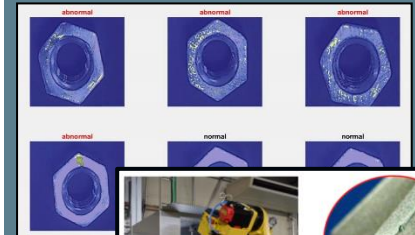
Power



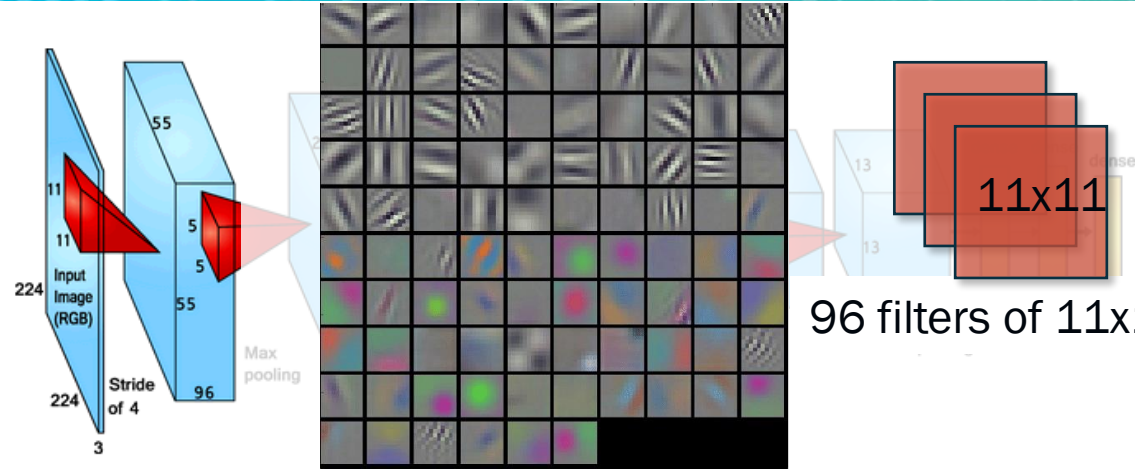
Hardware/Software  
Engineers



## Industrial Inspection

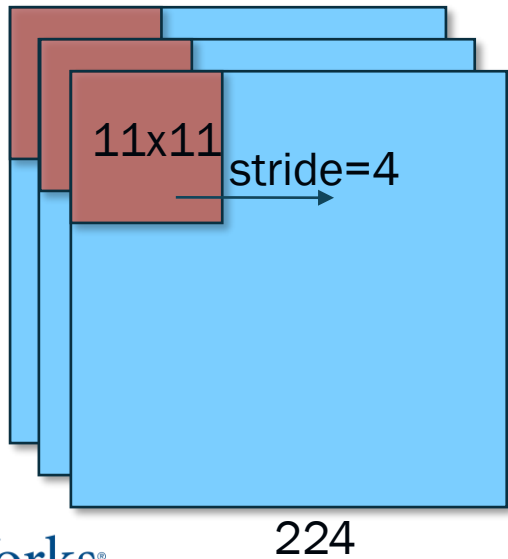


# Challenges of Deploying Deep Learning to FPGA Hardware: Convolution



96 filters of 11x11x3 of 32-bit parameters → 140k bytes

Figure 1. Illustration of AlexNet.

