# tinyML Talks

*Enabling Ultra-low Power Machine Learning at the Edge*

## "Deploying AI to Embedded Systems"

Bernhard Suhm – MathWorks

April 13, 2021

TINY
ML

www.tinyML.org

# tinyML Talks Sponsors

TINY ML TALKS webcast

**arm**
*tinyML Strategic Partner*

Deeplite

EDGE IMPULSE

maxim integrated™

Qeexo
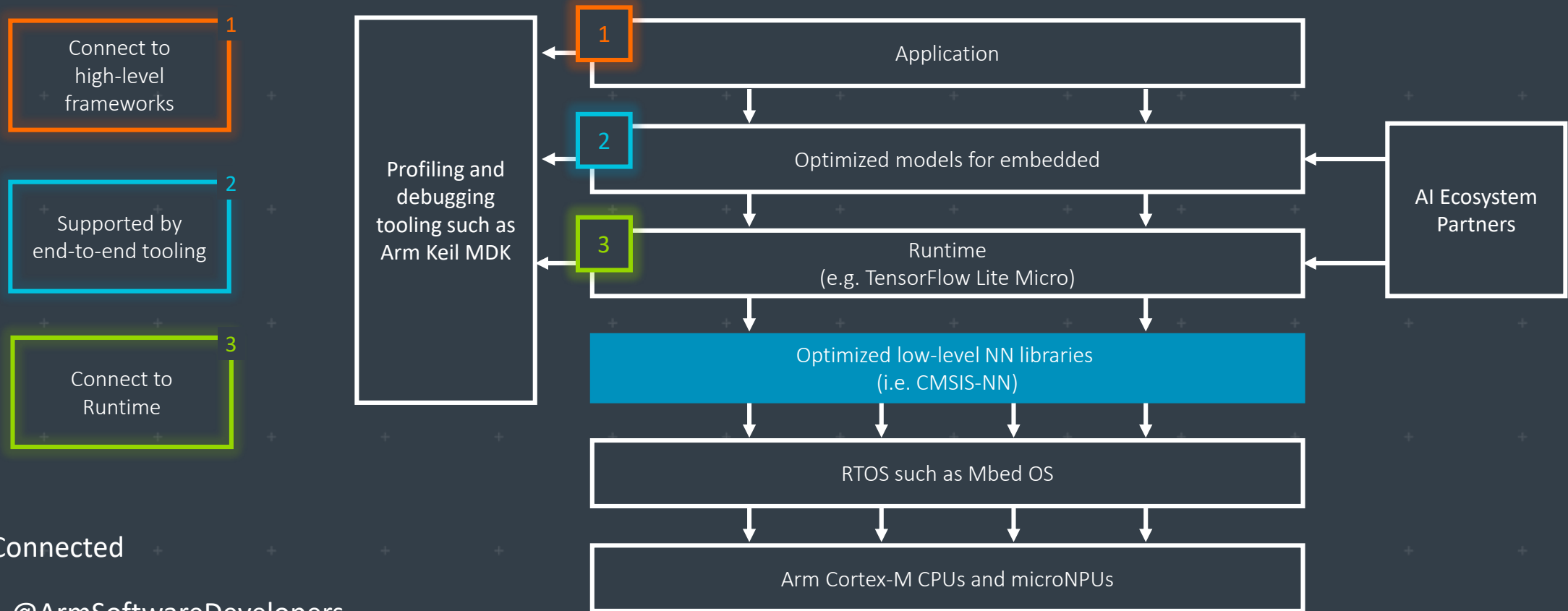
RealityAI

SynSense

Additional Sponsorships available – contact Olga@tinyML.org for info

# Arm: The Software and Hardware Foundation for tinyML

| | 1 |
|---|---|
| Connect to high-level frameworks | |

| | 2 |
|---|---|
| Supported by end-to-end tooling | |

| | 3 |
|---|---|
| Connect to Runtime | |

**Stay Connected**

▶ @ArmSoftwareDevelopers

🐦 @ArmSoftwareDev

Resources: developer.arm.com/solutions/machine-learning-on-arm

Profiling and debugging tooling such as Arm Keil MDK

**1** Application

**2** Optimized models for embedded

**3** Runtime (e.g. TensorFlow Lite Micro)

Optimized low-level NN libraries (i.e. CMSIS-NN)

RTOS such as Mbed OS

Arm Cortex-M CPUs and microNPUs

AI Ecosystem Partners

arm

# TinyML for all developers

**Dataset**

Acquire valuable training data securely

Enrich data and train ML algorithms

**Impulse**

Test impulse with real-time device data flows

**Test**

Embedded and edge compute deployment options

**Edge Device**

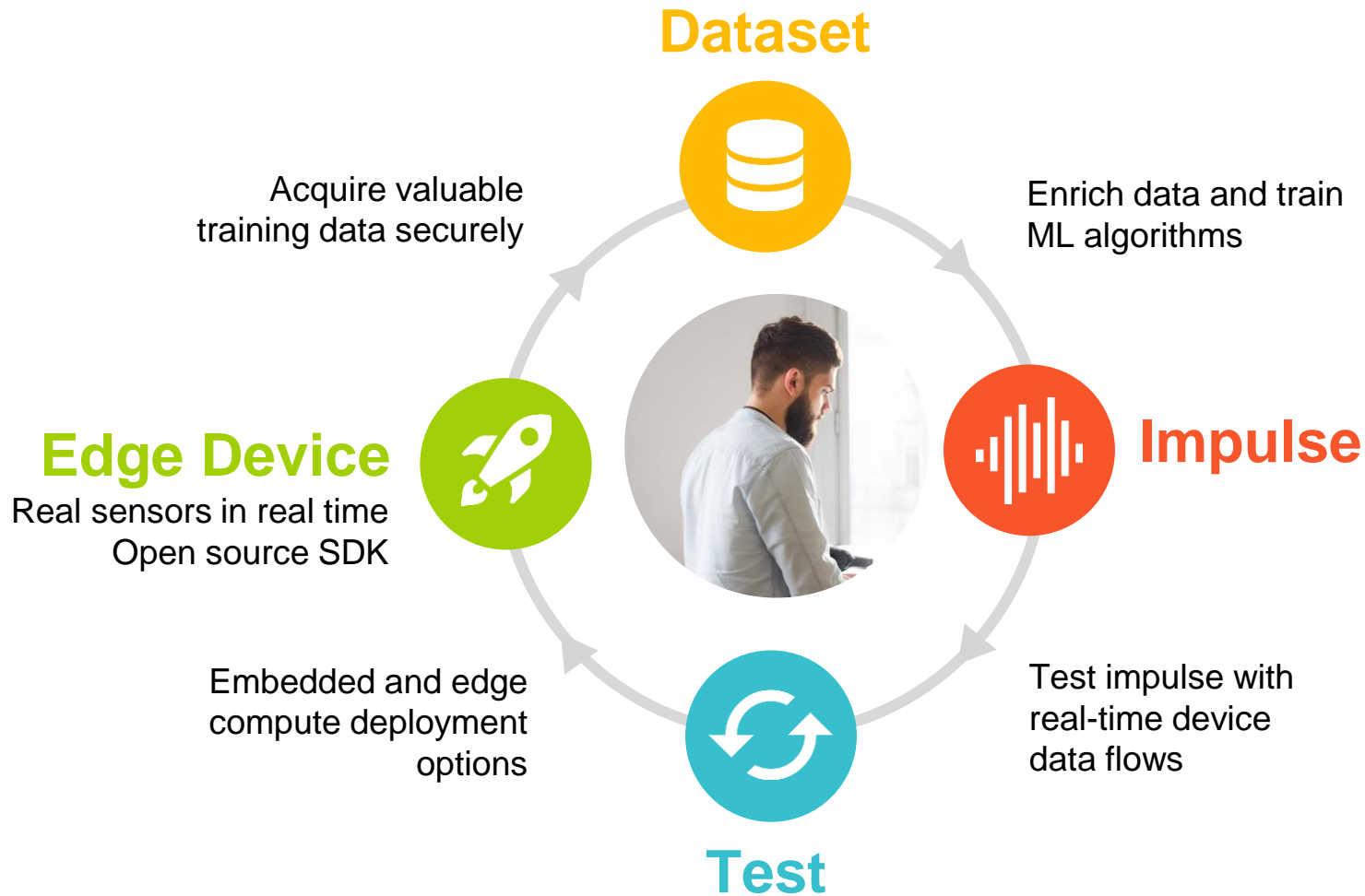Real sensors in real time
Open source SDK
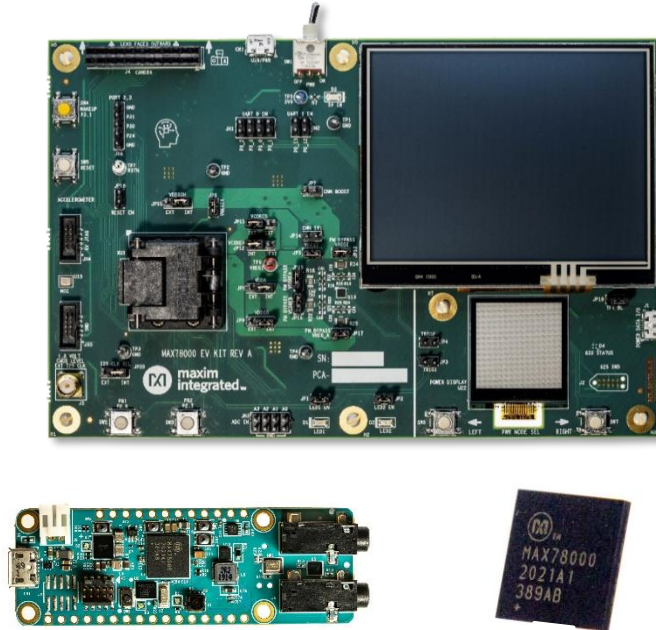
C++ library

Arduino library

WebAssembly

www.edgeimpulse.com
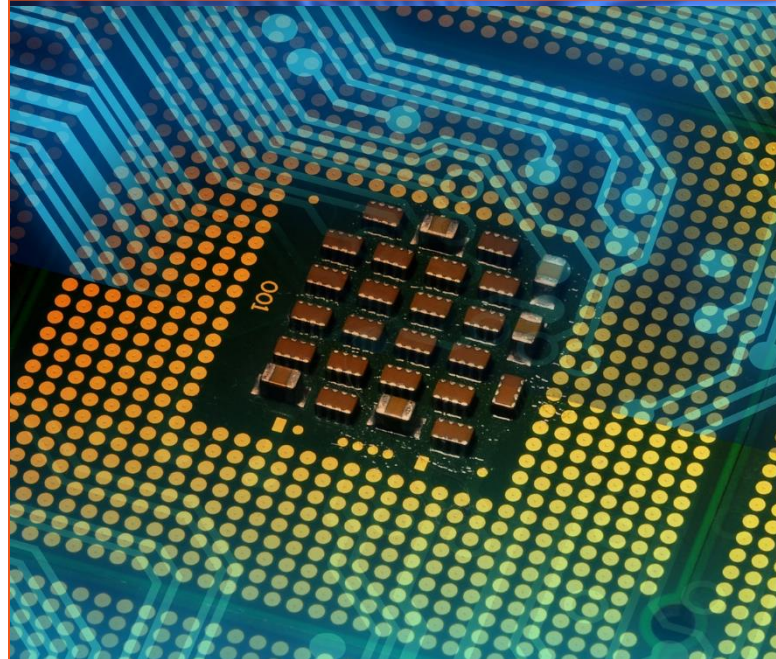
# Maxim Integrated: Enabling Edge Intelligence

