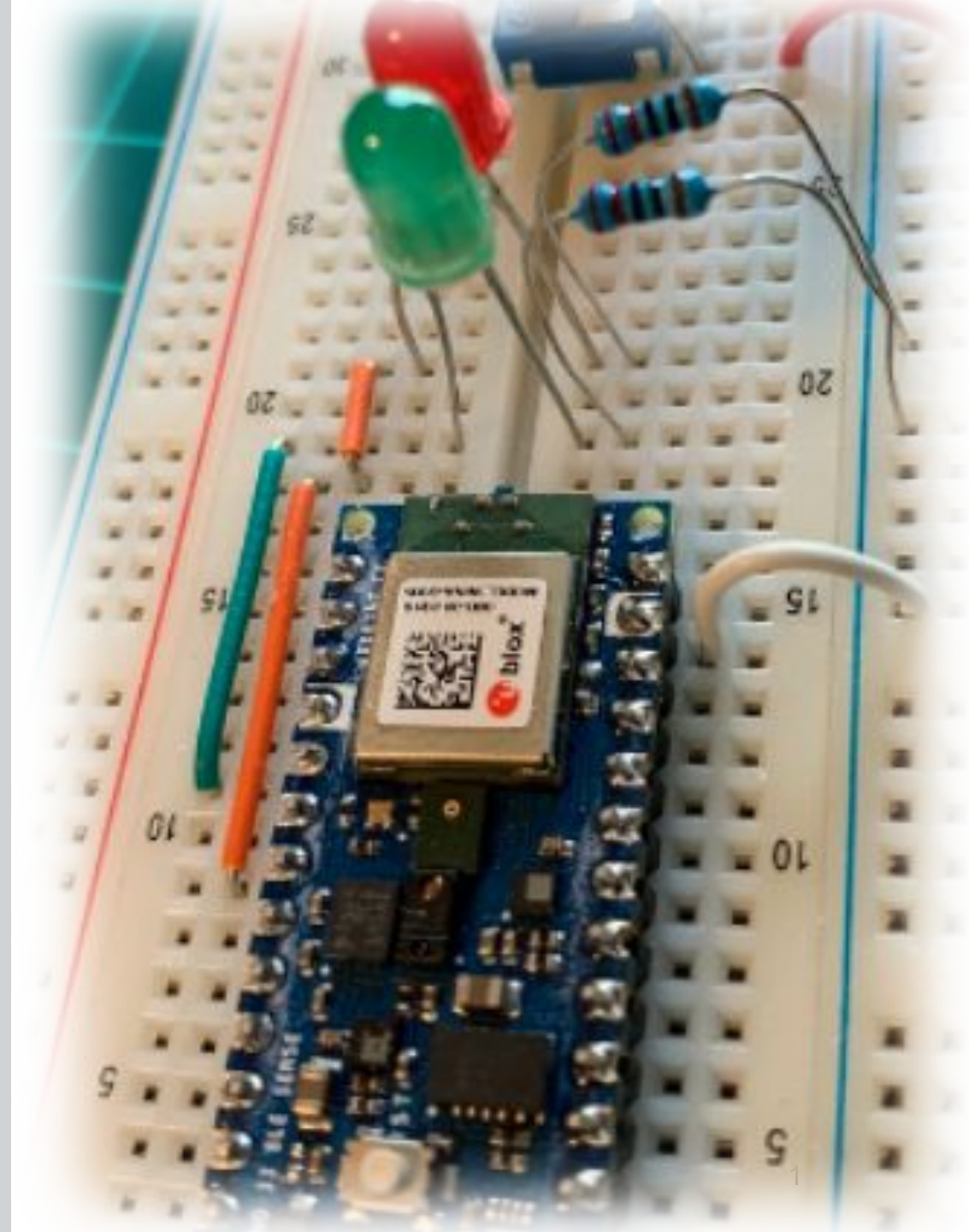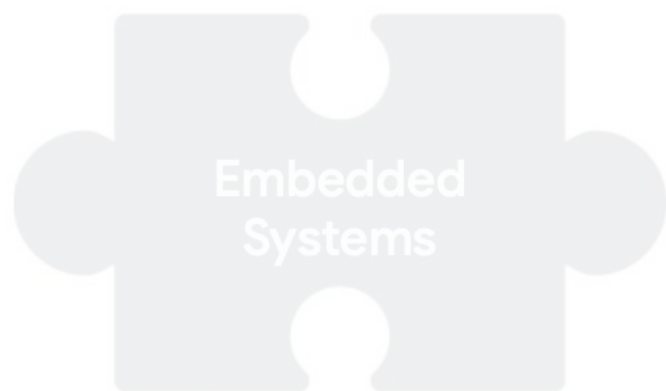# IESTI01 – TinyML

