# Embedded ML (TinyML) Intro & Applications

**Prof. Marcelo J. Rovai**

rovai@unifei.edu.br

UNIFEI - Federal University of Itajuba, Brazil
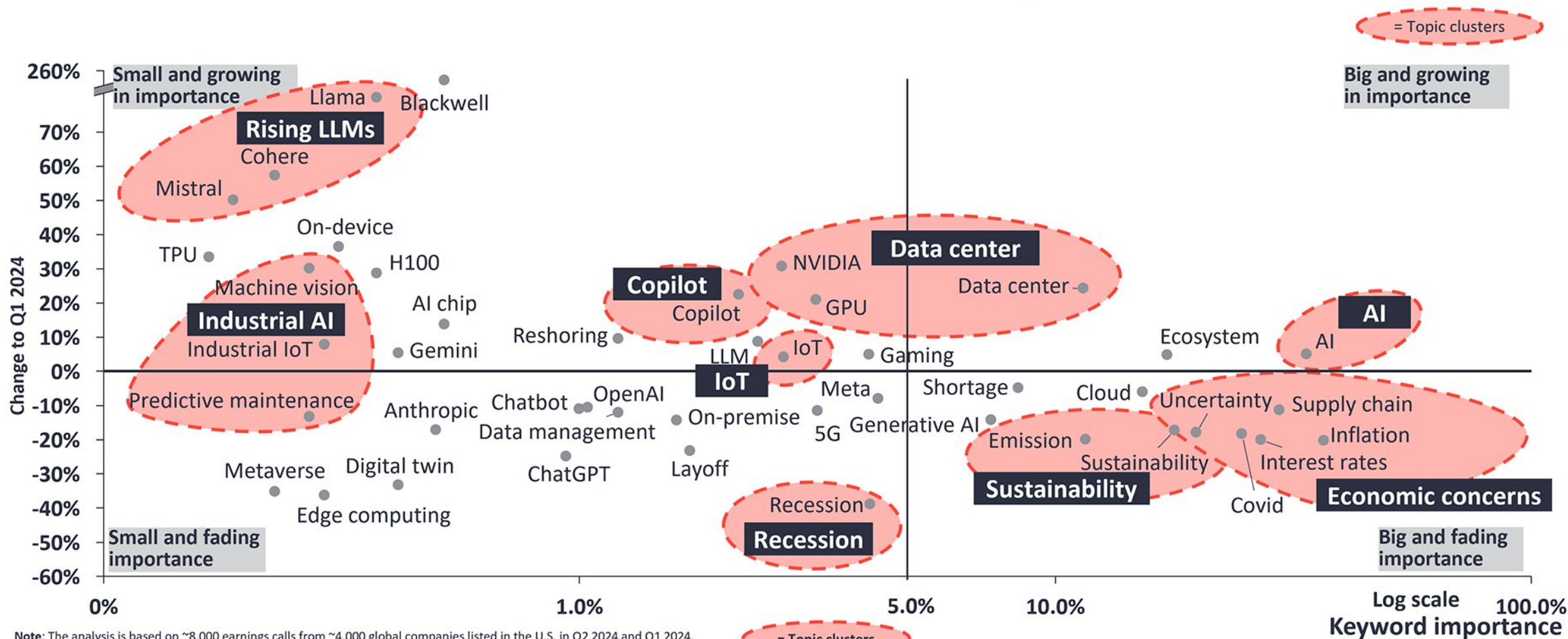TinyML4D Academic Network Co-Chair

1

# Internet of Things (IoT)

What CEOs talked about in Q2 2024 (vs. Q1 2024)

https://iot-analytics.com/what-ceos-talked-about-in-q2-2024/

IOT ANALYTICS

**What CEOs talked about**

"LLM model size competition is intensifying. … backwards!"

