# MATLAB EXPO

## 임베디드 시스템 적용을 위한 AI 개발

*신행재 부장, 매스웍스코리아*
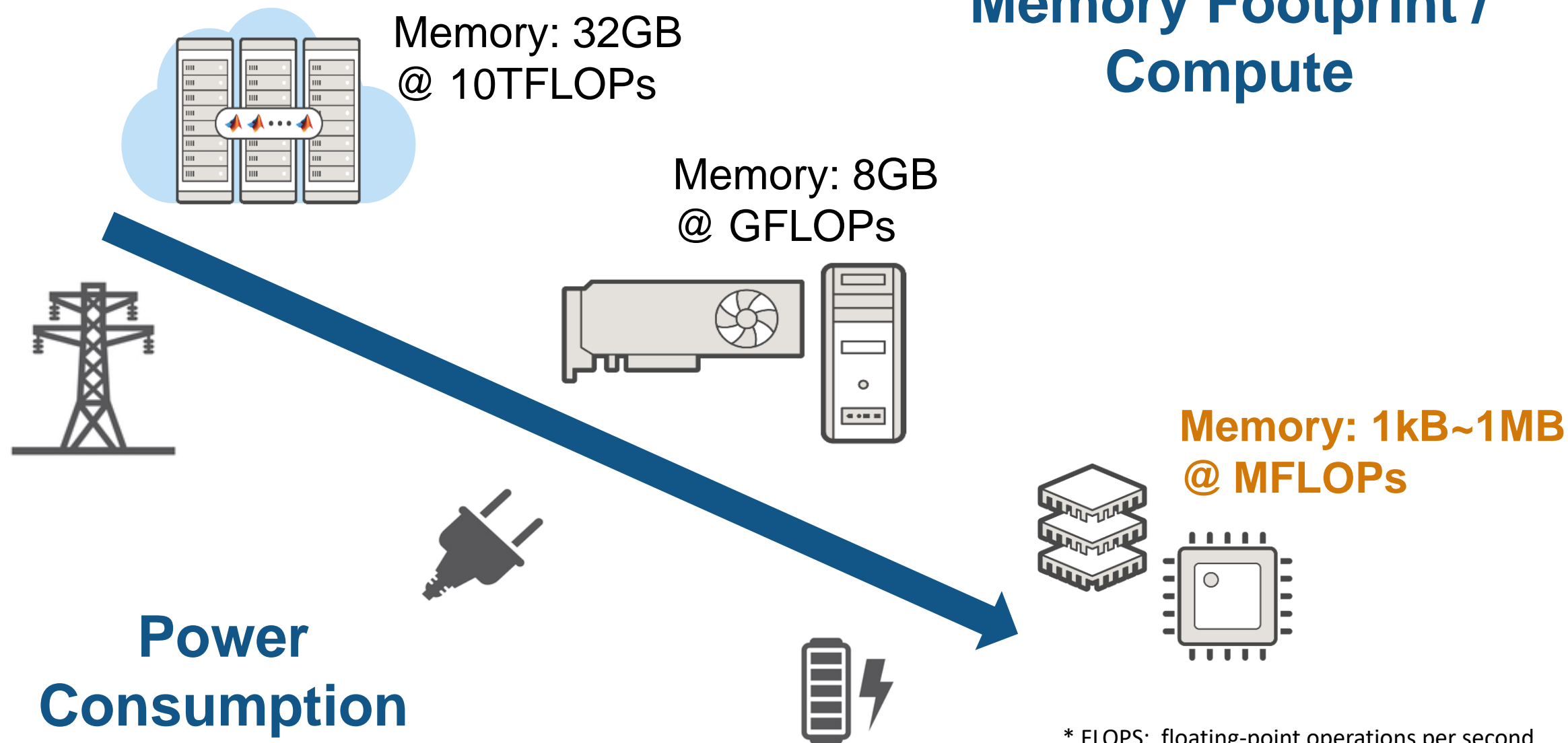
# Edge AI innovates many industries!
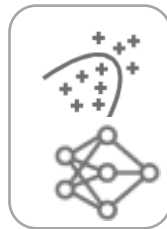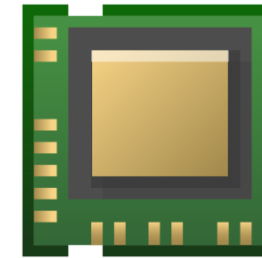
# Hardware Constraints