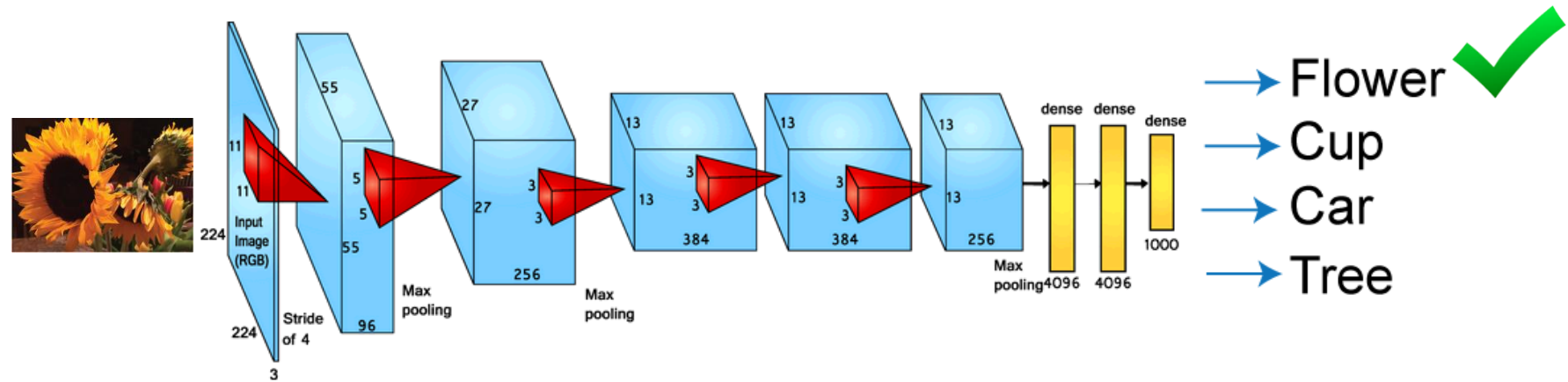


Each stride is an 11x11x3 matrix multiply-accumulate

→ 1.16M bytes of activations

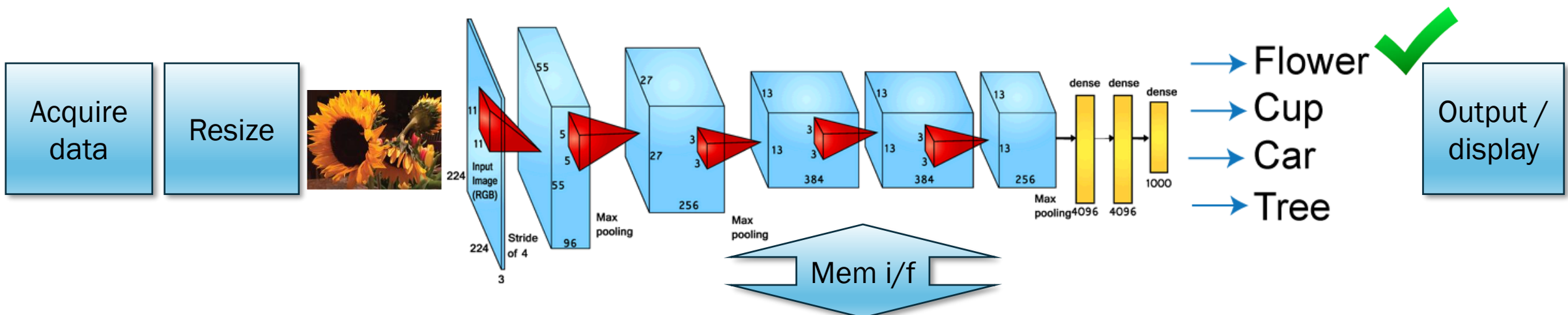
→ 105M floating-point multiply operations!

# Challenges of Deploying Deep Learning to FPGA Hardware




	input	conv 1	conv 2	conv 3	conv 4	conv 5	fc6	fc7	fc8	Total	
<b>Parameters (Bytes)</b>	n/a	140K	1.2M	3.5M	5.2M	1.8M	148M	64M	16M	<b>230 M</b>	➔ Off-chip RAM
<b>Activations (Bytes)</b>	588K	1.1M	728K	252K	252K	168K	16K	16K	4K	<b>3.1 M</b>	➔ Block RAM
<b>FLOPs</b>	n/a	105M	223M	149M	112M	74M	37M	16M	4M	<b>720 M</b>	➔ DSP Slices

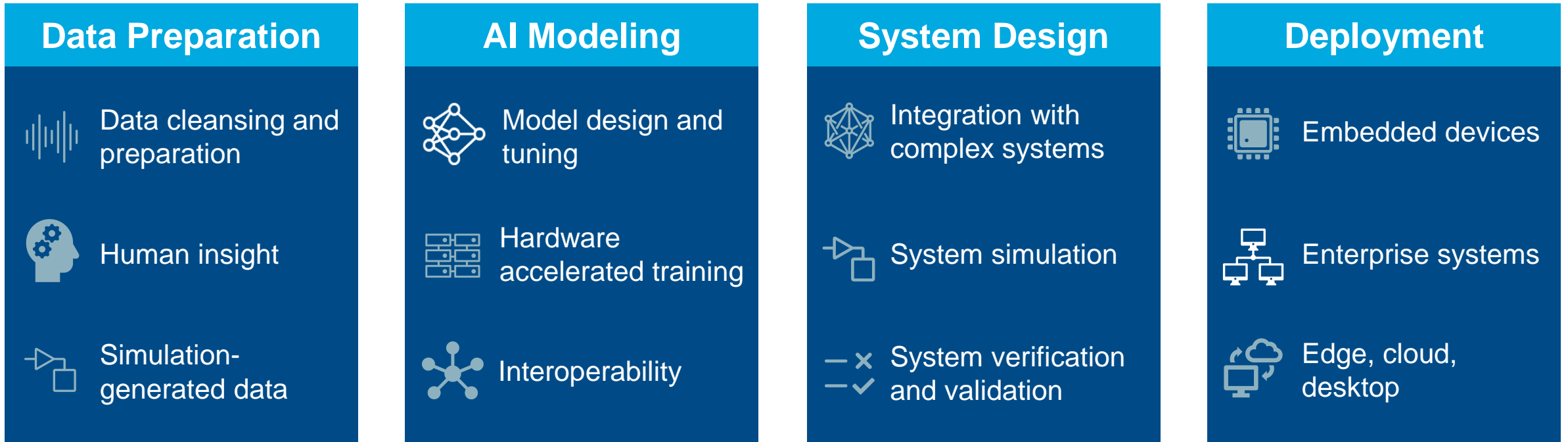
# Deploying Deep Learning to FPGA Hardware Requires Collaboration



	input	conv 1	conv 2	conv 3	conv 4	conv 5	fc6	fc7	fc8	Total
--	-------	--------	--------	--------	--------	--------	-----	-----	-----	-------

Parameters (Bytes)	<p style="text-align: center;"><b>Optimize</b></p> <ul style="list-style-type: none"> <li>• Network /layers</li> <li>• Fixed-point quantization</li> <li>• Processor micro-architecture</li> </ul> 									
Activations (Bytes)										
FLOPs										