## Embedded
## Machine Learning

4. TinyML Challenges:
 - Machine Learning

Prof. Marcelo Rovai

UNIFEI
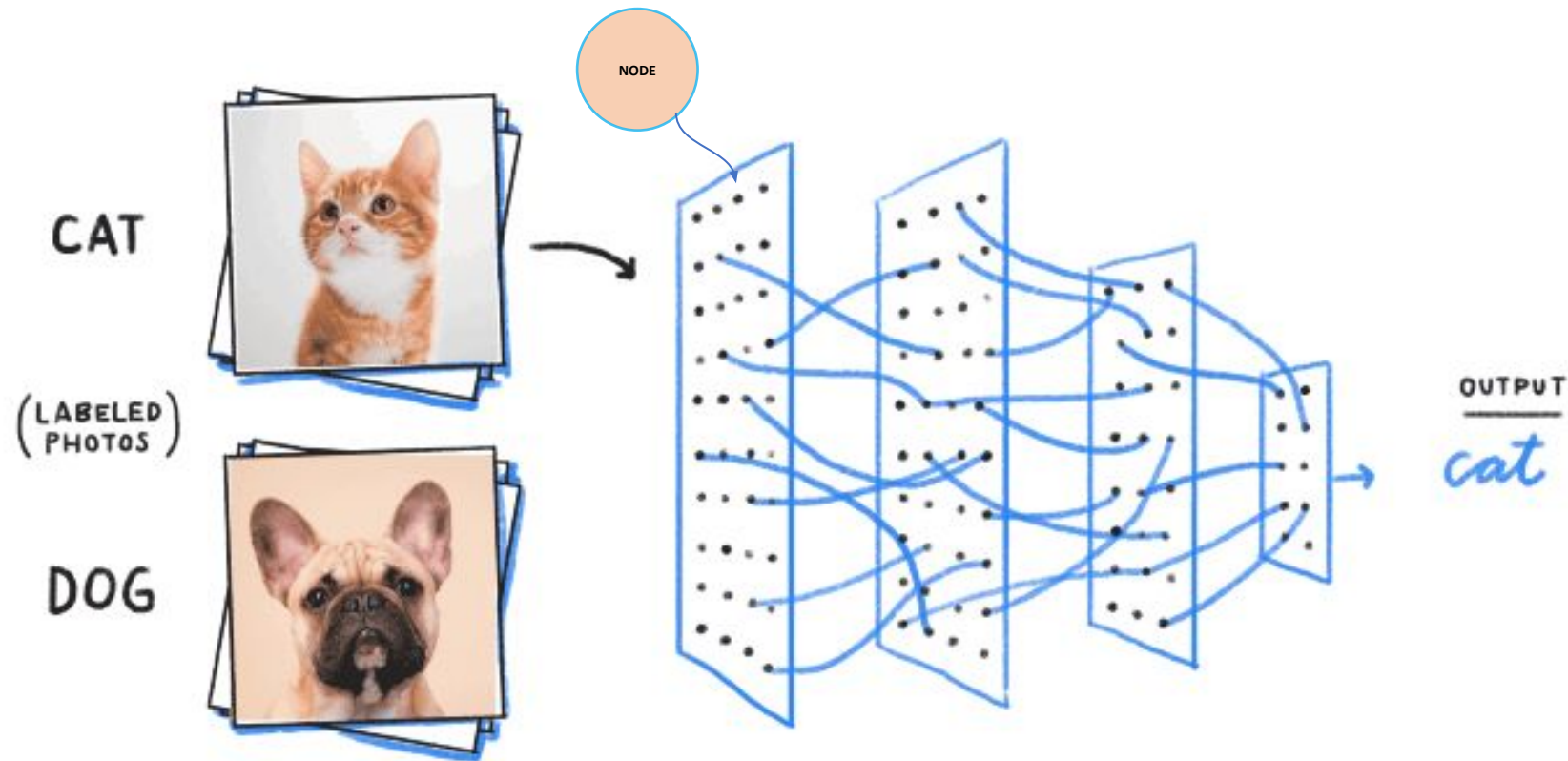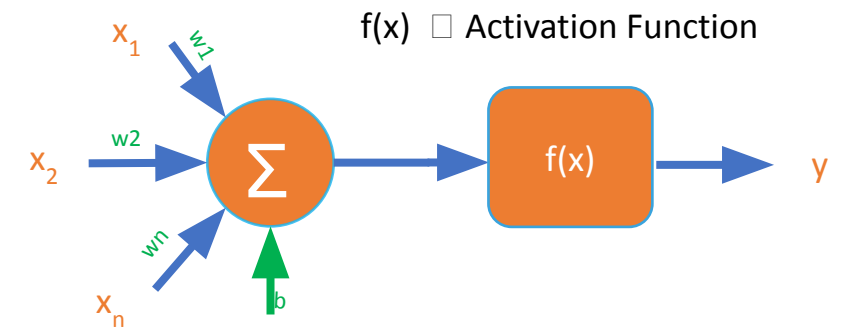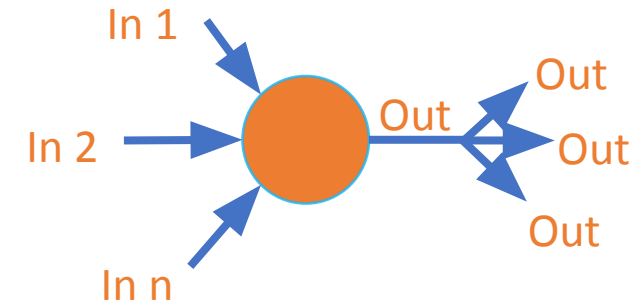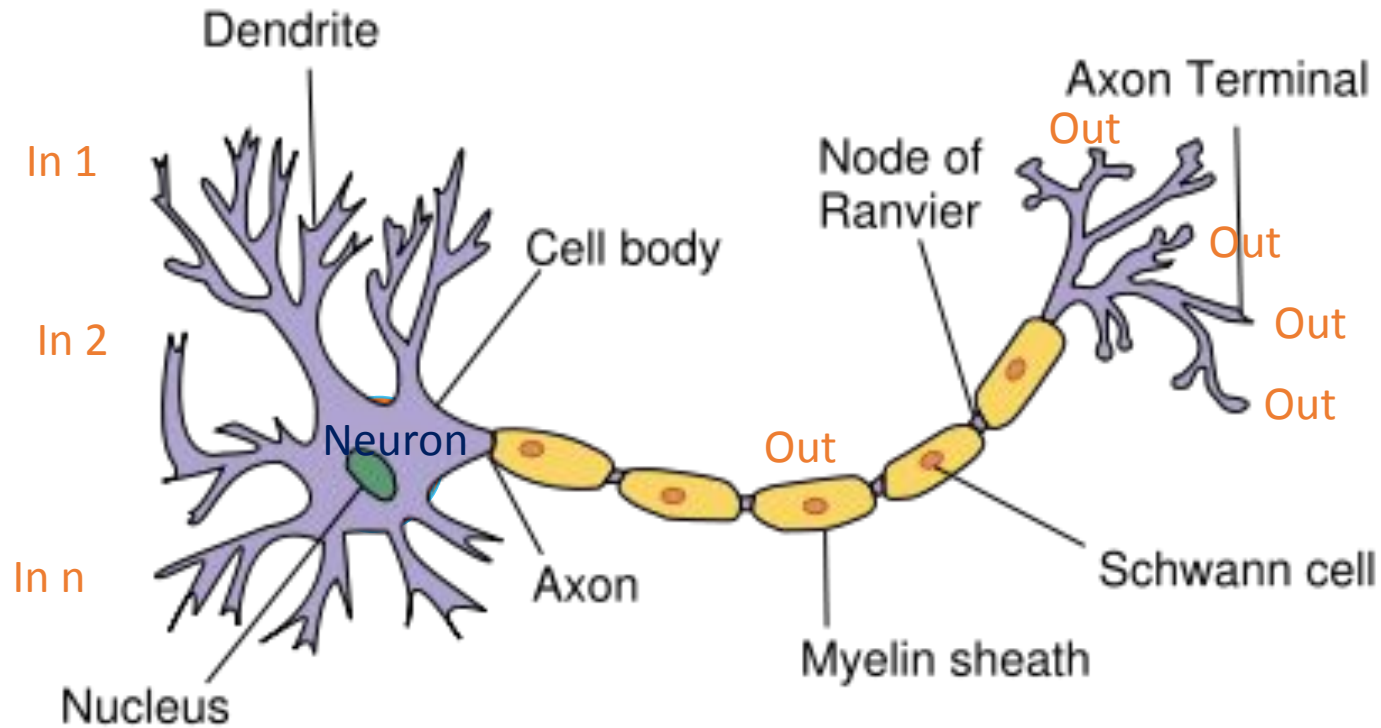
Embedded Systems + Machine Learning → TinyML