## Advanced AI Acceleration IC



The new MAX78000 implements AI inferences at low energy levels, enabling complex audio and video inferencing to run on small batteries. Now the edge can see and hear like never before.

www.maximintegrated.com/MAX78000

## Low Power Cortex M4 Micros



Large (3MB flash + 1MB SRAM) and small (256KB flash + 96KB SRAM, 1.6mm x 1.6mm) Cortex M4 microcontrollers enable algorithms and neural networks to run at wearable power levels.

www.maximintegrated.com/microcontrollers

## Sensors and Signal Conditioning



Health sensors measure PPG and ECG signals critical to understanding vital signs. Signal chain products enable measuring even the most sensitive signals.

www.maximintegrated.com/sensors
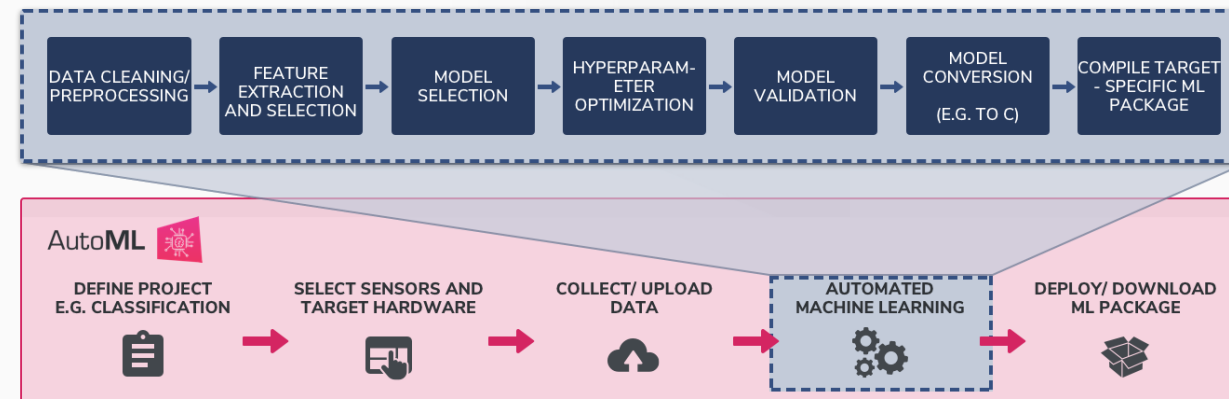
# Qeexo AutoML

Automated Machine Learning Platform that builds tinyML solutions for the Edge using sensor data

## Key Features

- Supports 17 ML methods:
  - Multi-class algorithms: GBM, XGBoost, Random Forest, Logistic Regression, Gaussian Naive Bayes, Decision Tree, Polynomial SVM, RBF SVM, SVM, CNN, RNN, CRNN, ANN
  - Single-class algorithms: Local Outlier Factor, One Class SVM, One Class Random Forest, Isolation Forest
- Labels, records, validates, and visualizes time-series sensor data
- On-device inference optimized for low latency, low power consumption, and small memory footprint applications
- Supports Arm® Cortex™- M0 to M4 class MCUs

## End-to-End Machine Learning Platform

DATA CLEANING/ PREPROCESSING → FEATURE EXTRACTION AND SELECTION → MODEL SELECTION → HYPERPARAM-ETER OPTIMIZATION → MODEL VALIDATION → MODEL CONVERSION (E.G. TO C) → COMPILE TARGET - SPECIFIC ML PACKAGE

AutoML

DEFINE PROJECT E.G. CLASSIFICATION → SELECT SENSORS AND TARGET HARDWARE → COLLECT/ UPLOAD DATA → AUTOMATED MACHINE LEARNING → DEPLOY/ DOWNLOAD ML PACKAGE

**For more information, visit: www.qeexo.com**

## Target Markets/Applications

- Industrial Predictive Maintenance
- Smart Home
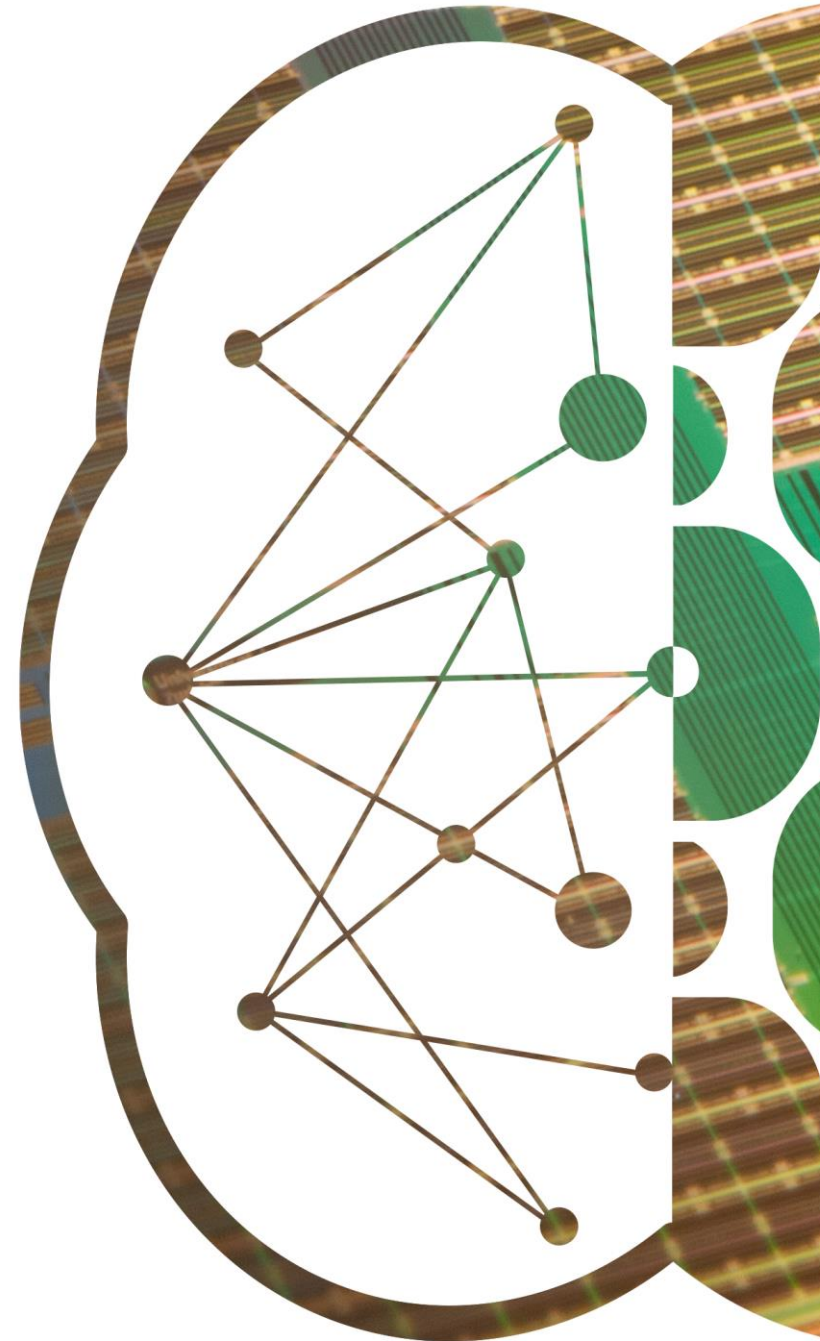- Wearables
- Automotive
- Mobile
- IoT

# SynSense

**SynSense** builds **sensing and inference** hardware for **ultra-low-power** (sub-mW) **embedded, mobile and edge** devices. We design systems for **real-time always-on smart sensing**, for audio, vision, IMUs, bio-signals and more.
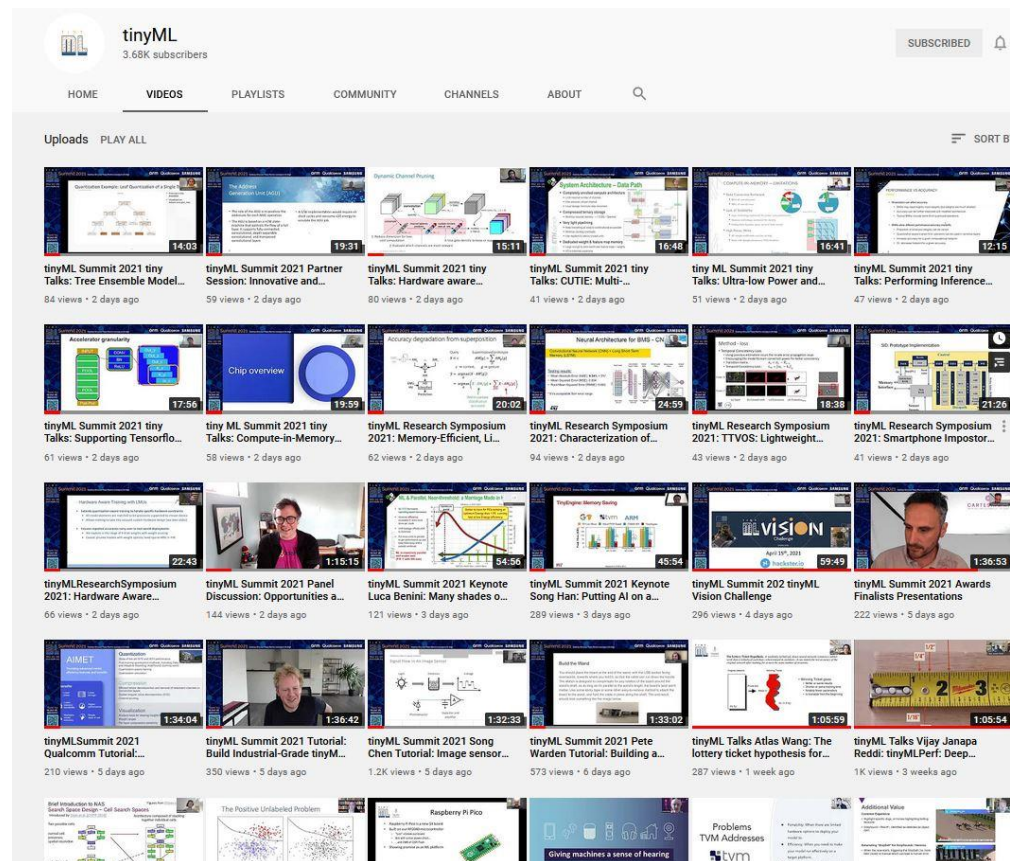
https://SynSense.ai

# Successful tinyML Summit 2021:

- **5** days of tutorials, talks, panels, breakouts, symposium
  - **4** tutorials
  - **6** keynotes & **6** plenary tinyTalks (more in breakouts)
  - **2** panel discussions
  - **5** disruptive news presentations
  - **17** breakout/partner sessions
  - **6** Best Product and Innovation Award Finalists & Presentations
  - **89** Speakers

- **5006** registered attendees representing:
  - **104** countries, **1000+** companies and **400+** academic institutions

- **26** Sponsoring companies



[www.youtube.com/tinyML](www.youtube.com/tinyML) with 150+ videos

tinyML Summit-2022, January 24-26, Silicon Valley, CA

June 7-10, 2021 (virtual, but LIVE)