# A Collaborative AI Workflow



 **Iteration and Refinement** 

# Design and Analyze Your Networks in MATLAB

## AI Modeling



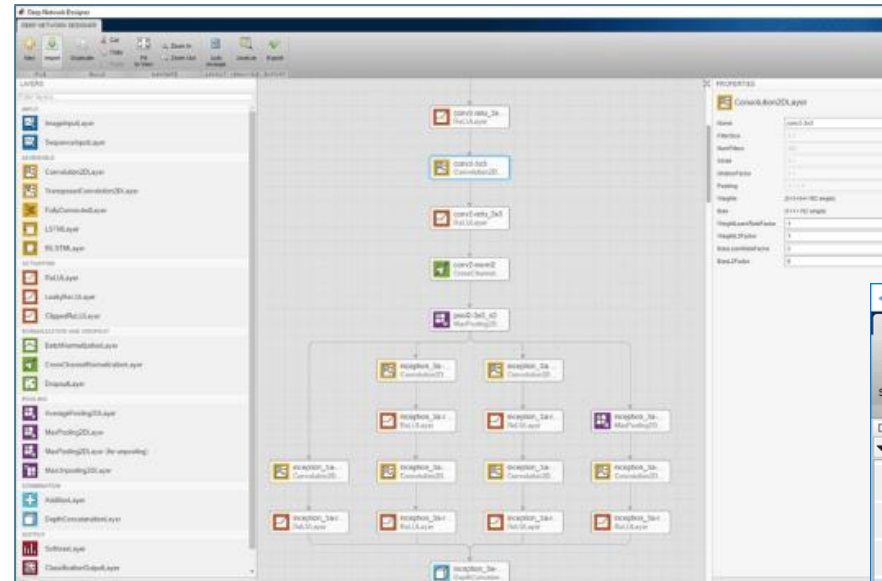
Model design and tuning



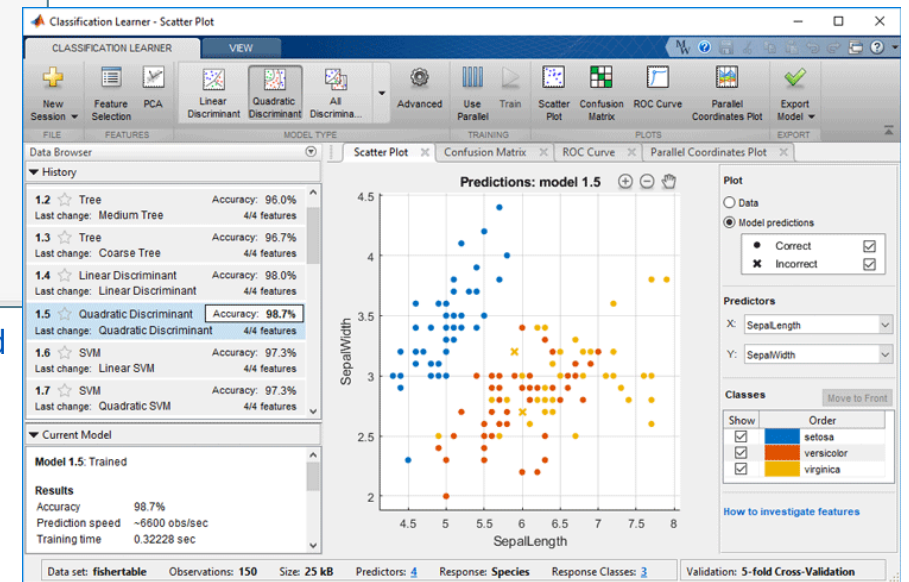
Hardware accelerated training



Interoperability



Deep Network Designer app to build, visualize, and edit deep learning networks



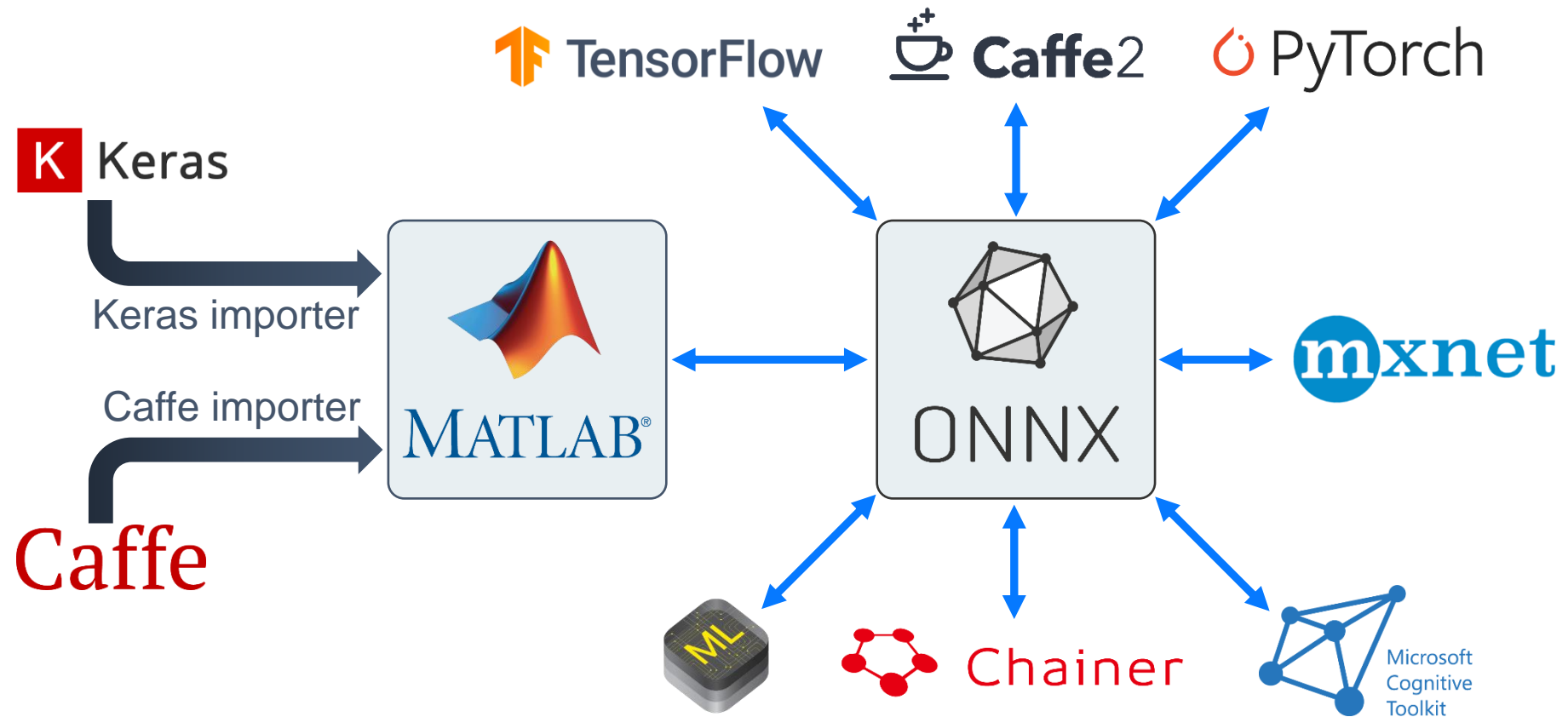
Classification Learner app to try different classifiers and find the best fit for your data set