# (Deep) Machine Learning

Deep Learning: Subset of Machine Learning in which multilayered neural networks learn from vast amounts of data

# Neuron (Perceptron)



Dendrite

Cell body

Node of Ranvier

Axon Terminal

In 1

In 2

Neuron

In n

Nucleus

Axon

Myelin sheath

Schwann cell

Out

Out

Out

Out

Out

In 1

In 2

In n

Out

Out

Out

Out

$x_1$  $w1$

$x_2$  $w2$

$x_n$  $wn$  $b$

$\Sigma$

$f(x)$

$y$

f(x) ☐ Activation Function

Parameters

$$y = f\left(\sum_{i=1}^{n} x_i w_i + b\right)$$

# The Neural Network Model Architecture



A mostly complete chart of **Neural Networks**

©2019 Fjodor van Veen & Stefan Leijnen    asimovinstitute.org

**Legend:**
- Input Cell
- Backfed Input Cell
- Noisy Input Cell
- Hidden Cell
- Probablistic Hidden Cell
- Spiking Hidden Cell
- Capsule Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- Gated Memory Cell
- Kernel
- Convolution or Pool

**Network types:**
- Perceptron (P)
- Feed Forward (FF)
- Radial Basis Network (RBF)
- Deep Feed Forward (DFF)
- Recurrent Neural Network (RNN)
- Long / Short Term Memory (LSTM)
- Gated Recurrent Unit (GRU)
- Auto Encoder (AE)
- Variational AE (VAE)
- Denoising AE (DAE)
- Sparse AE (SAE)
- Markov Chain (MC)
- Hopfield Network (HN)
- Boltzmann Machine (BM)
- Restricted BM (RBM)
- Deep Belief Network (DBN)
- Deep Convolutional Network (DCN)
- Deconvolutional Network (DN)
- Deep Convolutional Inverse Graphics Network (DCIGN)
- Generative Adversarial Network (GAN)
- Liquid State Machine (LSM)
- Extreme Learning Machine (ELM)
- Echo State Network (ESN)
- Deep Residual Network (DRN)
- Differentiable Neural Computer (DNC)
- Neural Turing Machine (NTM)
- Capsule Network (CN)
- Kohonen Network (KN)
- Attention Network (AN)