Deadline for abstracts: May 1

tinyML EMEA Technical Forum 2021

Enabling ultra-Low Power Machine Learning at the Edge

June 7-10, 2021

Inaugural tinyML EMEA Technical Forum

Venue

Virtual - online

Sponsorships are being accepted: sponsorships@tinyML.org

April 15th, 2021 Launch

# Next tinyML Talks

| Date | Presenter | Topic / Title |
|------|-----------|---------------|
| Tuesday, April 27 | **Michael Jo and Xingheng Lin** Rose-Hulman Institute of Technology | Train-by-weight (TBW): Accelerated Deep Learning by Data Dimensionality Reduction |

Webcast start time is 8 am Pacific time

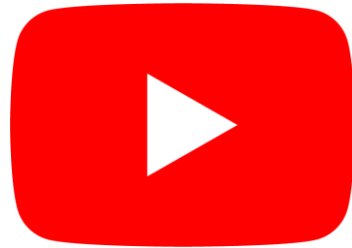Please contact talks@tinyml.org if you are interested in presenting
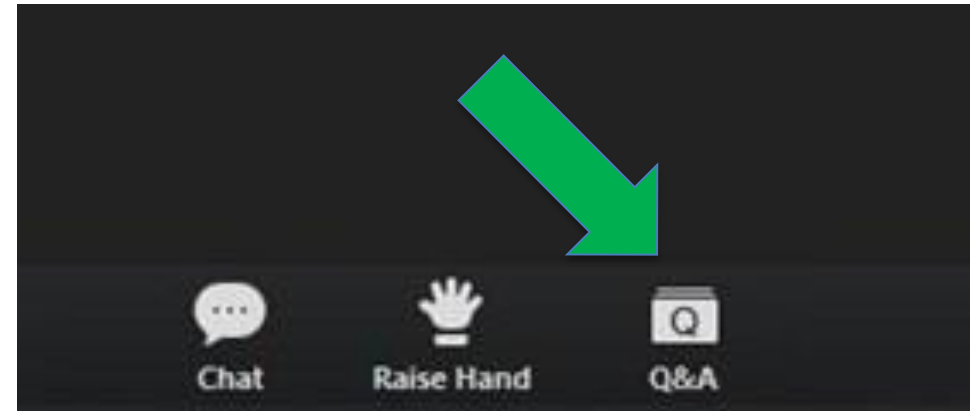
# Reminders

Slides & Videos will be posted tomorrow

Please use the Q&A window for your questions
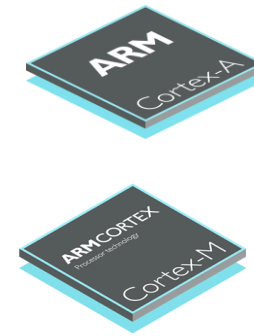
tinyml.org/forums     youtube.com/tinyml

# Bernhard Suhm

Bernhard Suhm is the product manager for Machine Learning at MathWorks. He works closely with customer facing and development teams to address customer needs and market trends in our machine learning related products, primarily the Statistics and Machine Learning toolbox. Prior to joining MathWorks Bernhard led a team of analysts consulting call centers on optimizing the delivery of customer service. He also held positions at a usability consulting company and Carnegie Mellon University. He received a PhD in Computer Science specializing in speech user interfaces from Karlsruhe University in Germany.
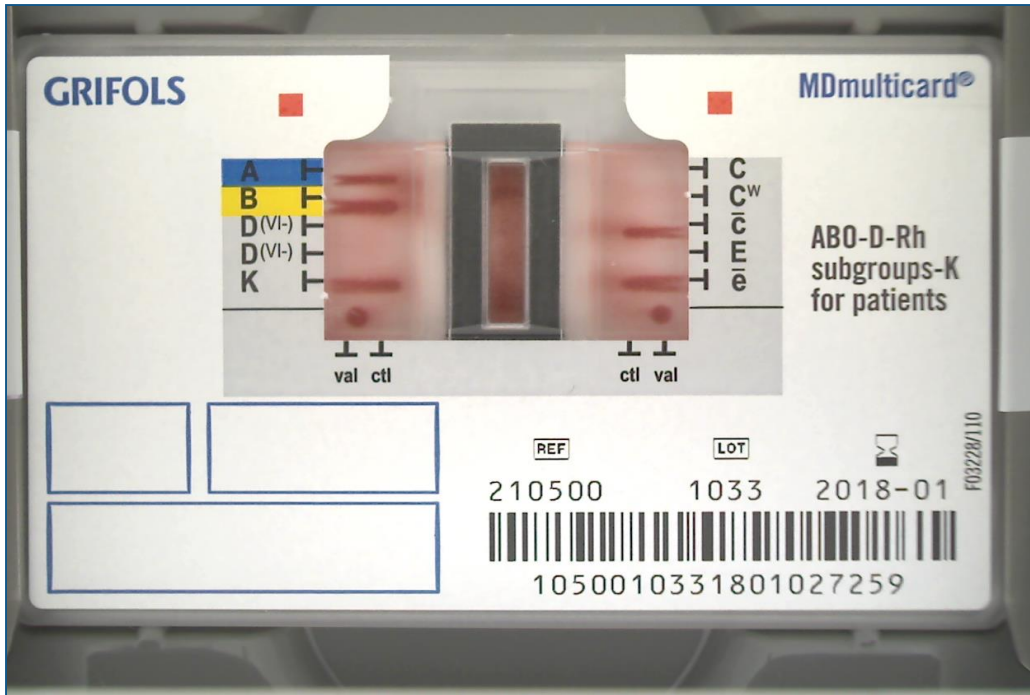
# Deploying AI to Embedded Systems

Bernhard Suhm

Product Manager – Machine Learning

# Examples for Embedded Deployment



**Card to Classify Blood Type**
**IDNEO**



**Oversteering Detection**

# Agenda

Deploying AI is difficult

Four specific challenges:

1. Integrate AI model with an embedded (or industrial) system
2. Fit large AI models on limited hardware
3. Test & Verify before deployment to production
4. Ongoing changes in environment or system behavior

is a **Leader** in the Gartner Magic Quadrant for 2021 Data Science and Machine Learning Platforms for the Second Year in a Row

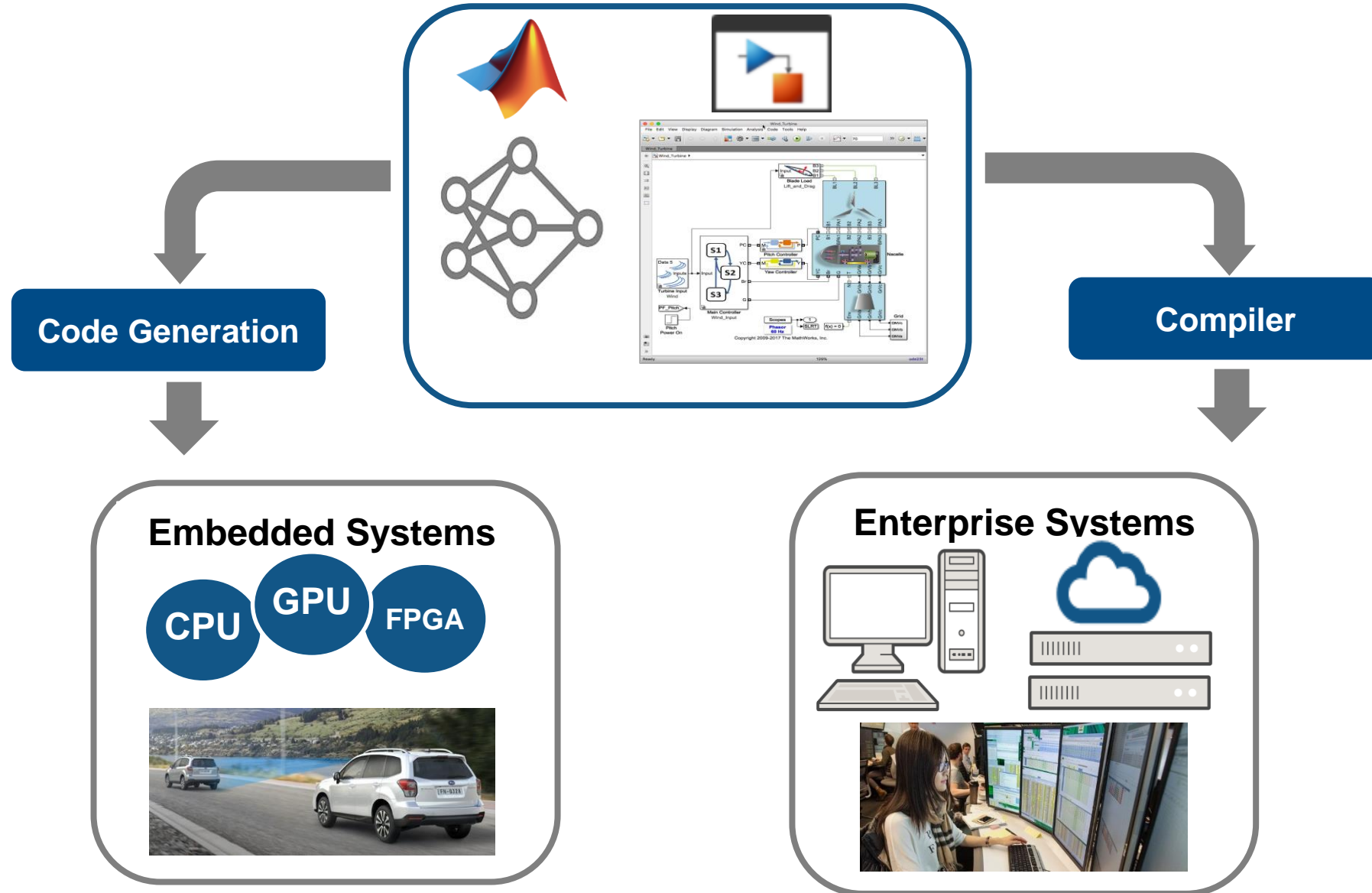Figure 1: Magic Quadrant for Data Science and Machine Learning Platforms
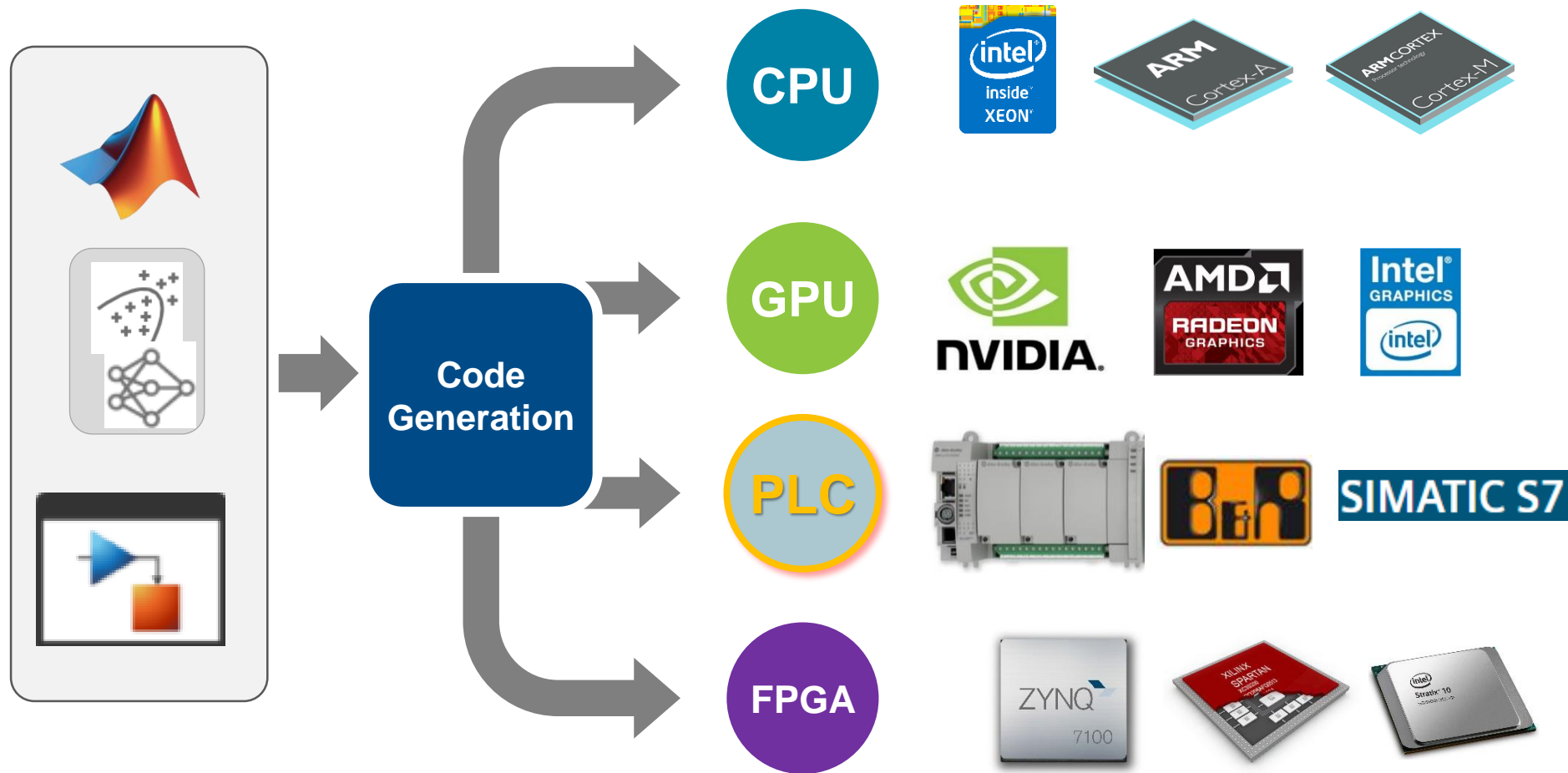
Source: Gartner (March 2021)