*Andrej Katpathy*

Note: The analysis is based on ~8,000 earnings calls from ~4,000 global companies listed in the U.S. in Q2 2024 and Q1 2024. The mentions of the selected keywords in each call were counted in each quarter.
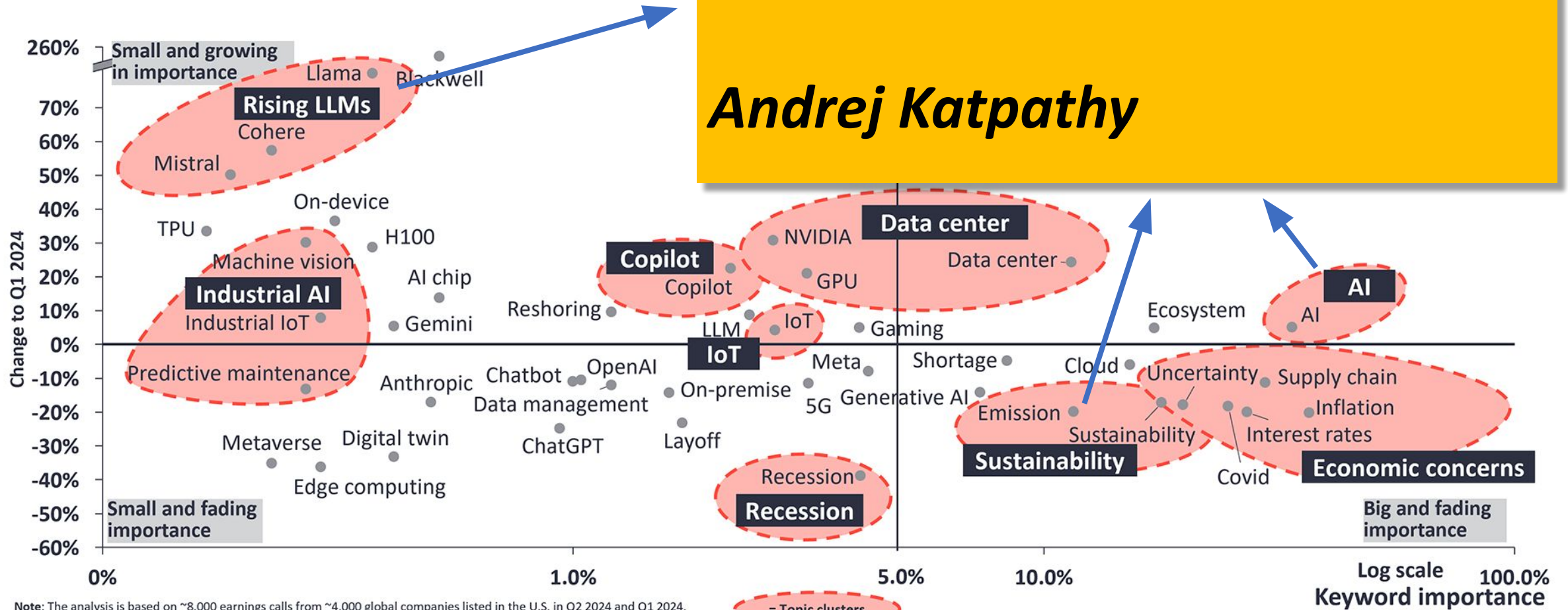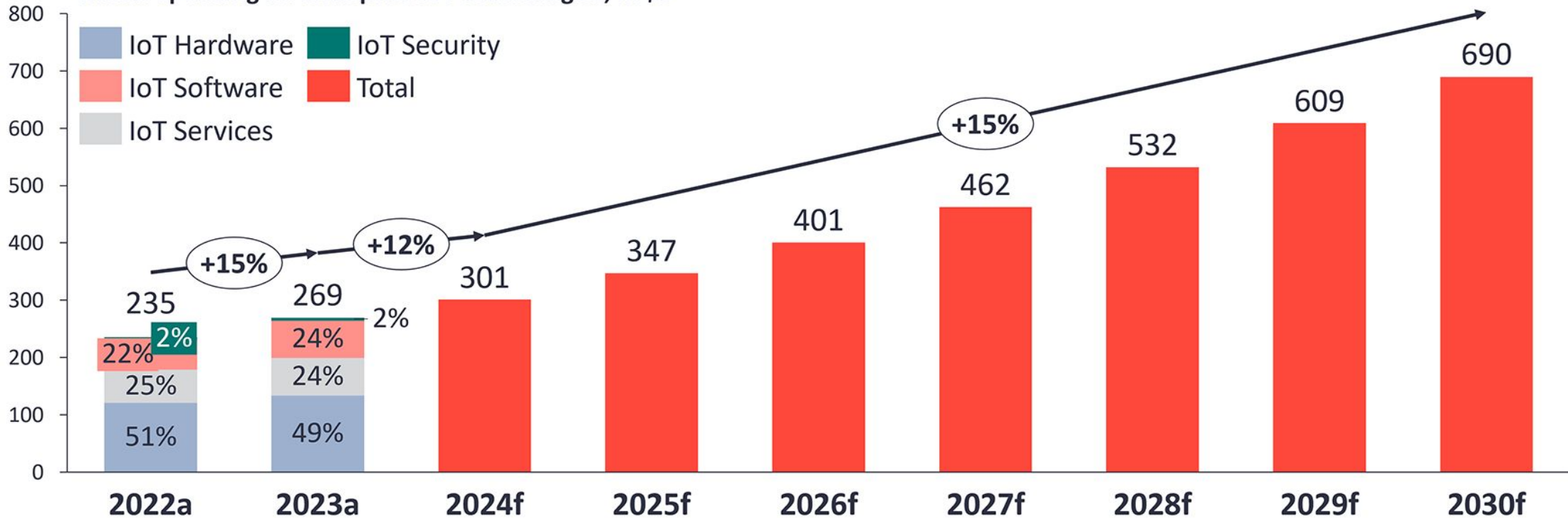Source: IoT Analytics Research 2024 – We welcome republishing of images but ask for source citation with a link to the original post and company website.