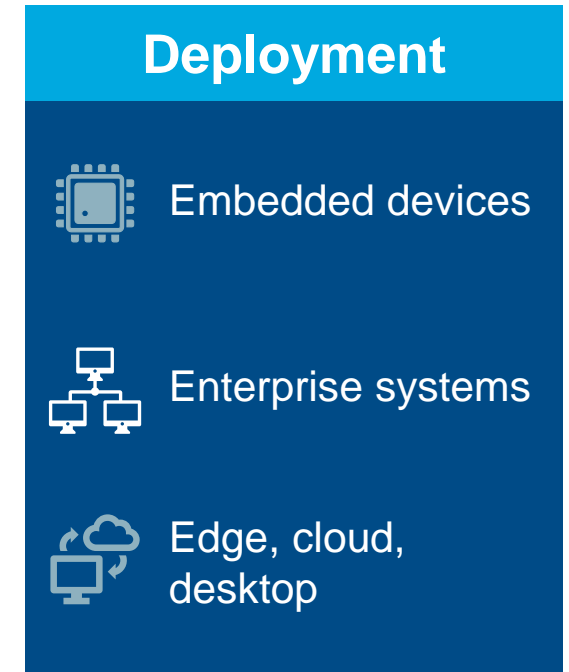
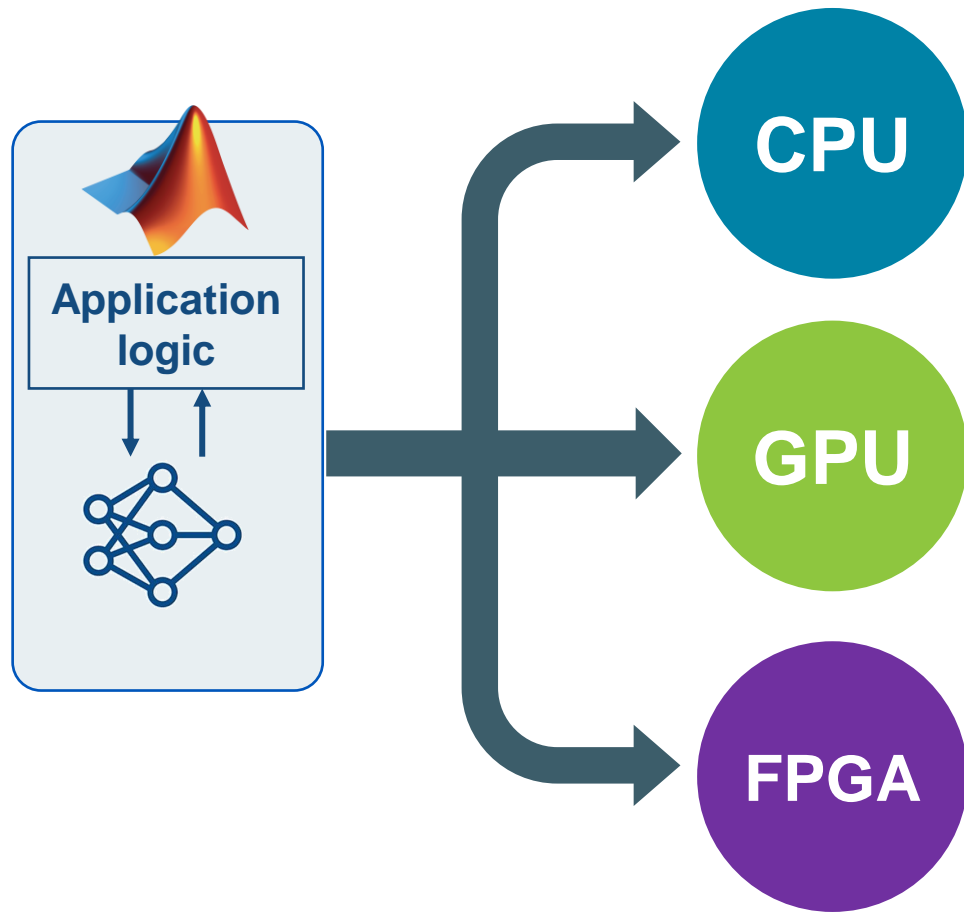
# MATLAB Interoperates with Other AI Frameworks