Gartner Magic Quadrant for Data Science and Machine Learning Platforms, Peter Krensky, Erick Brethenoux, Pieter denHamer, Farhan Choudhary, Afraz Jaffri, Subhangi Vashisth, 1st March 2021.

# Two Approaches for integrating AI with Larger System



**Code Generation**

**Compiler**

**Embedded Systems**

CPU  GPU  FPGA

**Enterprise Systems**

# One Codebase – Many Deployment targets

*NVIDIA, the NVIDIA logo, and TensorRT are registered trademarks of NVIDIA Corporation*
*Intel logo is a registrered trademark of Intel Coroporation*

# Human Activity Classification using Smartphones

Data:

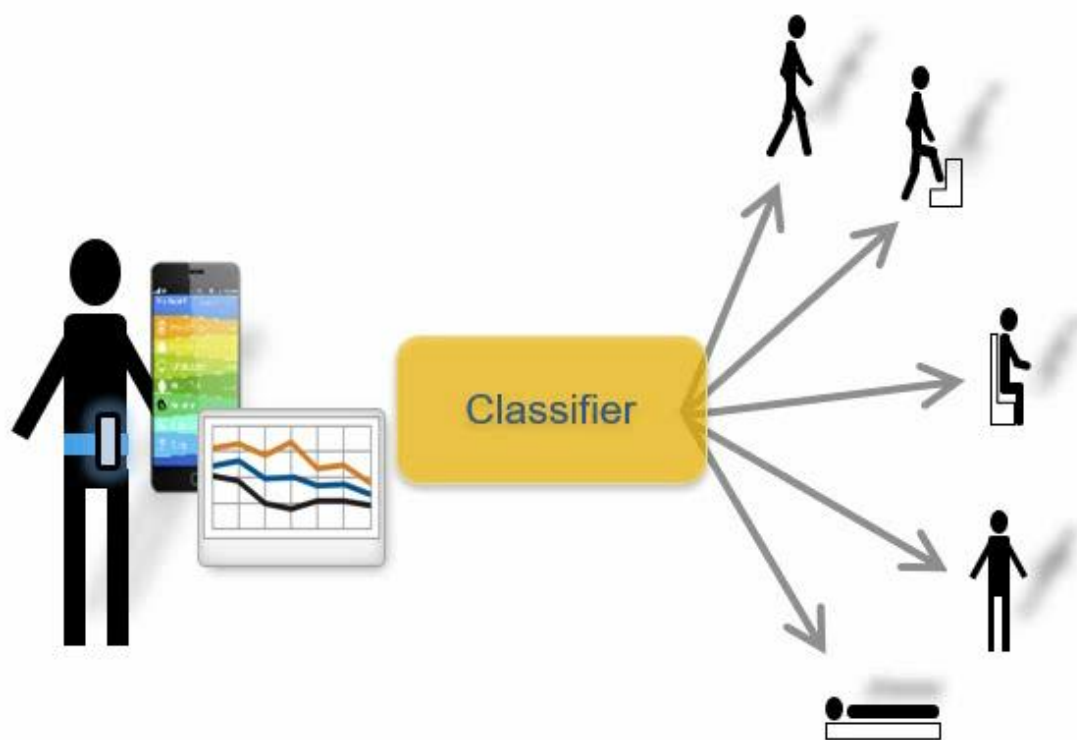Accelerometer from mobile

~7K observations from 30 subjects

Steps:

– Interactively build classifiers

– Integrate model with system

– Generate inference C-code

– Deploy to Android using Simulink

Classifier

**Walking**

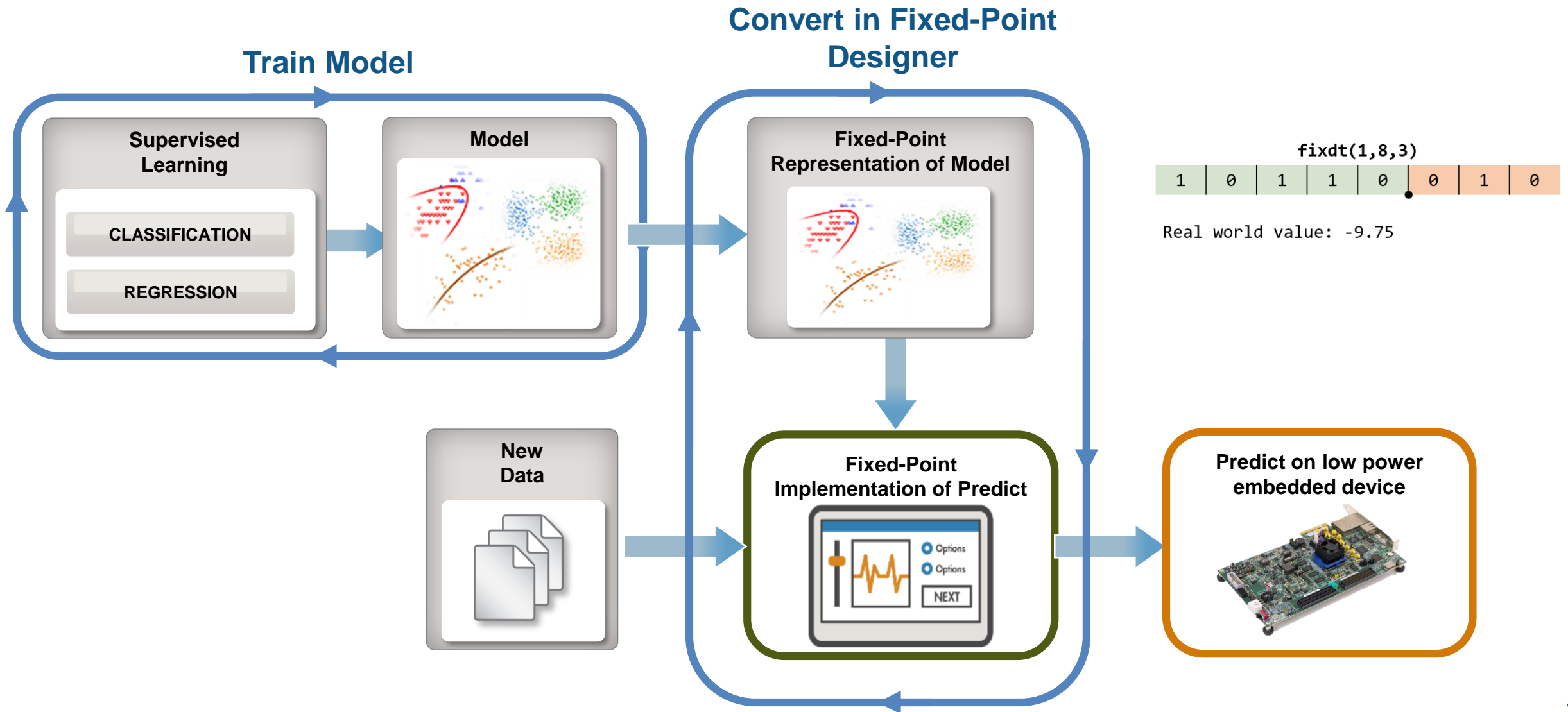**Idle**

# Human Activity Learning using Smartphones



## Human Activity Recognition Data

- Accelerometer from mobile
- 7K observations
- 66 manual features, reduced to 2
- 6 classes reduced to "Walking" vs "Idle"
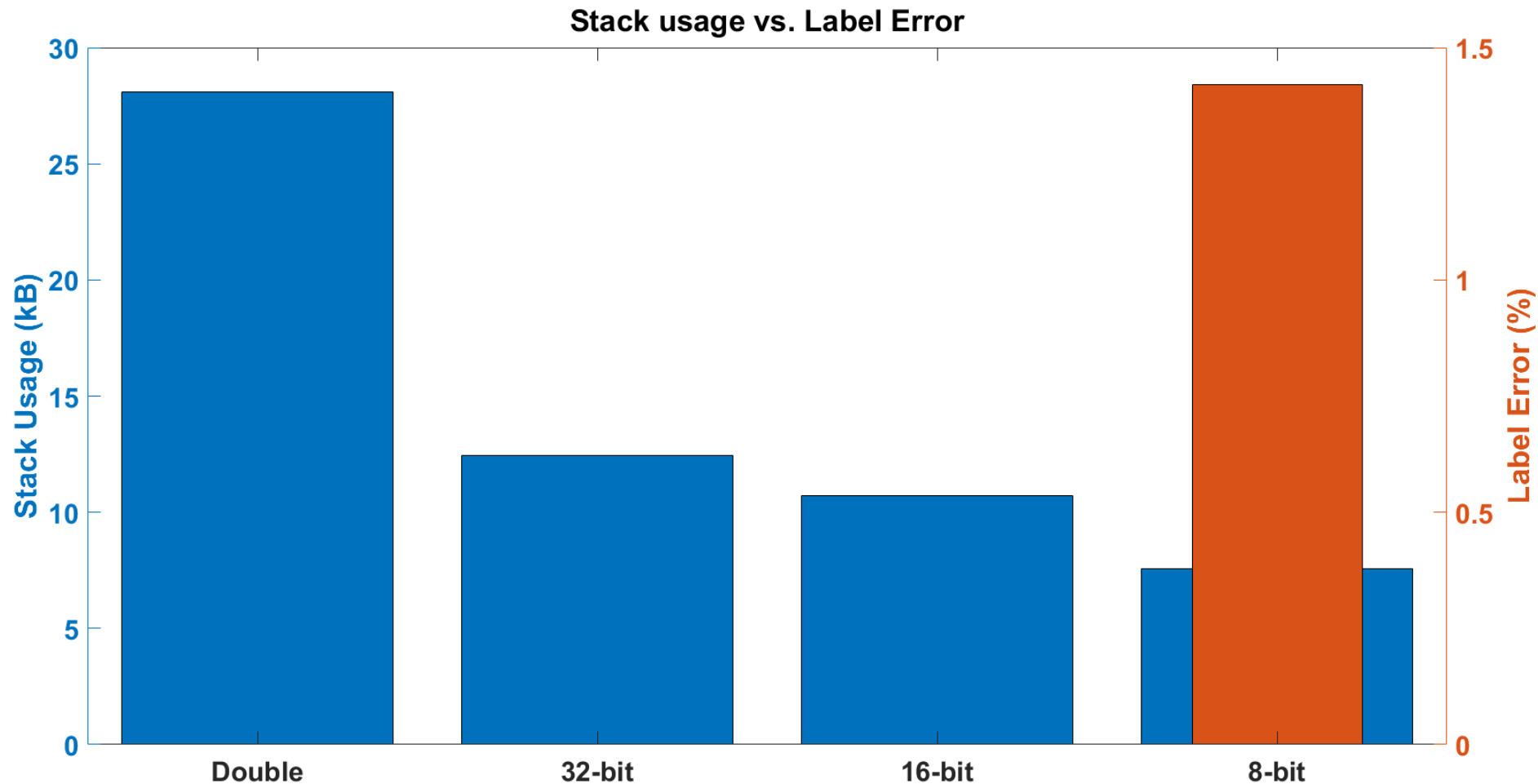
**Dataset courtesy of**:
Davide Anguita, Alessandro Ghio, Luca Oneto, Xavier Parra and Jorge L. Reyes-Ortiz. *Human Activity Recognition on Smartphones using a Multiclass Hardware-Friendly Support Vector Machine*. International Workshop of Ambient Assisted Living (IWAAL 2012). Vitoria-Gasteiz, Spain. Dec 2012 http://archive.ics.uci.edu/ml/datasets/Human+Activity+Recognition+Using+Smartphones