(Share of companies that mentioned the keyword in Q2 2024 at least once)

https://iot-analytics.com/what-ceos-talked-about-in-q2-2024/

The enterprise IoT market by technology 2023–2030

# Typical IoT Project



* "Things"

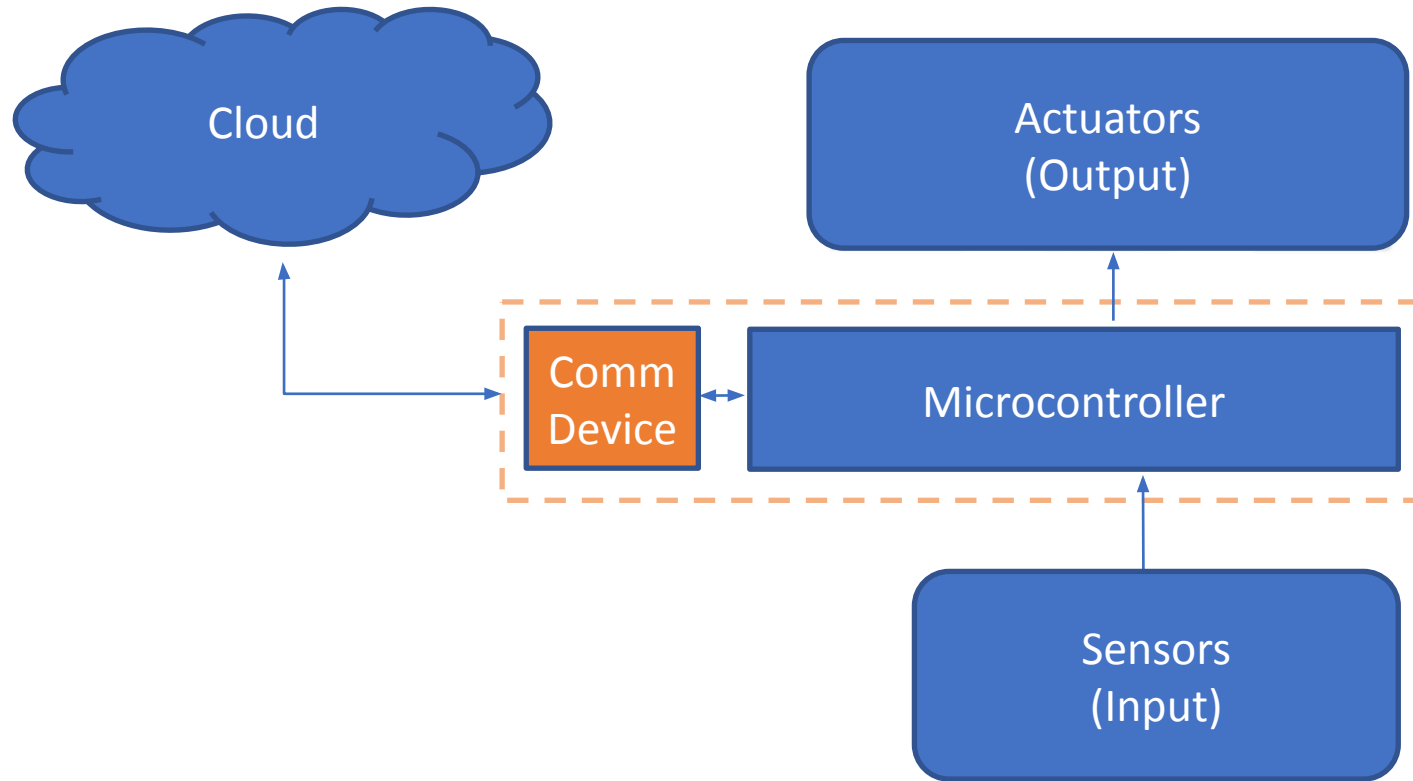# Typical IoT Project



**5 Quintillion**
bytes of data produced every day by IoT

**<1%**
of unstructured data is analyzed or used at all

Source: Harvard Business Review, What's Your Data Strategy?, April 18, 2017
Cisco, Internet of Things (IoT) Data Continues to Explode Exponentially. Who Is Using That Data and How?, Feb 5, 2018

# Typical AIoT Project

# Typical AIoT Project … … Issues

# Typical AIoT Project ...                              ... **Issues**

**AI**

Cloud

Actuators
(Output)

Comm
Device

Microcontroller

Sensors
(Input)

**Bandwidth**

**Latency**

# Typical AIoT Project …                    … **Issues**



**Bandwidth**

**Latency**

**Energy**

# Typical AIoT Project ...



## ... **Issues**

**B**andwidth

**L**atency

**E**nergy

**R**eliability

# Typical AIoT Project …                    … **Issues**

**AI**

Cloud

Comm Device ↔ Microcontroller

Actuators (Output)

Sensors (Input)