# The Neural Network Model Architecture



A mostly complete chart of
## Neural Networks
©2019 Fjodor van Veen & Stefan Leijnen  asimovinstitute.org

https://www.asimovinstitute.org/neural-network-zoo/

# ML Model Size Growth

# ML Compute Needs (2012 to Present Day)



AlexNet to AlphaGo Zero: A 300,000x Increase in Compute (Log Scale)

Petaflop/s-days

Source: https://openai.com

In recent years, **computing needs grew by 300,000x** to train the machine learning models that are widely deployed in the industry

# ML Compute Needs (from the 1960s)



**Two Distinct Eras of Compute Usage in Training AI Systems**

Petaflop/s-days

1e+4

1e+2 — AlphaGoZero

Neural Machine Translation

1e+0 — TI7 Dota 1v1

VGG

1e-2 — ResNets

AlexNet

3.4-month doubling

1e-4 — Deep Belief Nets and layer-wise pretraining

DQN

1e-6

TD-Gammon v2.1

1e-8 — BiLSTM for Speech

LeNet-5

NETtalk

1e-10 — RNN for Speech

ALVINN

1e-12

2-year doubling (Moore's Law)

1e-14 — Perceptron

← First Era    Modern Era →

1960   1970   1980   1990   2000   2010   2020

Source: https://openai.com

**In recent years,** the amount of computing needed has grown remarkably fast.

Compute requirements are **doubling nearly every 3 to 4 months**

Source: Google


Source: Google

Cloud TPU

TinyML

# ML Model Evolution

**Source:** S. Bianco, R. Cadene, L. Celona, and P. Napoletano, "Benchmark analysis of representative deep neural network architectures," *IEEE Access*, vol. 6, pp. 64 270–64 277, 2018
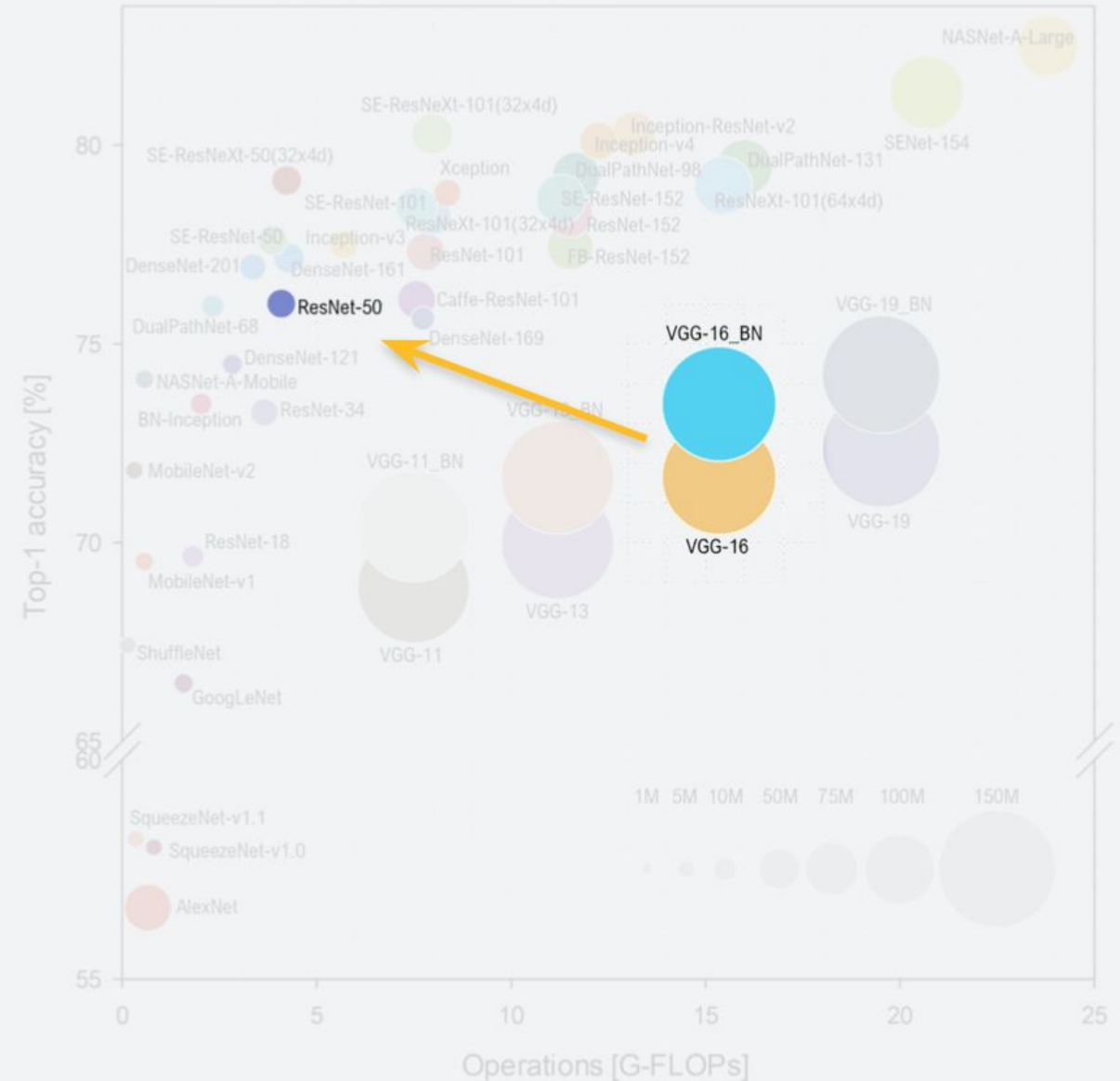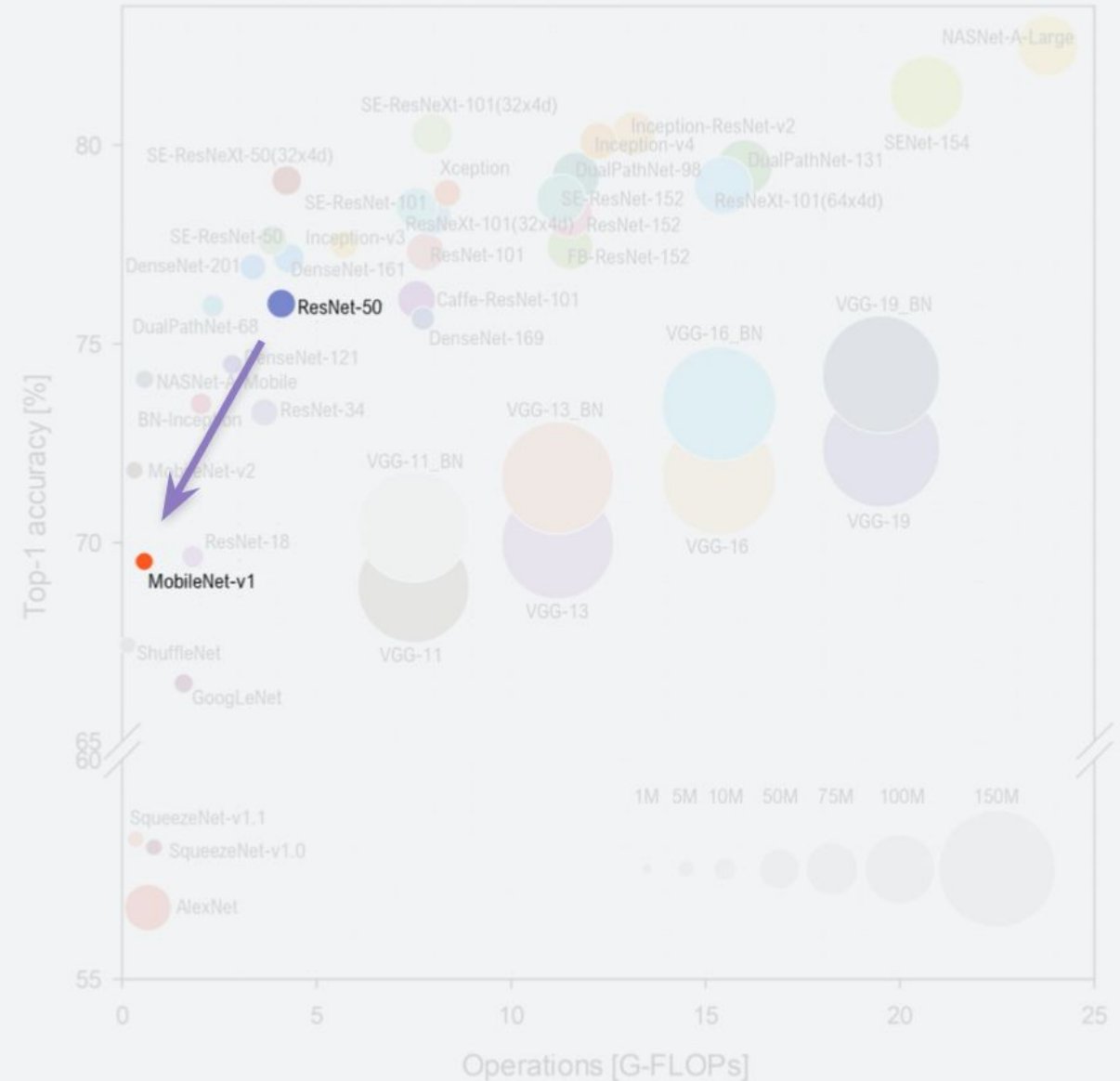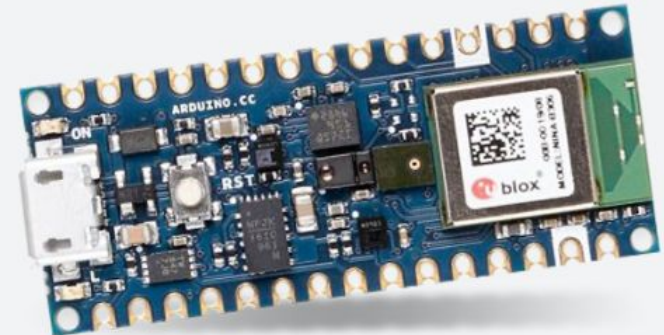
# ML Model Evolution