**Memory Footprint / Compute**

Memory: 32GB
@ 10TFLOPs

Memory: 8GB
@ GFLOPs

**Memory: 1kB~1MB
@ MFLOPs**

**Power
Consumption**

* FLOPS: floating-point operations per second

# What is "Edge" (Embedded) AI?
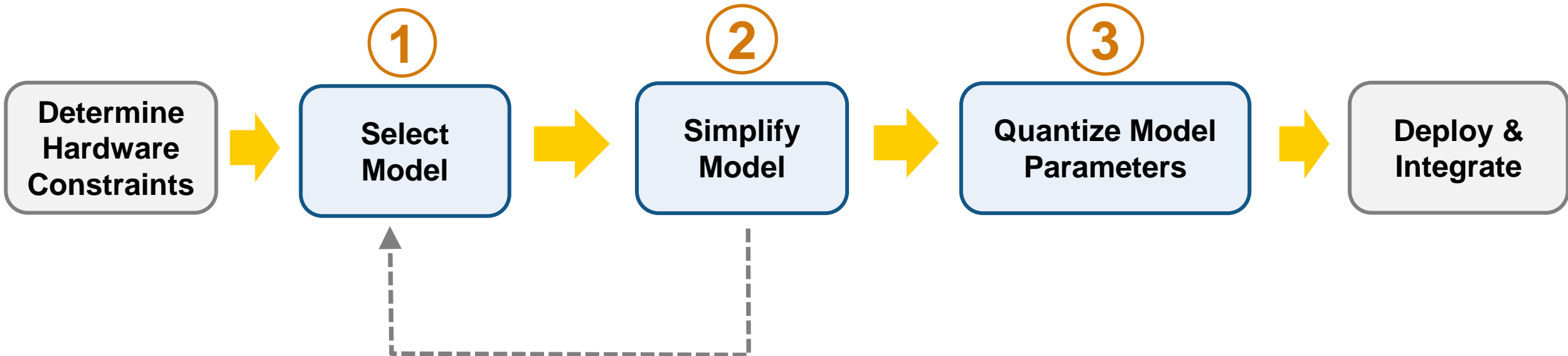
# Why is Edge AI (Model Compression) difficult?



AI is often big



Knowledge Gap

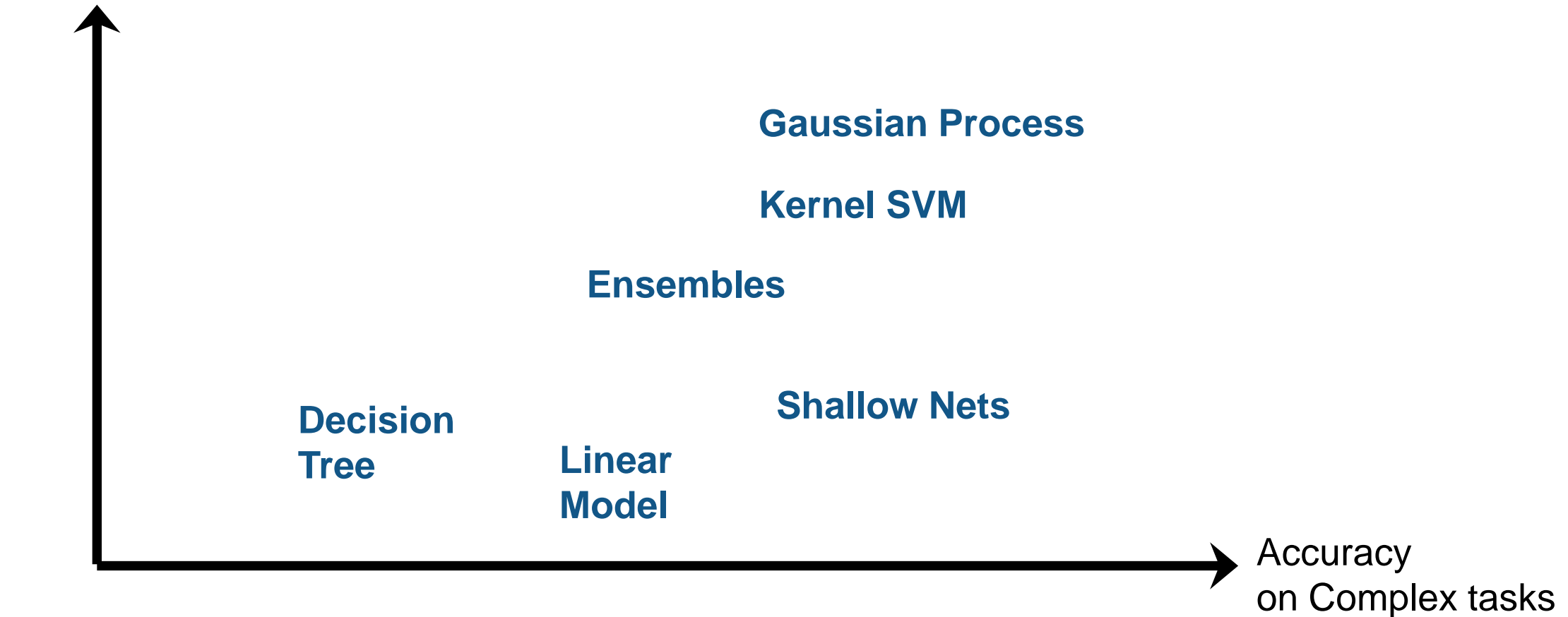# Model Compression Workflow

# Compressing Machine Learning

# Step ① Size aware model selection

Small    Large

Size /
Execution Time

Deep Neural Net

Gaussian Process

Kernel SVM

Ensembles

Shallow Nets

Decision
Tree

Linear
Model

Accuracy
on Complex tasks

# Model Compression Workflow for Machine Learning

Small    Large

① **Classification / Regression Learner**

② **In-App Feature Selection**

**Bayesopt**

③ **Fixed Point Designer / Native Simulink Block**

**Simplify Model**

**Select Features**

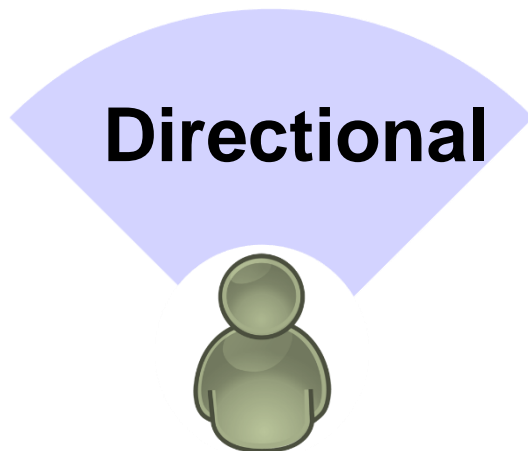**Determine Hardware Constraints** → **Select (Initial) Model** → **Tune Hyper-parameters** → **Quantize Model Parameters** → **Deploy & Integrate**