**AI Modeling**

- Model design and tuning
- Hardware accelerated training
- Interoperability

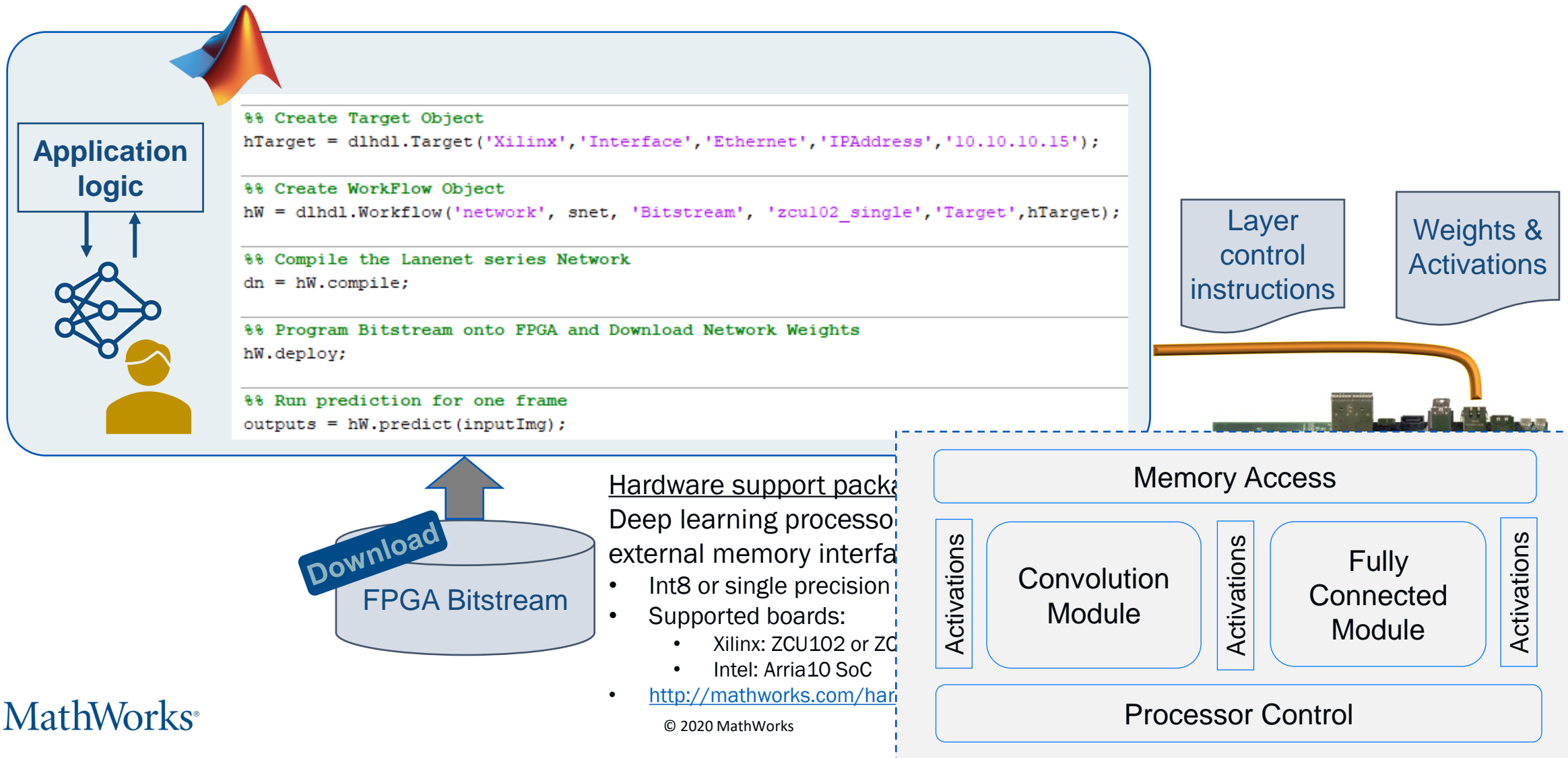


# Deploy from MATLAB to a Variety of Hardware Platforms

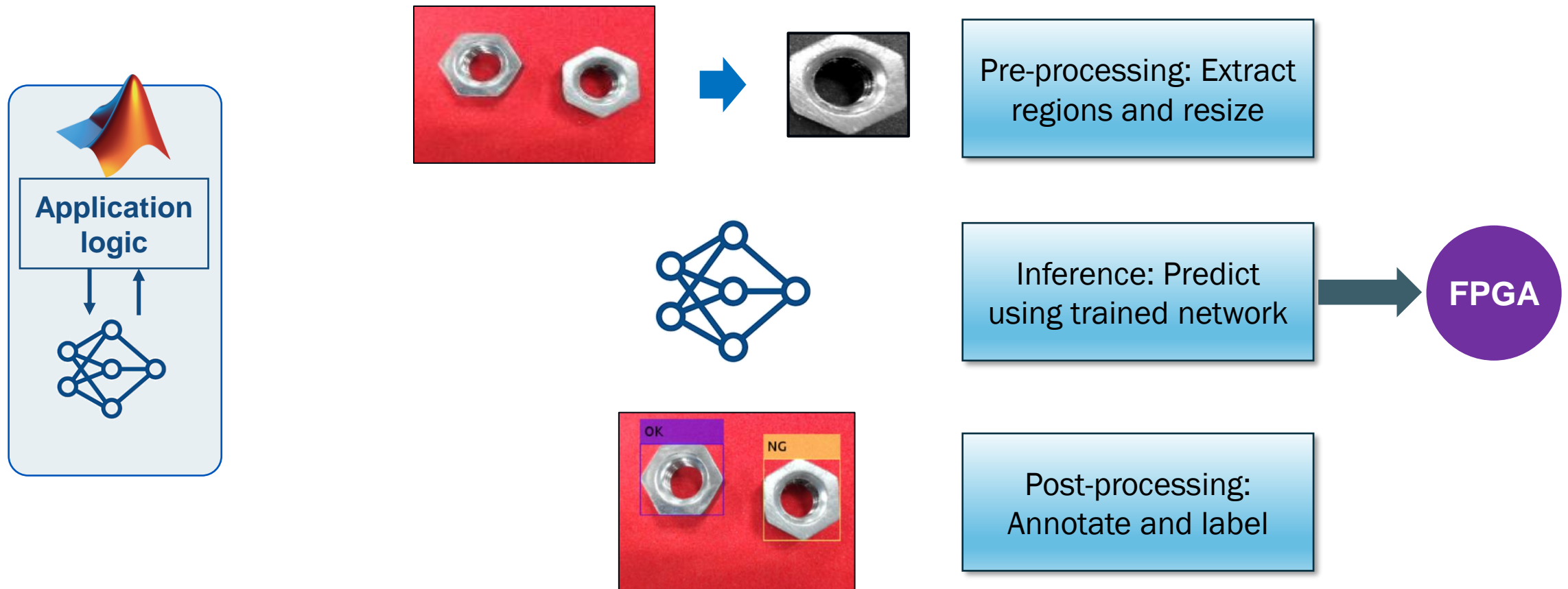


# FPGA Deployment from MATLAB

# Get Started Prototyping on FPGA with Deep Learning HDL Toolbox™



# Defect Detection Example



# Run Deep Learning on FPGA from MATLAB

HOME PLOTS APPS LIVE EDITOR INSERT VIEW Search Documentation Sign In

C:\Users\jerickso\Documents\MATLAB\DefectDetection\dnnfpga\_defectdetection.1

Live Editor - C:\Users\jerickso\Documents\MATLAB\DefectDetection\DefectDetectionExample.mlx

## Defect Detection

**Prerequisites**

- Xilinx ZCU102 SoC development kit
- Deep Learning HDL Toolbox™ Support Package for Xilinx FPGA and SoC
- Deep Learning Toolbox™
- Deep Learning HDL Toolbox™