- **AlexNet** (2012)
  - 57.1% accuracy
  - 61MB in size



**Source:** S. Bianco, R. Cadene, L. Celona, and P. Napoletano, "Benchmark analysis of representative deep neural network architectures," *IEEE Access*, vol. 6, pp. 64 270–64 277, 2018
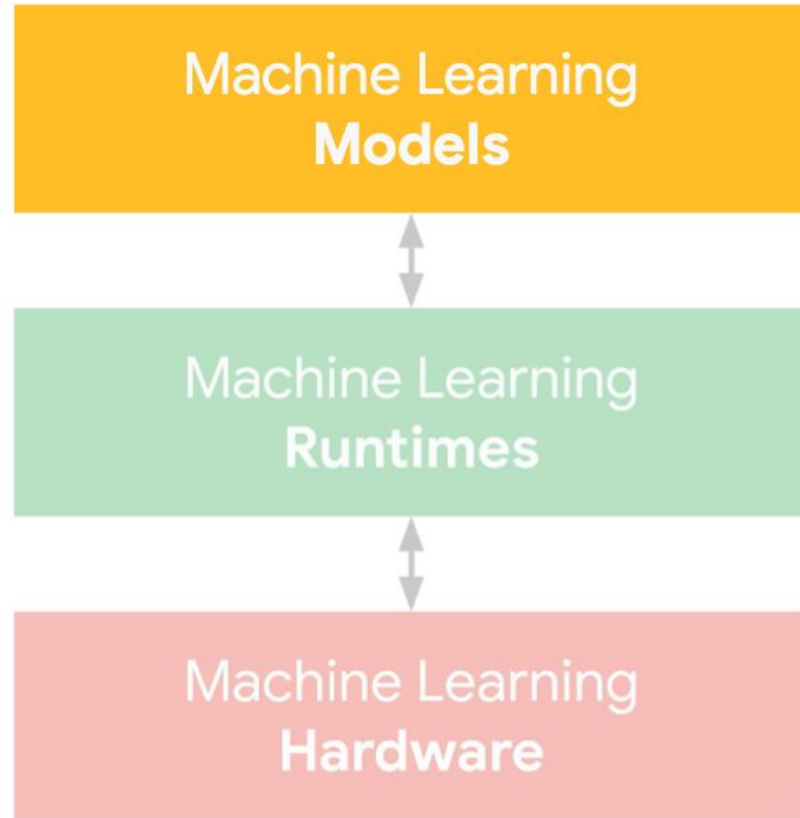
# ML Model Evolution

- **VGGNet** (2014) [VGG-16]
  - **71.5%** accuracy
  - **528MB** in size

**Source:** S. Bianco, R. Cadene, L. Celona, and P. Napoletano, "Benchmark analysis of representative deep neural network architectures," *IEEE Access*, vol. 6, pp. 64 270–64 277, 2018

# ML Model Evolution

- **ResNet** (2015)
  - **75.8%** accuracy
  - **22.7MB** in size

**Source:** S. Bianco, R. Cadene, L. Celona, and P. Napoletano, "Benchmark analysis of representative deep neural network architectures," *IEEE Access*, vol. 6, pp. 64 270–64 277, 2018

# ML Model Evolution

- **MobileNet** (2015)
  - *MobileNetv1*
    - **70.6%** accuracy
    - **16.9MB** in size

**Source:** S. Bianco, R. Cadene, L. Celona, and P. Napoletano, "Benchmark analysis of representative deep neural network architectures," *IEEE Access*, vol. 6, pp. 64 270–64 277, 2018

# ML Model Evolution

- **MobileNet** (2015)
  - *MobileNetv1*
    - 70.6% accuracy
    - 16.9MB in size

**Source:** S. Bianco, R. Cadene, L. Celona, and P. Napoletano, "Benchmark analysis of representative deep neural network architectures," *IEEE Access*, vol. 6, pp. 64 270–64 277, 2018

# Problem:

Our board (in your kit for Course )

only has **256KB** of RAM (memory)

yet **MobileNetv1** needs **16.9MB**!

Machine Learning **Models**

Machine Learning **Runtimes**

Machine Learning **Hardware**

# Model Compression Techniques

**Pruning**

Quantization

Knowledge Distillation

...

# Pruning



PRUNING
SYNAPSES

# Pruning



PRUNING
NEURONS

Machine Learning **Models**

Machine Learning **Runtimes**

Machine Learning **Hardware**

# Model Compression Techniques

Pruning

**Quantization**

Knowledge Distillation

...