# Demo: Embedding AI in an intelligent Hearing Aid



0.5 to 256 kB
on-chip memory

**Directional**

**All Around**

# Demo: Fit Machine Learning for Intelligent Hearing Aid

# **Machine Learning Demo** Size Reduction by factor 20



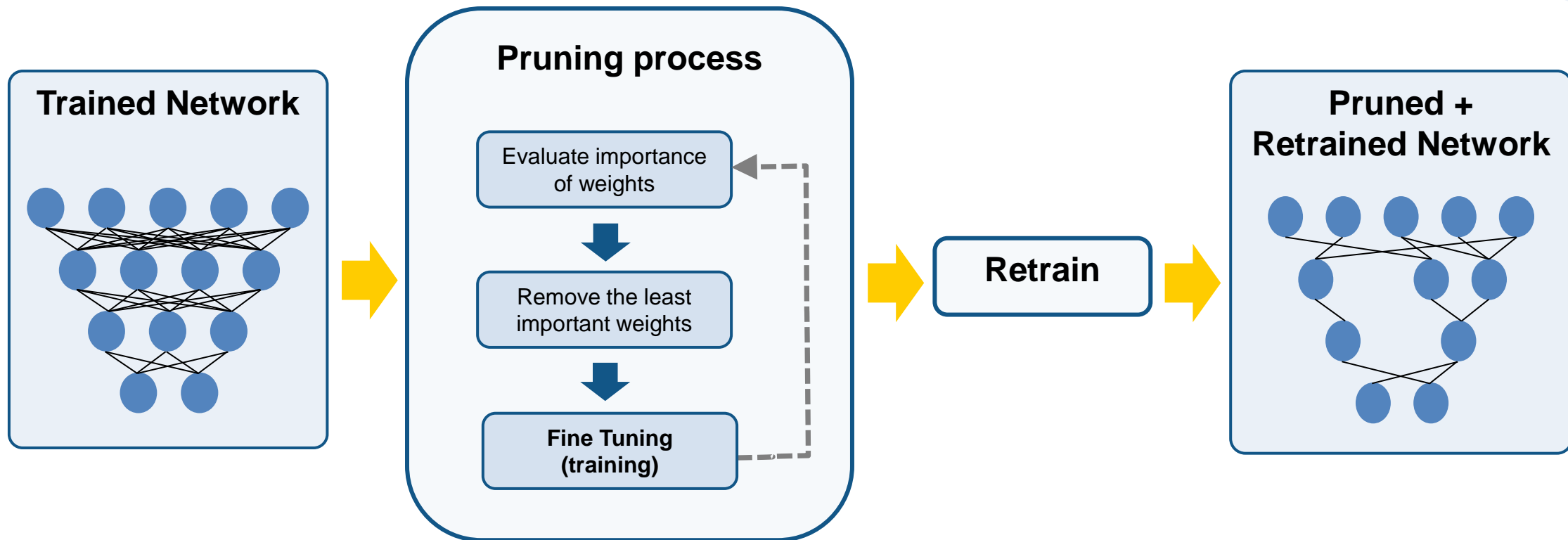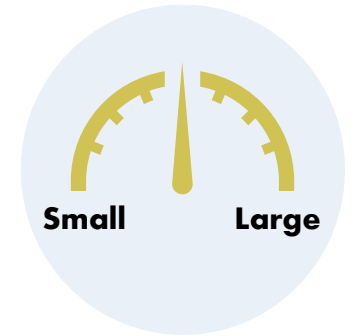Target Hardware
Size: 50 kB

# Compressing Deep Learning

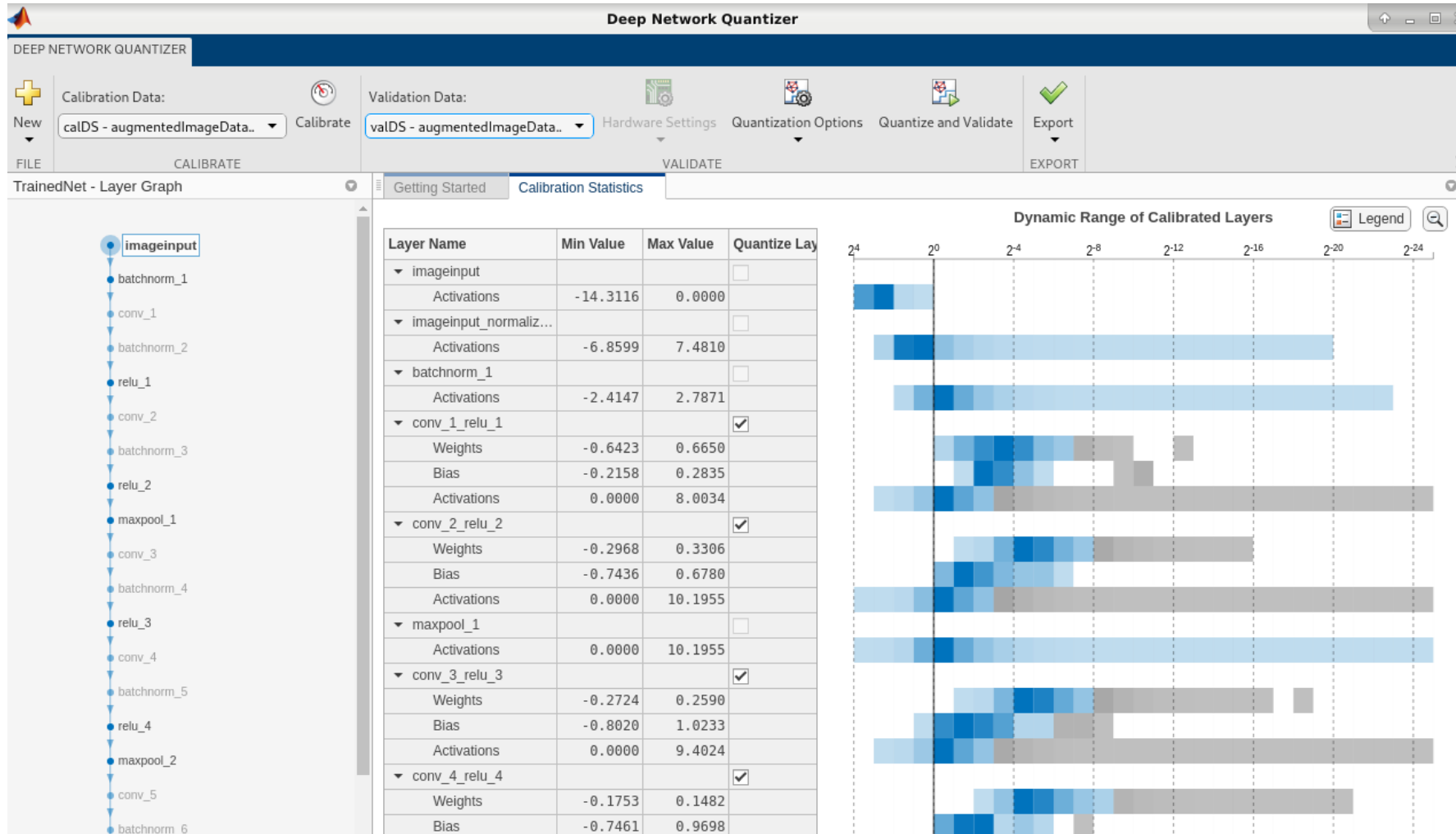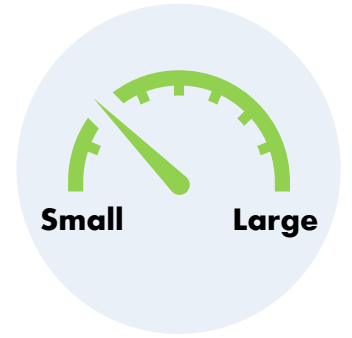# Step ① Size aware model selection