# Convert machine learning to fixed-point to reduce memory footprint

**Train Model**

**Convert in Fixed-Point Designer**

**Supervised Learning**

CLASSIFICATION

REGRESSION

**Model**

**Fixed-Point Representation of Model**

`fixdt(1,8,3)`

| 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 |
|---|---|---|---|---|---|---|---|

Real world value: -9.75

**New Data**

**Fixed-Point Implementation of Predict**

Options
Options
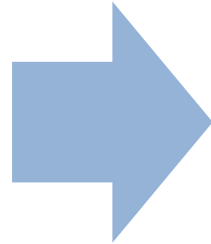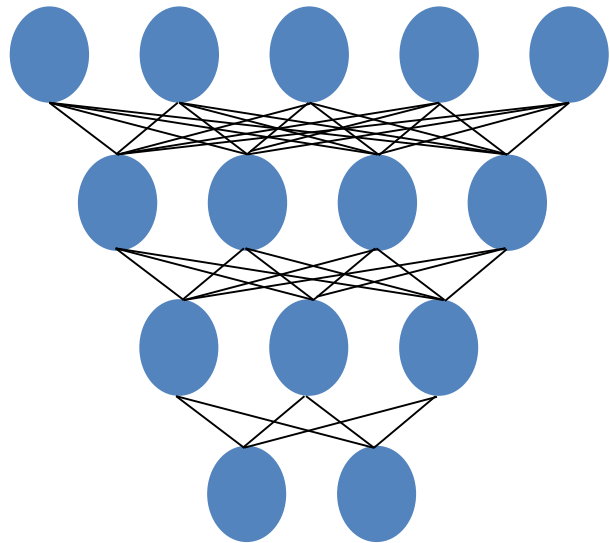NEXT

**Predict on low power embedded device**

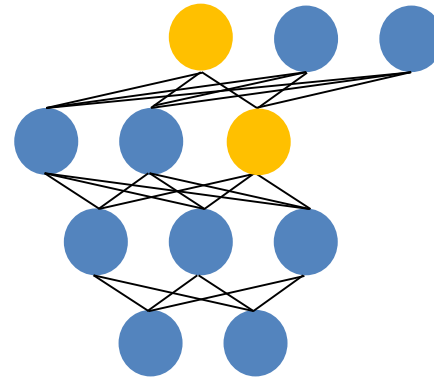# Fixed-point conversion is a trade-off between resource usage optimization and accuracy

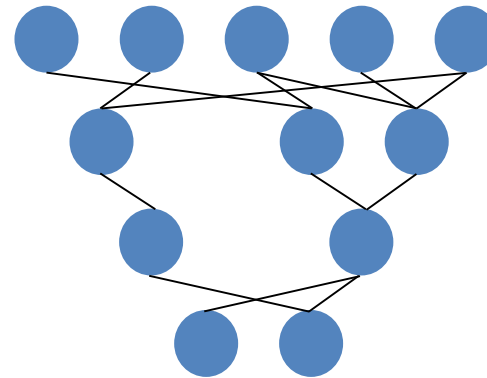# Three Approaches to reducing size of Deep Neural Nets
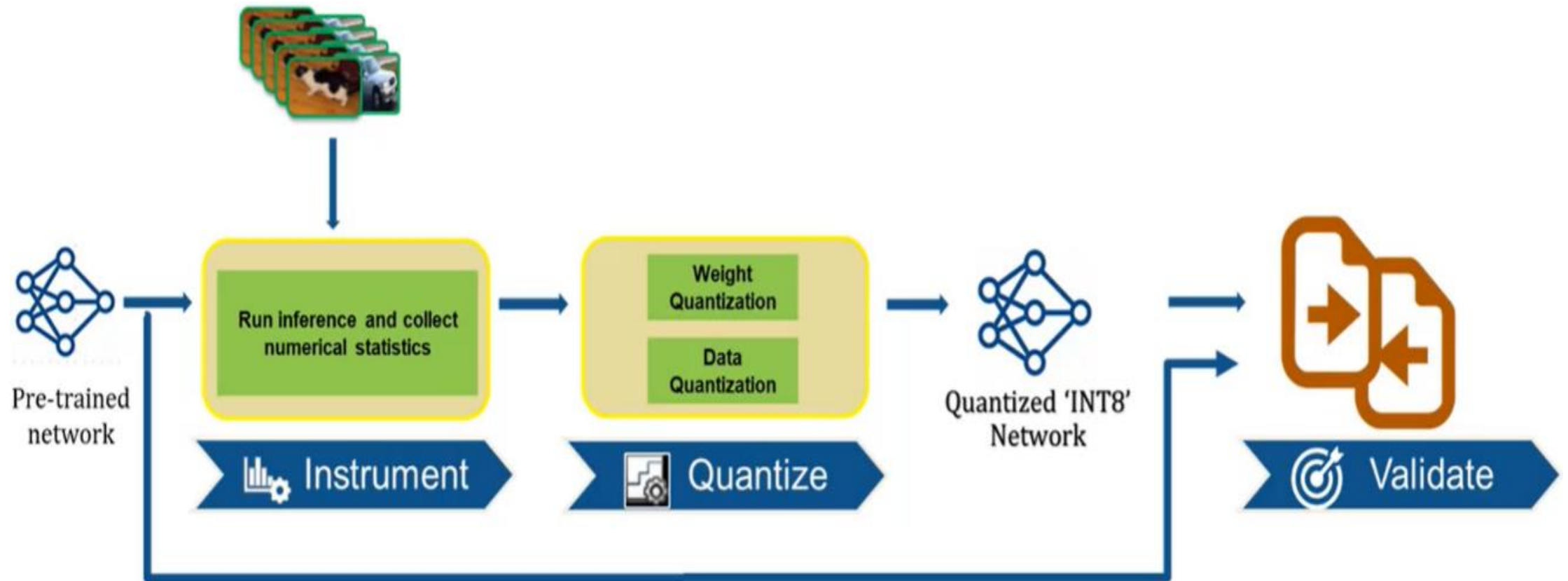
Original Deep Neural Net