# Quantization

# Knowledge Distillation



TEACHER → STUDENT

Machine Learning
**Models**

Machine Learning
**Runtimes**

Machine Learning
**Hardware**

TensorFlow

[TF Video]

Less memory

Less compute power

Only focused on *inference*

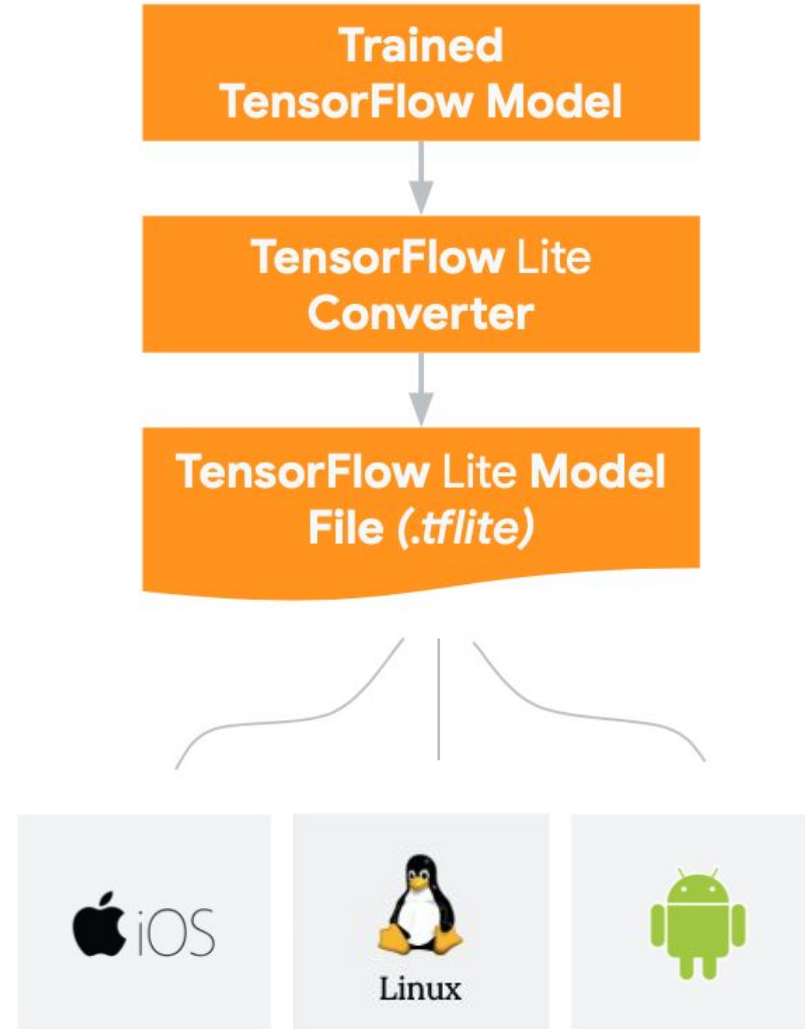TensorFlow Lite
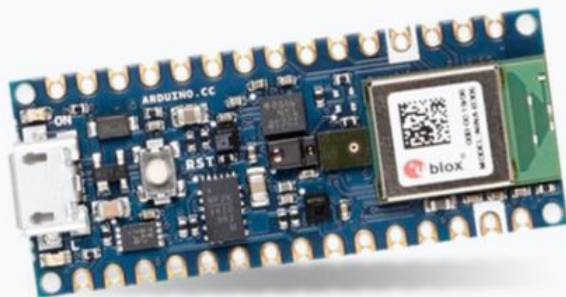
# Key Differences

|  | **TensorFlow** | **TensorFlow** Lite |
|---|---|---|
| **Topology** | **Variable** | **Fixed** |
| **Weights** | **Variable** | **Fixed** |
| **Binary Size** | **Unimportant** | **High Priority** |
| **Distributed Compute** | **Needed** | **Not Needed** |
| **Developer Background** | **ML Researcher** | **Application Developer** |

# Architecture



Trained
TensorFlow Model

↓

TensorFlow Lite
Converter

↓

TensorFlow Lite Model
File (.tflite)

iOS    Linux    Android

**Even less memory**

**Even less compute power**

**Also**, only focused on *inference*

**TensorFlow**

**TensorFlow** Lite

Train a model → Convert model → Optimize model → Deploy model at Edge → Make inferences at Edge

TensorFlow

TensorFlow Lite

Train a model → **Convert model** → **Optimize model** → **Deploy model at Edge** → **Make inferences at Edge**

TensorFlow

TensorFlow Lite

Train a model → Convert model → Optimize model → **Deploy model at Edge** → **Make inferences at Edge**

Raspberry Pi

Linux

iOS

Android

(TF Micro)

Microcontroller
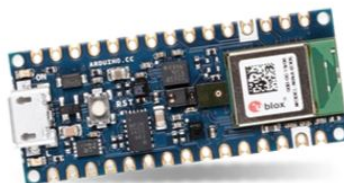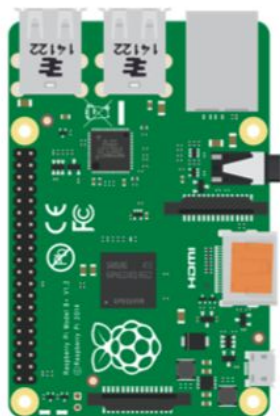
TensorFlow

TensorFlow Lite

Train a model  >  Convert model  >  Optimize model  >  Deploy model at Edge  >  **Make inferences at Edge**

TensorFlow Lite

**TFL Question and Answer**

Please select an article below.