# Step ② Smart pruning

Remove **unimportant** parts of the network

# Step ③ Quantize your model

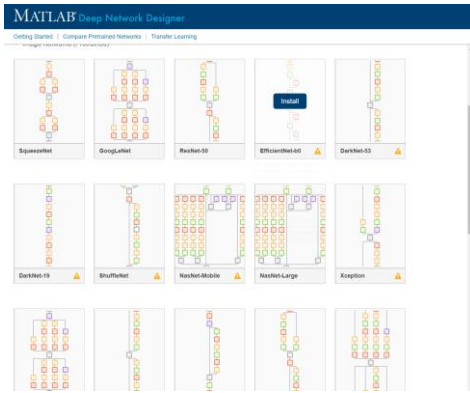# Deep Learning Demo: Scene classification

**Classify 10 classes**

More difficult problem → more complex model

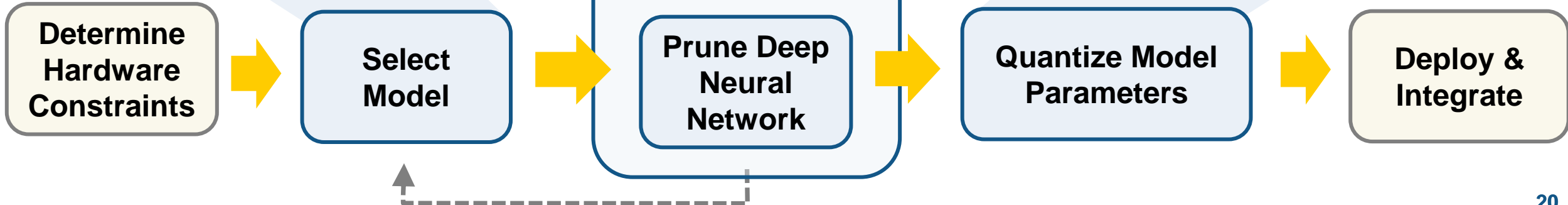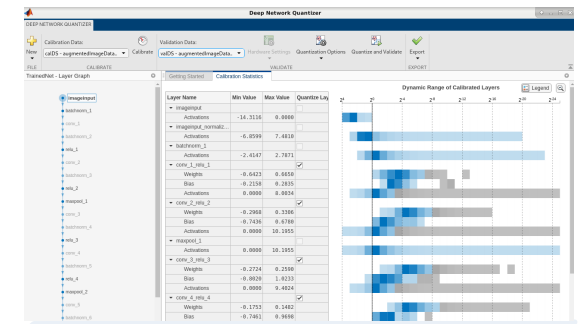# Functionality for Compressing Deep Neural Nets