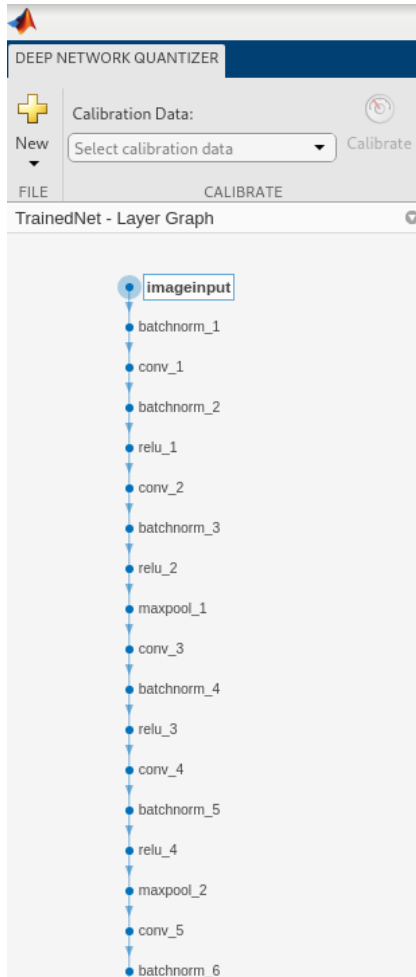1. Layer Fusion

2. Pruning

3. Quantization

# Model Quantization Workflow

# Deep Learning Quantization



**Use Deep Network Quantizer to Optimize the Inference Network**
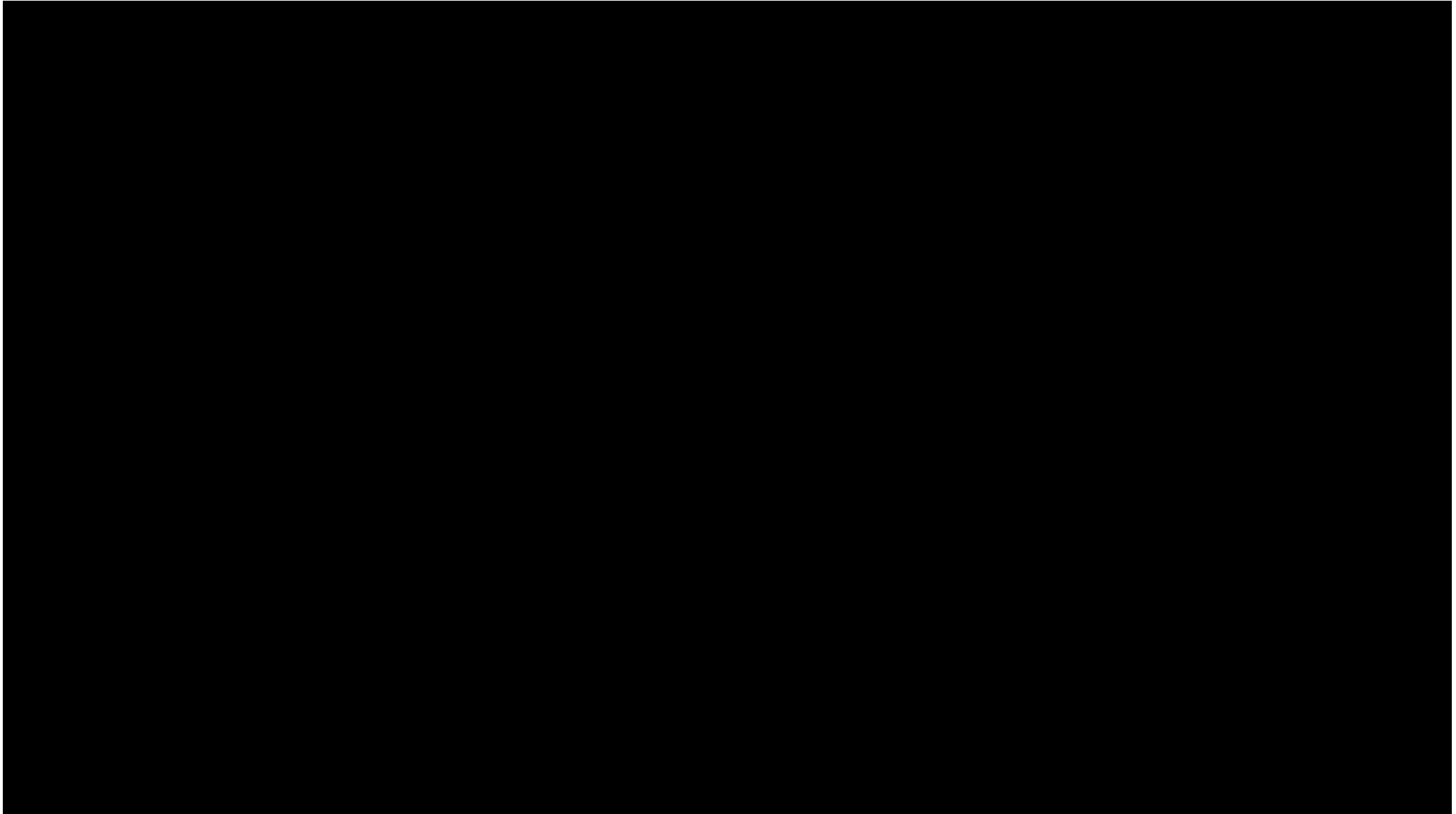
```
1   load('trainedNet');
2   analyzeNetwork(trainedNet);
3   numData = size(xTrain);
4   numData = numData(end);
5   augImds = augmentedImageDatastore(trainedNet.Layers(1).InputSize, xTrain, yTrain);
6   calDS = augImds.subset(1:floor(numData * 0.8));
7   valDS = augImds.subset(floor(numData * 0.8)+1:numData);
8   dq = dlquantizer(trainedNet, 'ExecutionEnvironment', 'GPU');
9   dq.calibrate(calDS)
```

- Load trained network
- Split data: calibration – 80%, validation – 20%
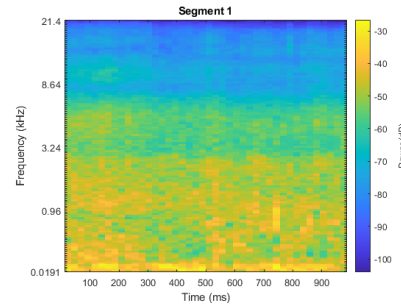- Launch Deep Network Quantizer App

*Video (start ~1:10)*
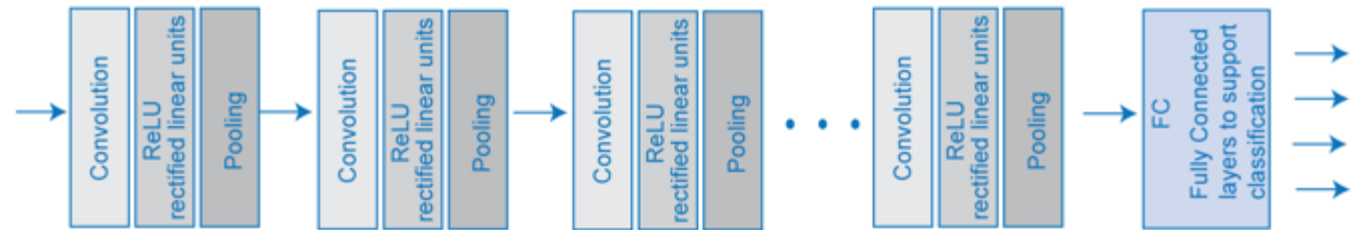
# Quiz: Which Sounds do you hear?

# Embedded Deployment of Acoustic Scene Recognition

Reformat the data

Convolutional Neural Networks (**CNN**)

SqueezeNet  ~5MB
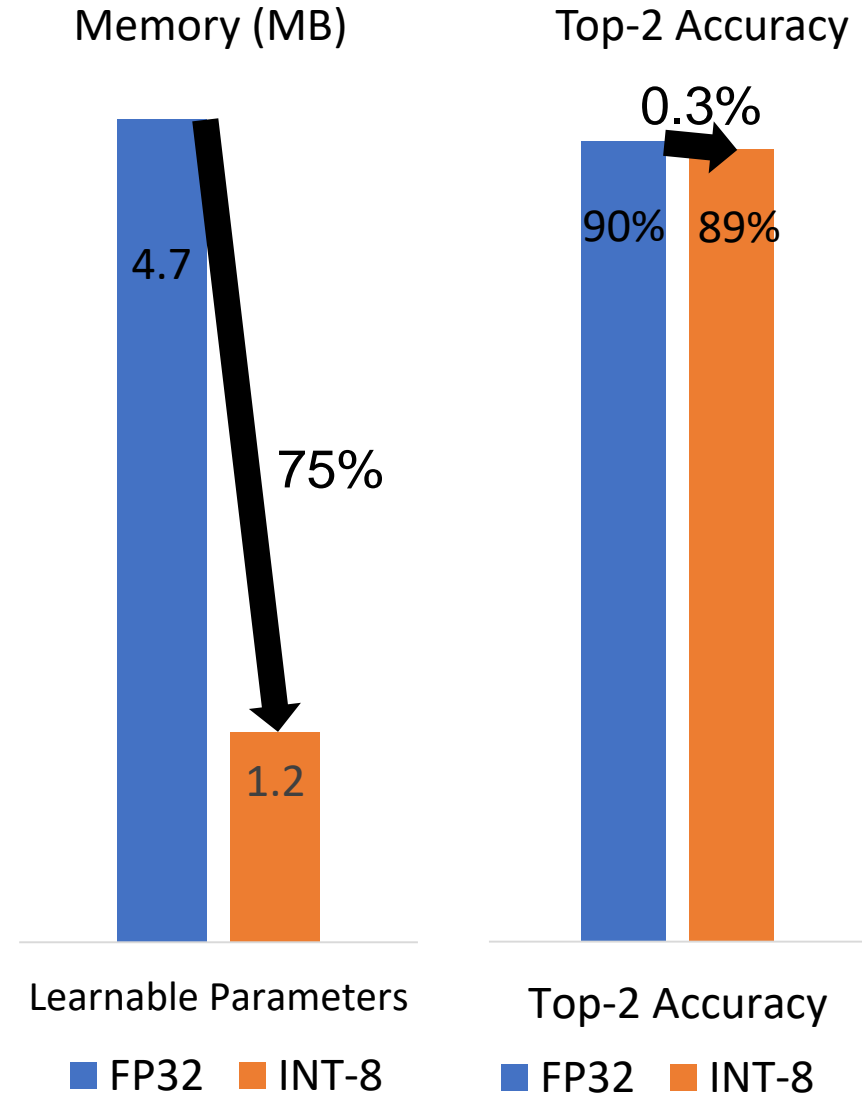ResNet-50  ~100MB

Limited resources

# 4x reduction in Memory – less than 5% increase in error rate

Validation Summary

✔ Validation Results

Number of samples: 2,280

| Metric | Floating-Point Network Results | Quantized Network Results | Perc |
|---|---|---|---|
| Learnable parameter memory (MB) | 4.7033 | 1.1903 | 74.6 |
| hComputeAccuracy | 0.9022 | 0.8991 | 0.34 |