**Create Folder and Copy Relevant Files**

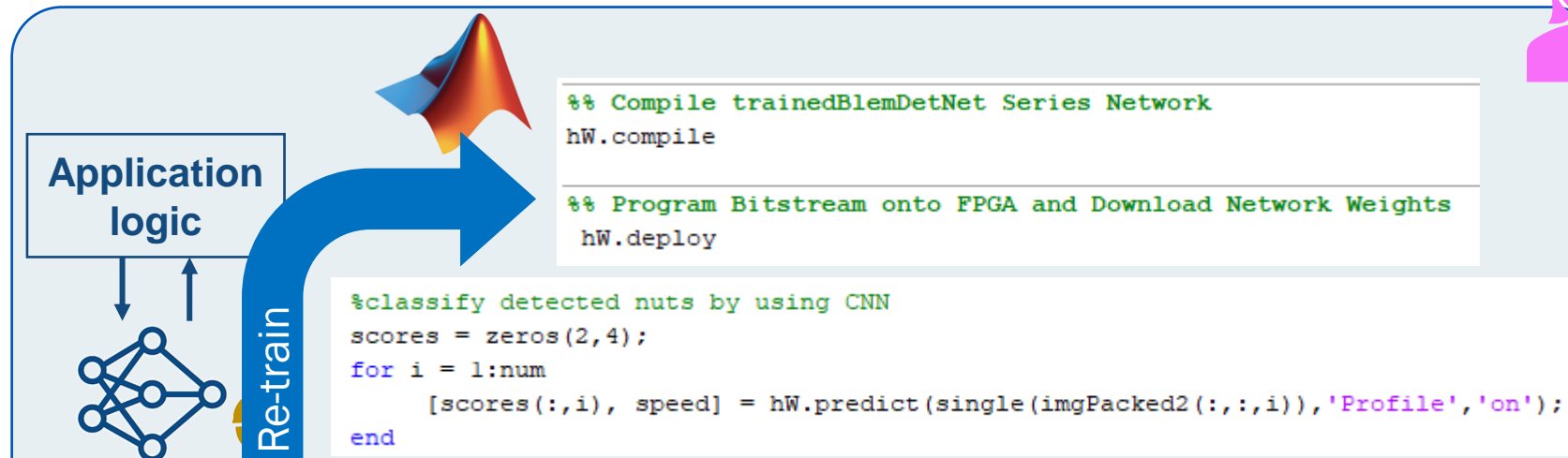
```
1 unzip('dnnfpga_defectdetection.zip');  
2 [newDir, origDir] = cloneSetupDir('dnnfpga_defectdetection');  
3 cd(newDir);
```

Command Window

```
fx >>  
<
```

UTF-8 script Ln 10 Col 31

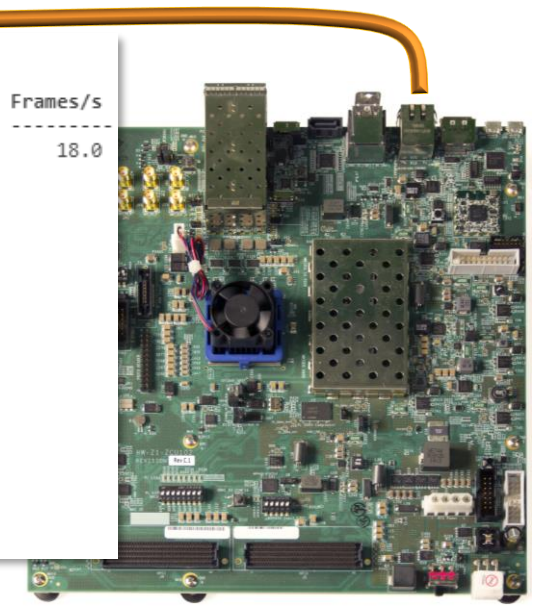
# Prototyping: Design Exploration and Customization



Deep Learning Processor Profiler Performance Results

	LastLayerLatency(cycles)	LastLayerLatency(seconds)	FramesNum	Total Latency	Frames/s
Network	12213262	0.05551	1	12213302	18.0
conv_module	3292045	0.01496			
conv1	412728	0.00188			
norm1	173252	0.00079			
pool1	58636	0.00027			
conv2	656582	0.00298			
norm2	128169	0.00058			
pool2	53269	0.00024			
conv3	780456	0.00355			
conv4	600050	0.00273			
conv5	408977	0.00186			
pool5	20059	0.00009			
fc_module	8921217	0.04055			
fc6	1759800	0.00800			
fc7	7030644	0.03196			
fc8	130771	0.00059			