**Google** Assistant

amazon

**B**andwidth

**L**atency

**E**nergy

**R**eliability

**P**rivacy

# Typical AIoT Project …                    … Issues

AI

Cloud

**B**andwidth

Actuators
(Output)

**L**atency

Comm
Device

Microcontroller

**E**nergy

Sensors
(Input)

**R**eliability

**P**rivacy

… **Solution ?**

# IoT 2.0 * – Edge AI/ML    * Intelligence of Things



**Cloud**

**Actuators (Output)**

**AI**

**LoRa Wan**

**Microcontroller**

Few and eventual data to be sent to the cloud

**Sensors (Input)**

**Low Power**

# … Solution -> ML goes close to data

# When to use an Edge AI/ML approach:



Cloud

Actuators
(Output)

AI

Comm
Device

Microcontroller

Sensors
(Input)

**B**andwidth

**L**atency

**E**nergy

**R**eliability

**P**rivacy

The Distributed Intelligence Spectrum

Cloud

On Premise Servers

Gateway

Intelligent Device

Ultra Low Powered Devices and Sensors

TinyML

EDGE AI

CLOUD AI

*Source: ABI Research: TinyML*

17

# Market Forecast



## Very Edge AI-Enabled Device Global Shipments by Vertical

Legend:
- Utilities
- Transport and Logistics
- Smart Cities
- Retail
- Industrial and Manufacturing
- Healthcare
- Consumer
- Banking and Finance
- Agriculture