**Memory (MB)**

4.7

75%

1.2

Learnable Parameters

■ FP32  ■ INT-8

**Top-2 Accuracy**
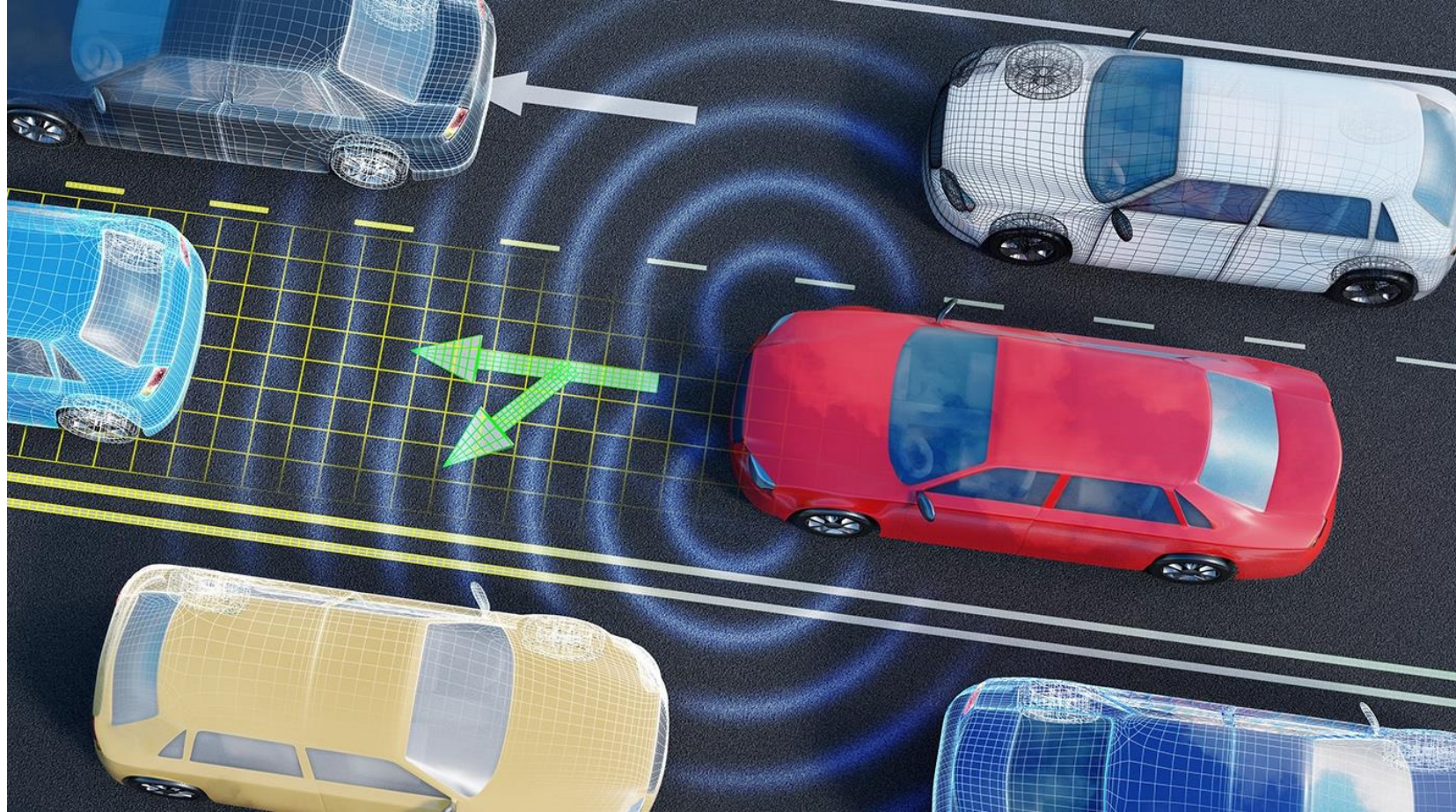
0.3%

90%  89%

Top-2 Accuracy

■ FP32  ■ INT-8

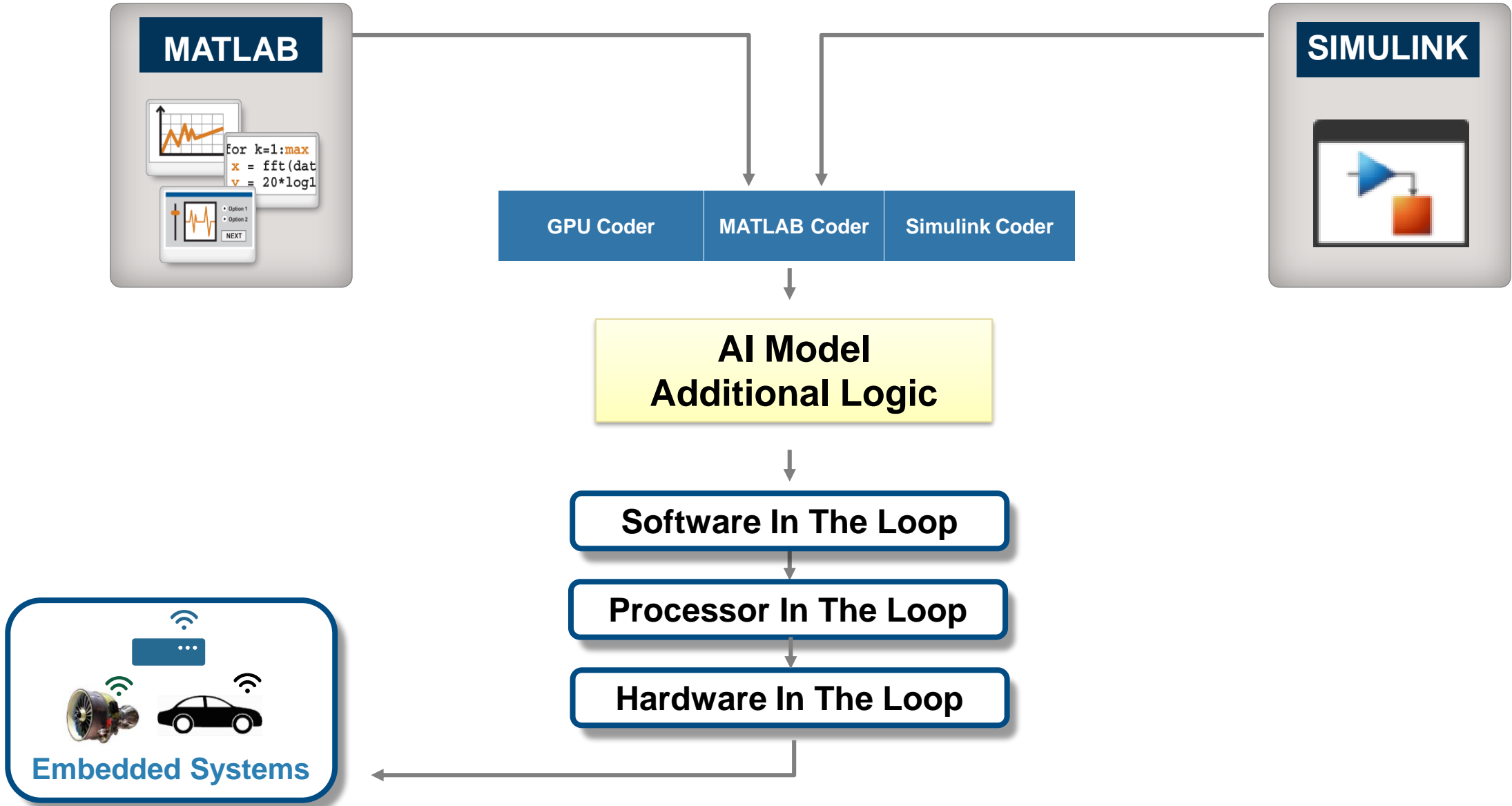# Agenda

Deploying AI is difficult

Four specific challenges:

✓ Integrate AI model with an industrial / embedded system

✓ Fit large AI models on limited hardware

3. Test & Verify before deployment on hardware

4. Ongoing changes in environment or system behavior

# Integrate AI into system-wide context, simulate before moving to hardware, and verify effectiveness

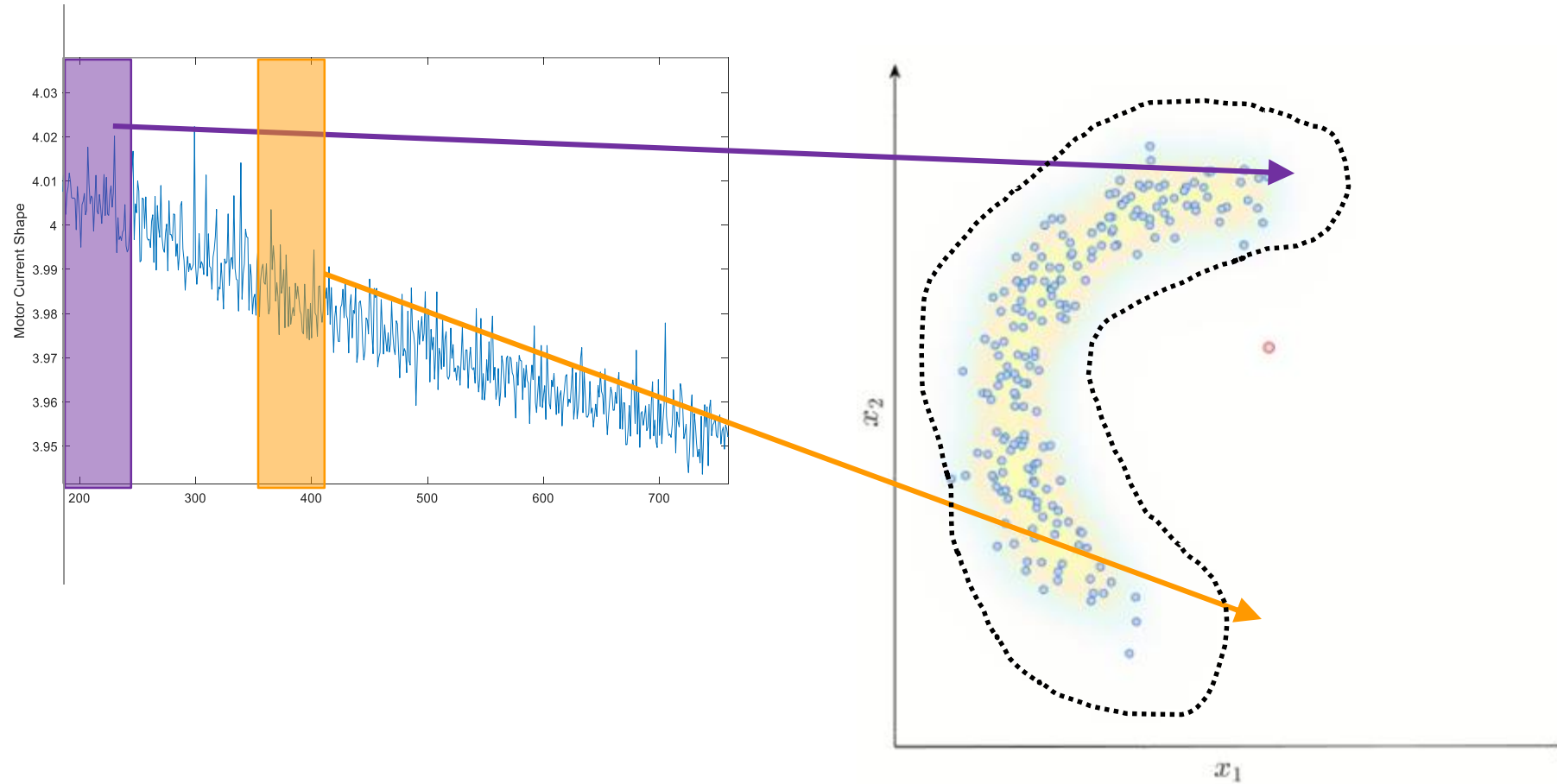# Facilitate Verification & Validation of your AI Application

# Agenda

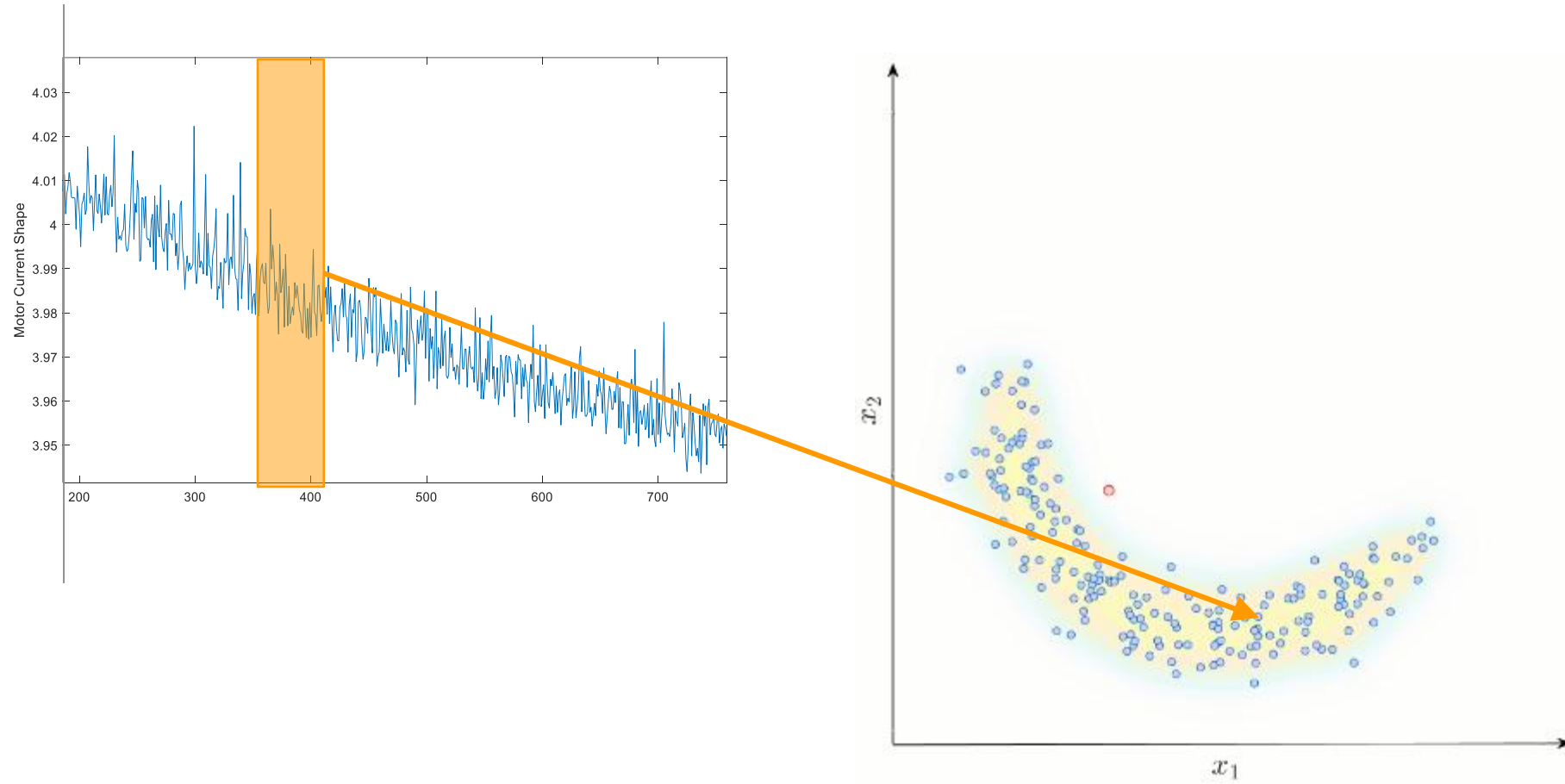Deploying AI is difficult

Three specific challenges:

✓ Integrate AI model with an industrial / embedded system

✓ Fit large AI models on limited hardware

✓ Test & Verify before deployment on hardware

4. Ongoing changes in environment or system behavior

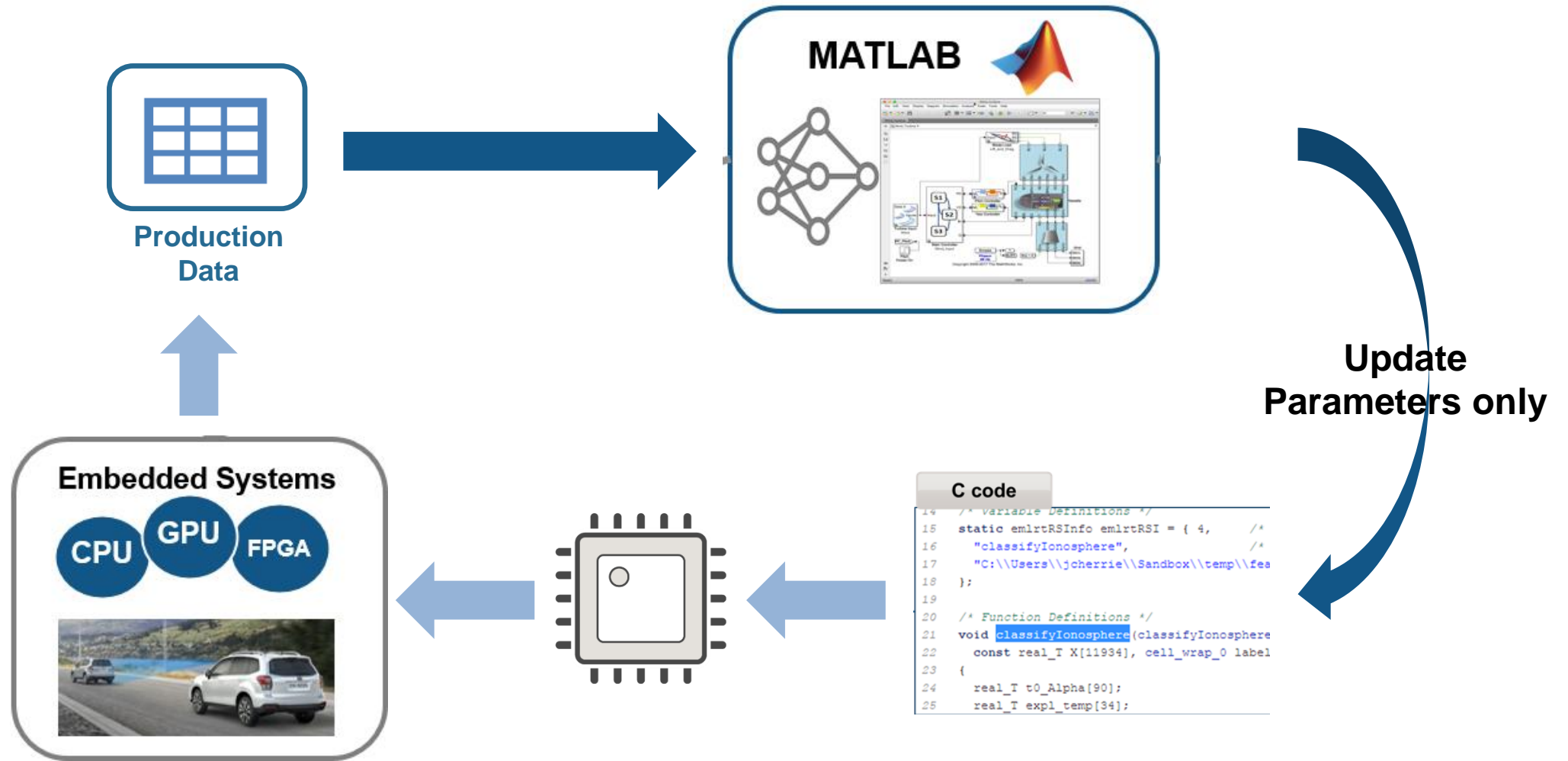# AI models reflect System behaviors and Environment



*(illustration only; not based on actual data)*
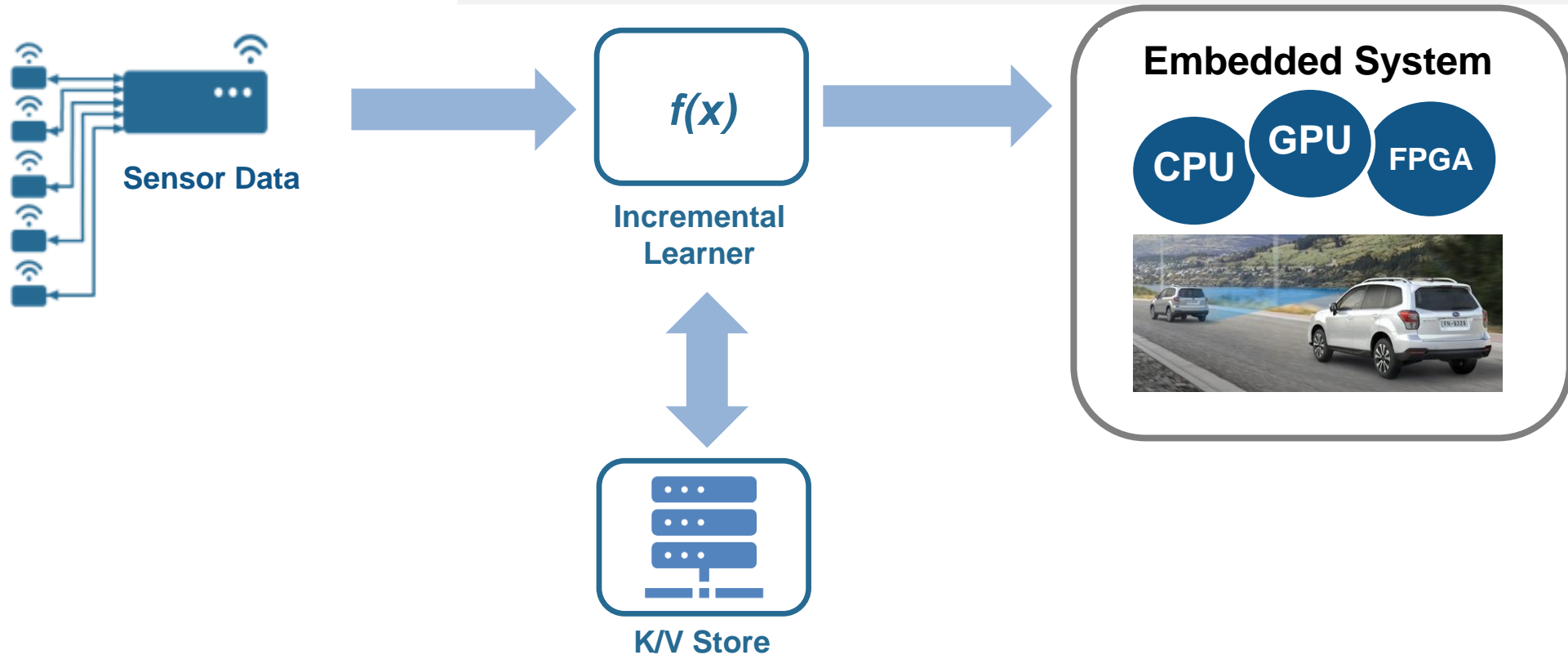
# Deployed Models Need to Adapt.

# "In-place" Model Updates



**Production Data**

MATLAB

**Update Parameters only**

**Embedded Systems**

CPU   GPU   FPGA

C code

```
14  /* Variable Definitions */
15  static emlrtRSInfo emlrtRSI = { 4,        /*
16    "classifyIonosphere",                   /*
17    "C:\\Users\\jcherrie\\Sandbox\\temp\\fea
18  };
19
20  /* Function Definitions */
21  void classifyIonosphere(classifyIonosphere
22    const real_T X[11934], cell_wrap_0 label
23  {
24    real_T t0_Alpha[90];
25    real_T expl_temp[34];
```

# Embedded Deployment is about finding the right trade-off

✓ Memory footprint

Energy Cost
- Computation
- Data Access

Computation speed

*Rough energy costs for various operations in 45nm 0.9V*

| Operation | Energy (pJ) |
|-----------|-------------|
| 8b Add | 0.03 |
| 32b FP Add | 0.9 |
| 8b Multiply | 0.2 |
| 32b FP Mult | 3.7 |

*Quantization greatly reduces energy consumption*

| | %Reduction |
|---|-----------|
| 32b FP to 8b Add | 97% |
| 32b FP to 8b Multiply | 95% |

*Source: Mark Horowitz "Computing's Energy Problem (and what we can do about it)", ISSCC 2014*

# Conclusions

Deploy to many targets from one codebase

Tools for handling these challenges in production deployments:

- Fit models to embedded hardware with Quantization / Fixed-Point conversion
- Simulink facilitates system integration, verification and testing
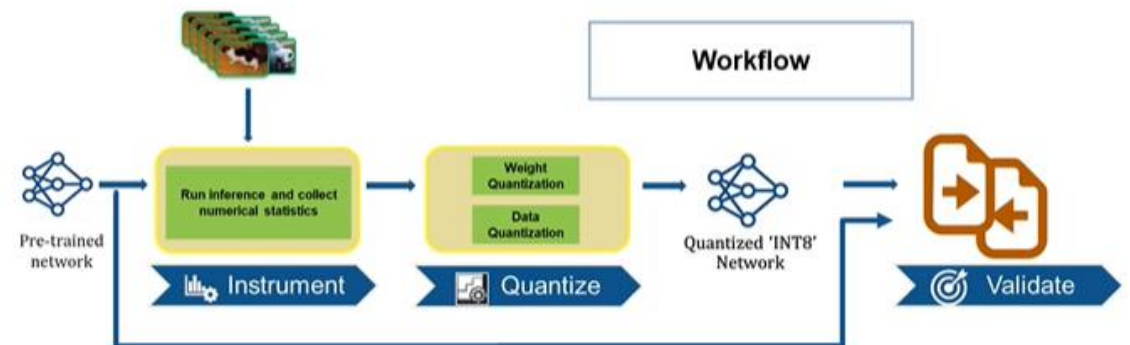- Incrementally adapt deployed models

# Learn More

## Machine Learning

- [Fixed-Point Prediction with SVM](#)
- [Instrument MATLAB Code](#)
- [Update Model Parameters for Code Generation](#)
- [Incremental Learning with Logistic Regression](#)
- [Hand Gesture recognition on Arduino Nano](#)

## Deep Learning

- [Quantize Residual Network Trained for Image Classification](#)
- [Video walk through Deep Network Quantizer](#)



Deep Learning Model Quantization in MATLAB

# Copyright Notice

# www.tinyML.org