\* The clock frequency of the DL processor is: 220MHz



# Design Exploration and Customization

The screenshot shows the MATLAB Live Editor interface with the following code blocks:

```
11 hT = dlhdl.Target('Xilinx','Interface','Ethernet','IPAddress','10.10.10.15')

12 hW = dlhdl.Workflow('Network',snet_defnet,'Bitstream','zcu102_single','Target',h

13 hW.compile

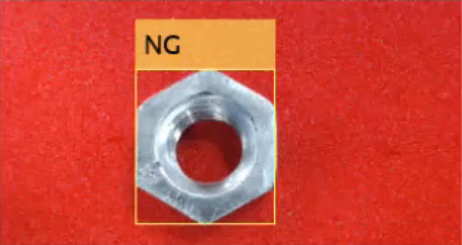
14 hW.deploy

15 unzip('testImages.zip')
16
17 filename=[pwd,'/testImages/ng1.png'];
18 img=imread(filename);
19 predictDefect(hW, img);
```

**Deep Learning Processor Profiler Performance Results**

Network	LastLayerLatency(cycles)	LastLayerLatency(seconds)
conv_module	12213262	0.05551
conv1	3292045	0.01496
conv2	412728	0.00188
conv3	173252	0.00079
conv4	58636	0.00027
conv5	656582	0.00298
conv6	128169	0.00058
conv7	53269	0.00024
conv8	780456	0.00355
conv9	600050	0.00273
conv10	408977	0.00186
conv11	20059	0.00009
conv12	8921217	0.04055
conv13	1759800	0.00800
conv14	7030644	0.03196
conv15	130771	0.00059

\* The clock frequency of the DL processor is: 220MHz





# Optimizing Deep Learning Applications Requires Collaboration

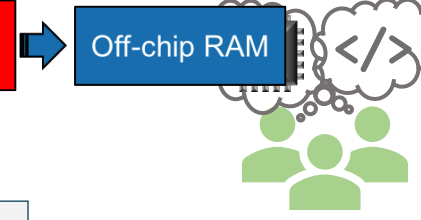


Systems  
Engineer

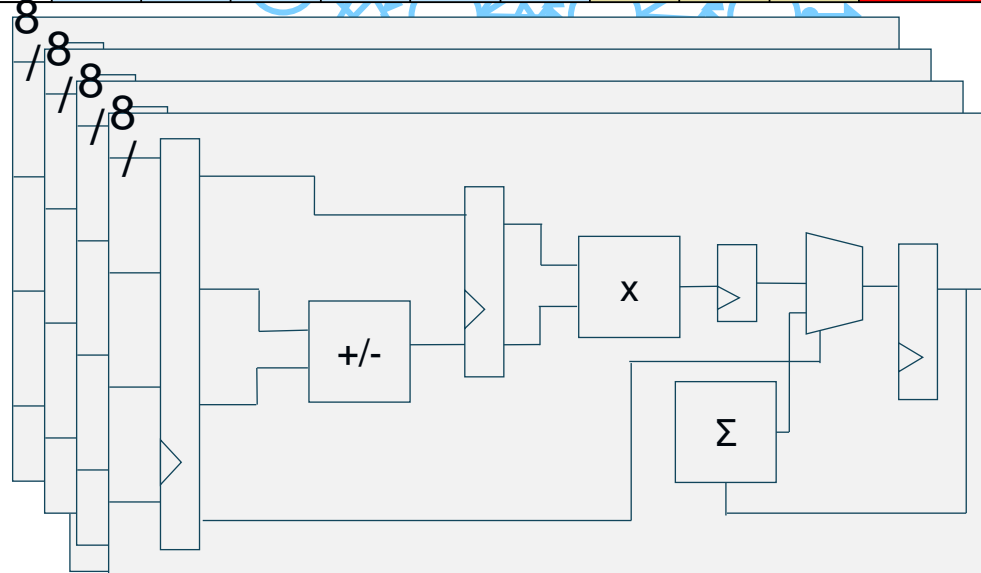


Deep Learning  
Practitioner

	input	conv 1	conv 2	conv 3	conv 4	conv 5	fc6	fc7	fc8	Total
Parameters (Bytes)	n/a	35K	0.4M	0.9M	1.3M	0.5M	37M	16M	4M	58 M



Hardware/Software  
Engineers



# INT8 Quantization

The screenshot displays the MATLAB Live Editor interface for a project named 'DefectDetection'. The editor shows a script with the following code:

```
22     websave('trainedBlemDetNet.mat',url);  
23 end  
24 net2 = load('trainedBlemDetNet.mat');  
25 snet_blemdetnet = net2.convnet  
26 analyzeNetwork(snet_blemdetnet)
```

The workflow steps are:

- 27 Create Workflow Object for trainedBlemDetNet Network  
`hw = dlhdl.Workflow('Network',snet_blemdetnet,'Bitstream','zcu102_single','Target');`
- 28 Compile trainedBlemDetNet Series Network  
`hw.compile`
- 29 Program Bitstream onto FPGA and Download Network Weights  
`hw.deploy`
- Run Prediction for One Image  
`filename=[pwd,'/testImages/ok1.png'];  
img=imread(filename);  
predictDefect(hw, img);`

The 'Performance Results' table shows the following data:

LastLayerLatency(seconds)	FramesNum	Total Latency	Frames/s
0.02222	1	4887512	45.0
0.00571			
0.00212			
0.00087			
0.00072			
0.00181			
0.00019			
0.01651			
0.01643			
0.00007			


The image prediction result shows a metal nut on a red background with a bounding box and the label 'OK'.