① **Deep Network Designer**

② **Taylor Pruning**

③ **Deep Network Quantizer**



```
taylorPrunableNetwork(net)
```

**Determine Hardware Constraints** → **Select Model** → **Simplify Model** **Prune Deep Neural Network** → **Quantize Model Parameters** → **Deploy & Integrate**

# Step 1: Select Model

①

Select
Model

## Load original trained CNN model and dataset

```
1    load('trained10classNetwork'); |
2    load('data')
```
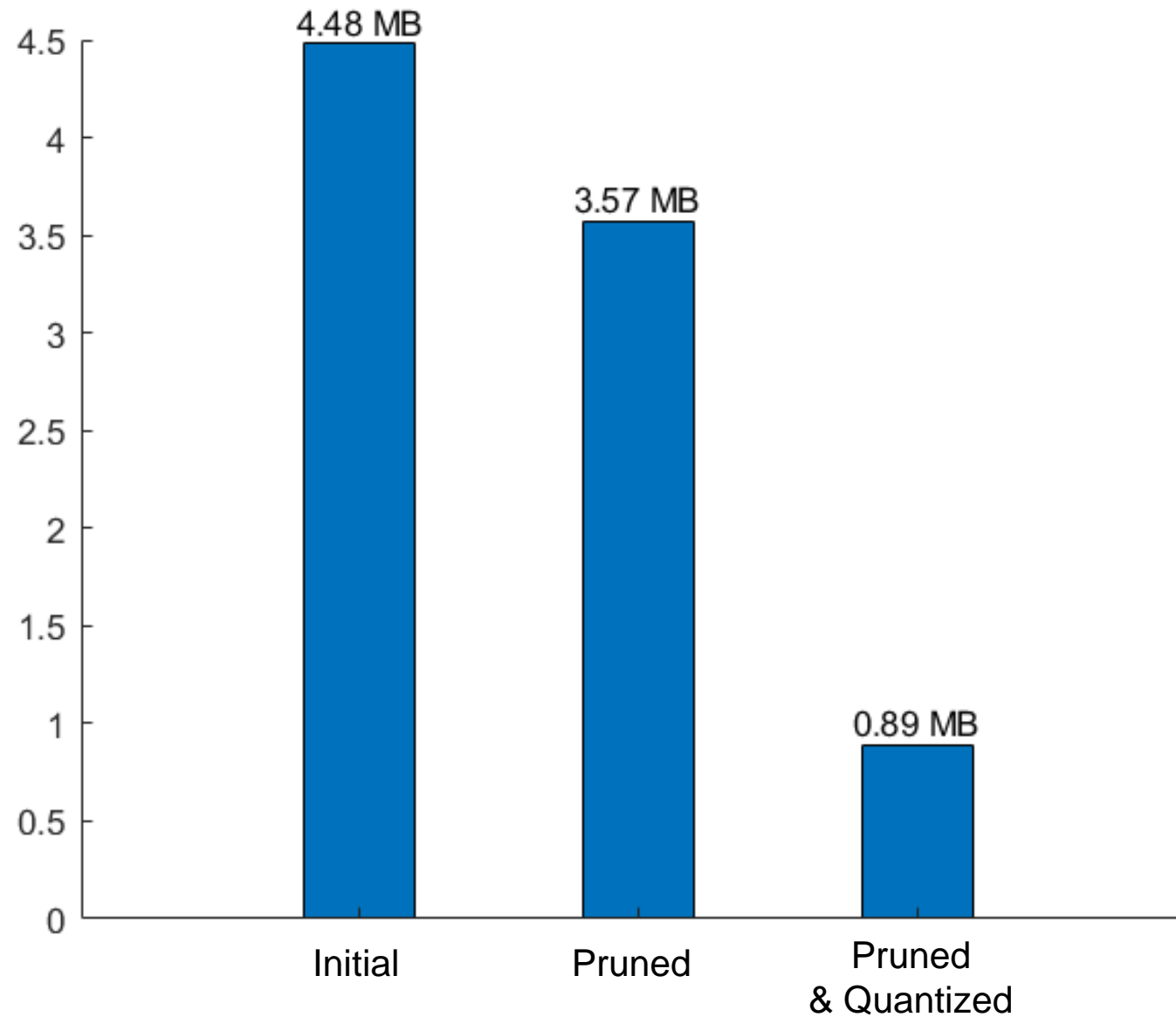
Note: Sounds have been converted to spectograms