TensorFlow

Google

Super_Bowl_50

Warsaw

Normans

Nikola_Tesla

Computational_complexity_theory
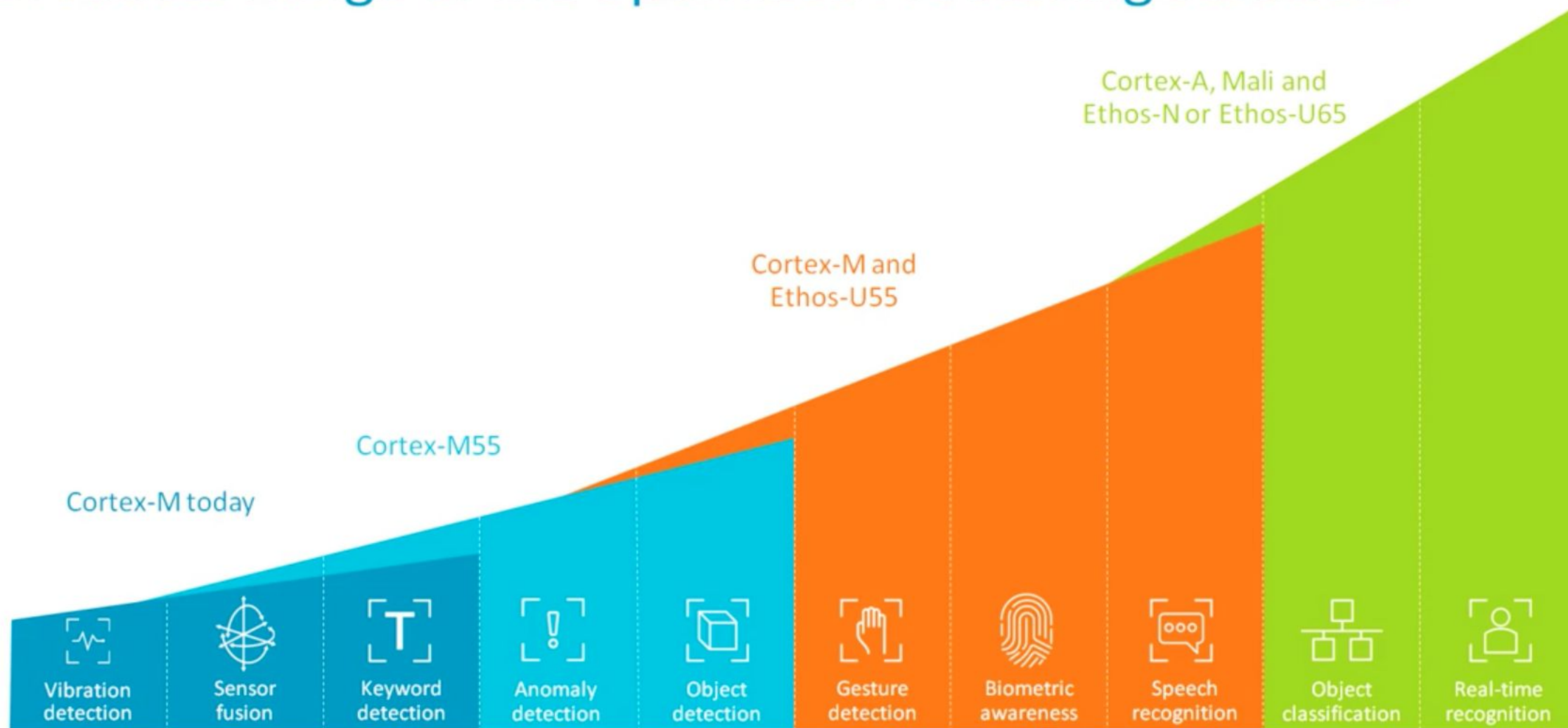
Teacher

Martin_Luther

"Energy-efficient On-device Processing for Next-generation Endpoint ML"

By Tomas Edso, Senior Principal Engineer (ML), ARM
At tinyML Summit 2020 presentation

# Reading Material

# Main references

- [Harvard School of Engineering and Applied Sciences - CS249r: Tiny Machine Learning](#)
- [Professional Certificate in Tiny Machine Learning (TinyML) – edX/Harvard](#)
- [Introduction to Embedded Machine Learning - Coursera/Edge Impulse](#)
- [Computer Vision with Embedded Machine Learning - Coursera/Edge Impulse](#)
- Fundamentals textbook: ["Deep Learning with Python" by François Chollet](#)
- Applications & Deploy textbook: ["TinyML" by Pete Warden, Daniel Situnayake](#)
- Deploy textbook ["TinyML Cookbook" by Gian Marco Iodice](#)

**I want to thank Shawn Hymel and Edge Impulse, Pete Warden and Laurence Moroney from Google, Professor Vijay Janapa Reddi and Brian Plancher from Harvard, and the rest of the TinyMLedu team for preparing the excellent material on TinyML that is the basis of this course at UNIFEI.**

The IESTI01 course is part of the **TinyML4D**, an initiative to make TinyML education available to everyone globally.

**Thanks**