Command Window: `fx >>`

File: UTF-8 | script | Ln 30 | Col 38

## Iterate and Converge on Deep Learning FPGA Deployment from MATLAB

**Application logic**



```

% Create target object
hTarget = dlhdl.Target(...);

% Create workflow object, using the target
hW = dlhdl.Workflow(...);

% Compile the network
hW.compile;

% Program the bitstream and deploy the compiled network and weights
hW.deploy;

% Run prediction
[score, speed] = hW.predict(img, 'Profile', 'on');
                
```

>> deepNetworkQuantizer

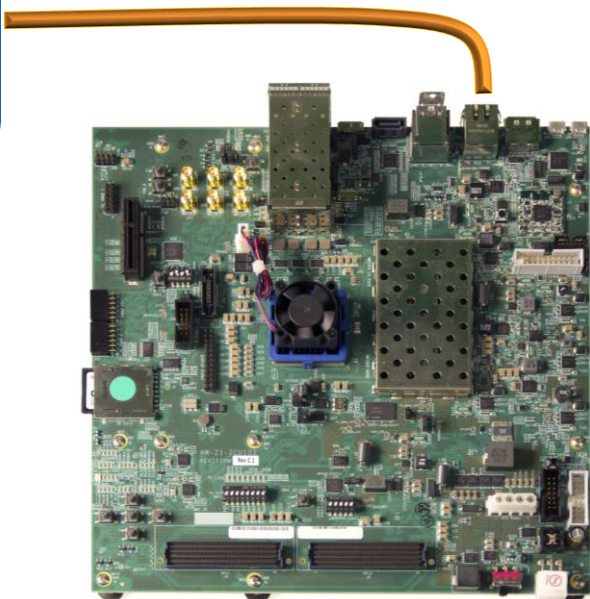
Layer control instructions

Weights & Activations

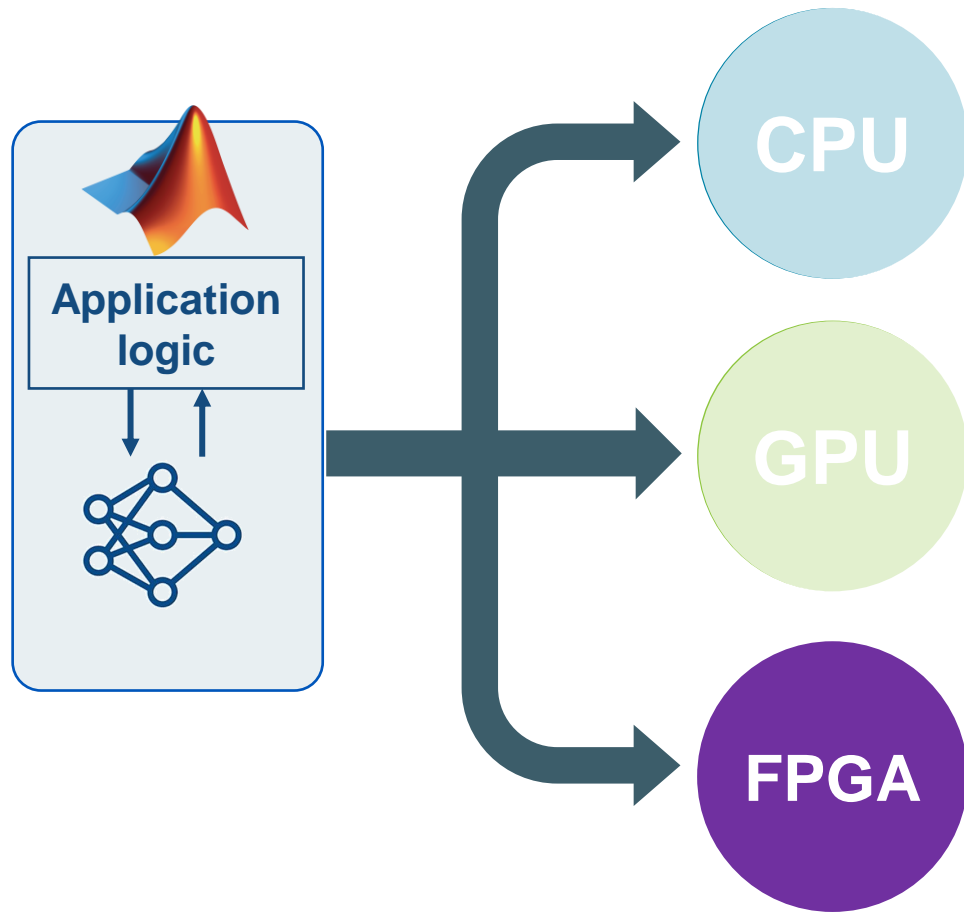
Quantize

Parameters	Speed
<b>140 MB</b>	<b>18 fps</b>
<b>84 MB</b>	<b>45 fps</b>
<b>68 MB</b>	<b>139 fps</b>




Generate HDL



# Deploy from MATLAB to a Variety of Hardware Platforms



### Deep Learning HDL Toolbox

-  Prototype from MATLAB
-  Tune for system requirements
-  Configure and generate RTL

Deep Learning Solutions in MATLAB

<https://www.mathworks.com/solutions/deep-learning.html>

Deep Learning HDL Toolbox

<https://www.mathworks.com/products/deep-learning-hdl.html>

Onramp: Deep Learning in MATLAB

<https://www.mathworks.com/learn/tutorials/deep-learning-onramp.html>

MathWorks FPGA Solutions Page

<https://www.mathworks.com/solutions/fpga-asic-soc-development.html>