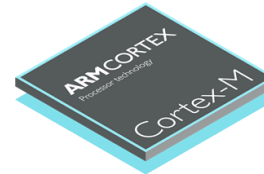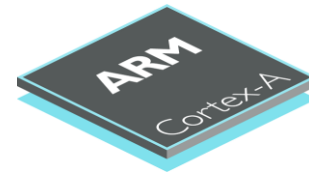Segment 1
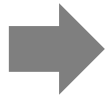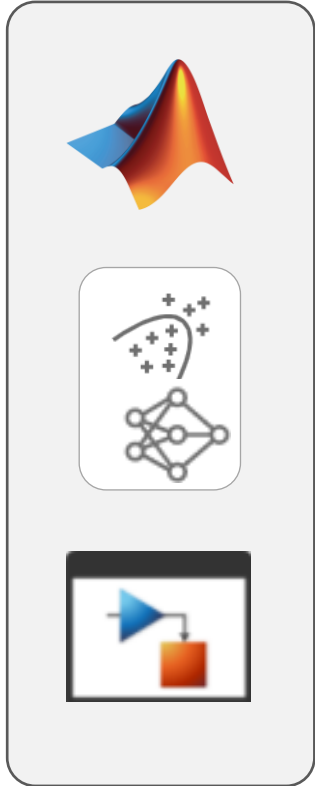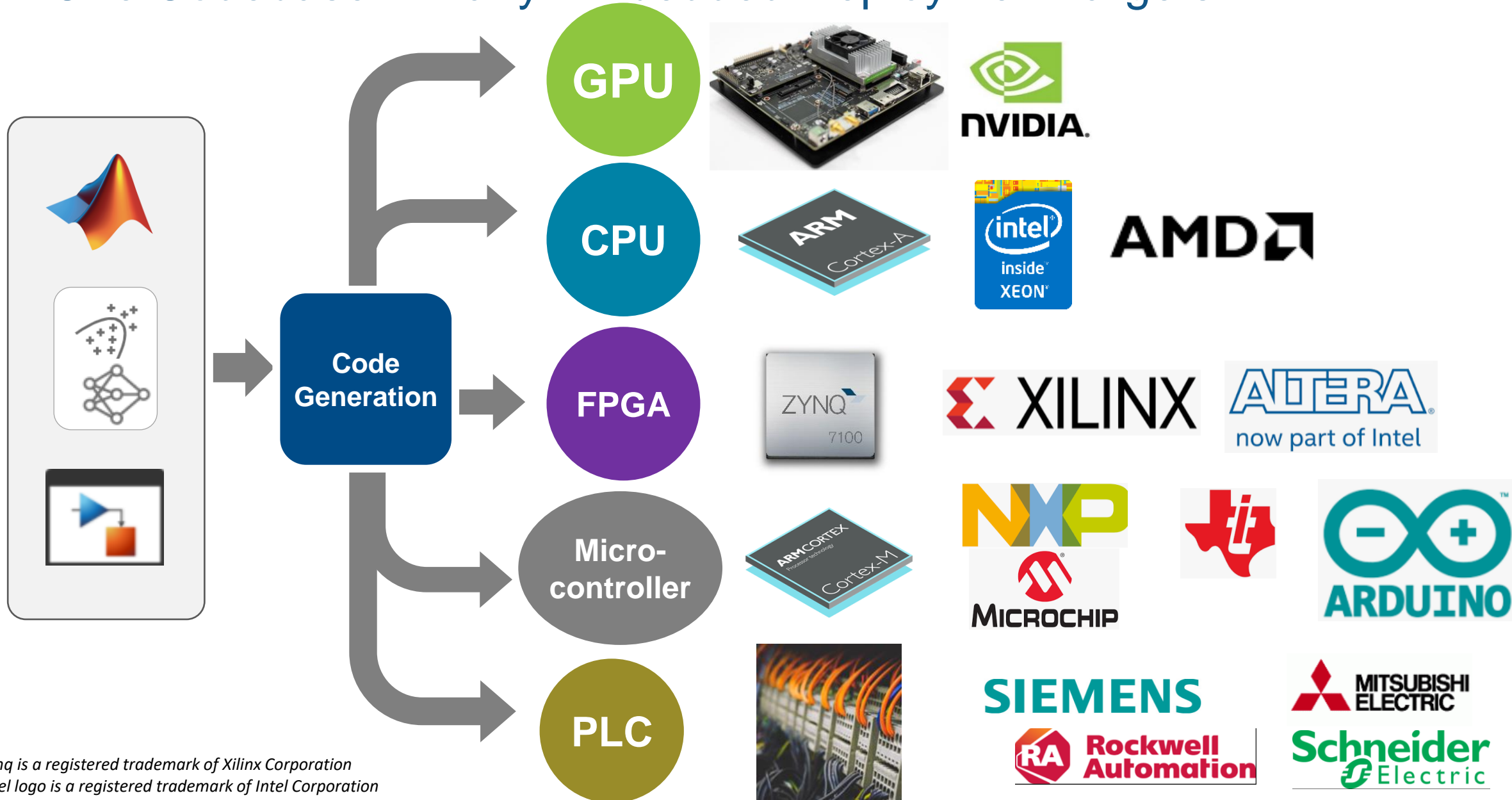
# **Deep Learning Demo** Size Reduction by factor 5

# One Codebase – Many Embedded Deployment targets

# One Codebase – Many Embedded Deployment targets
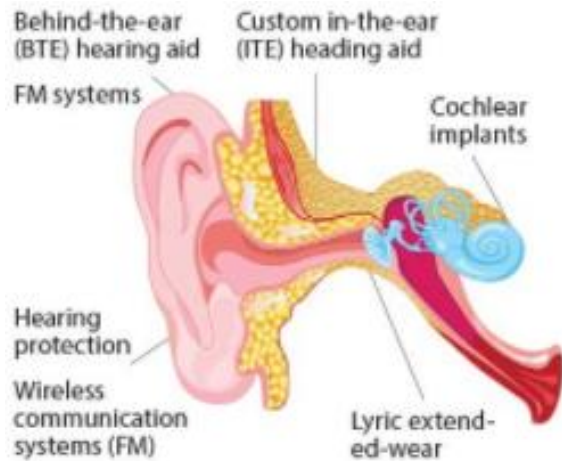
## Conclusions

You can fit AI for many applications onto limited hardware

MathWorks tools make fitting AI models on constrained hardware a lot easier

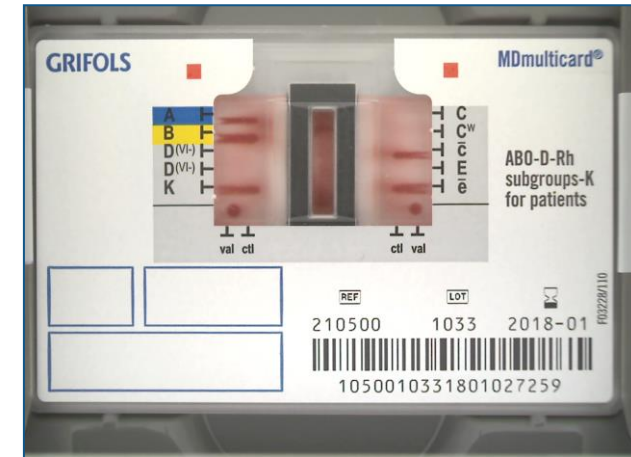Same high-level Workflow for any type of AI

# Learn More



**Hearing Implant using MBD**

**Autonomous Tractor**

**Card to Classify Blood Type**

To get your started:

Learn about Embedded Deployment

Quantization of classification SVM (Doc)

Deploy Hand-Gesture Classifier onto Arduino (Doc)

Generate C/C++ Code from Simulink (Video)

Quantizing a Deep Neural Network (Video)

# Thank you