Y-axis: (MILLIONS) — 0, 500, 1,000, 1,500, 2,000, 2,500, 3,000
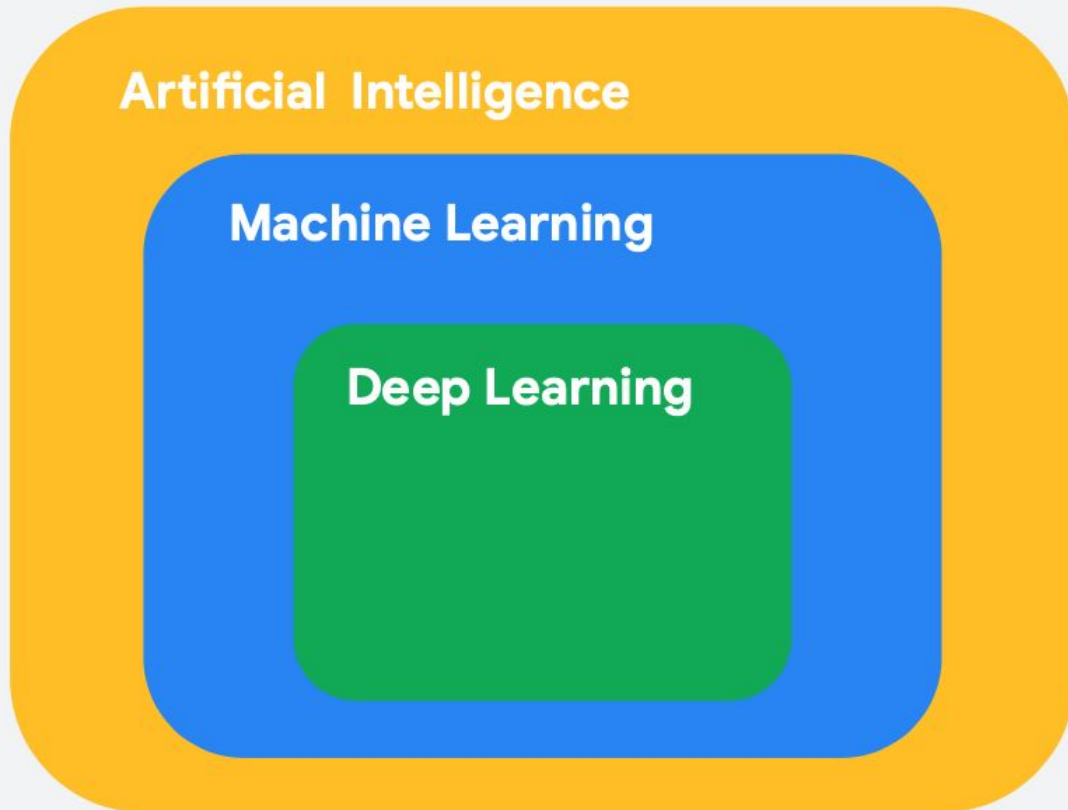
X-axis: 2021, 2024, 2027, 2030

*Source: ABI Research: TinyML*
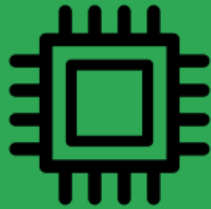
# Embedded ML (TinyML)

Introduction

**AI: Any technique that enables computers to mimic human behavior**

**ML: Ability to learn without explicitly being programed**

**DL: Extract patterns from data using neural networks**

**Edge AI (or Edge ML)** is the processing of Artificial Intelligence algorithms on edge, that is, on users' devices. The concept derives from **Edge Computing**, which starts from the same premise: data is stored, processed, and managed directly at the Internet of Things (IoT) endpoints.

**TinyML is a subset of EdgeML, where** sensors are generating data with ultra-low power consumption (batteries), so that we can ultimately deploy machine learning continuously ("always on devices")

# What is Tiny Machine Learning (**TinyML**)?

**TinyML**

Fastest-growing field of **ML**

# What is Tiny Machine Learning (**TinyML**)?

**TinyML**

Fastest-growing field of **ML**

Algorithms, hardware, software

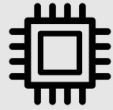# What is Tiny Machine Learning (TinyML)?
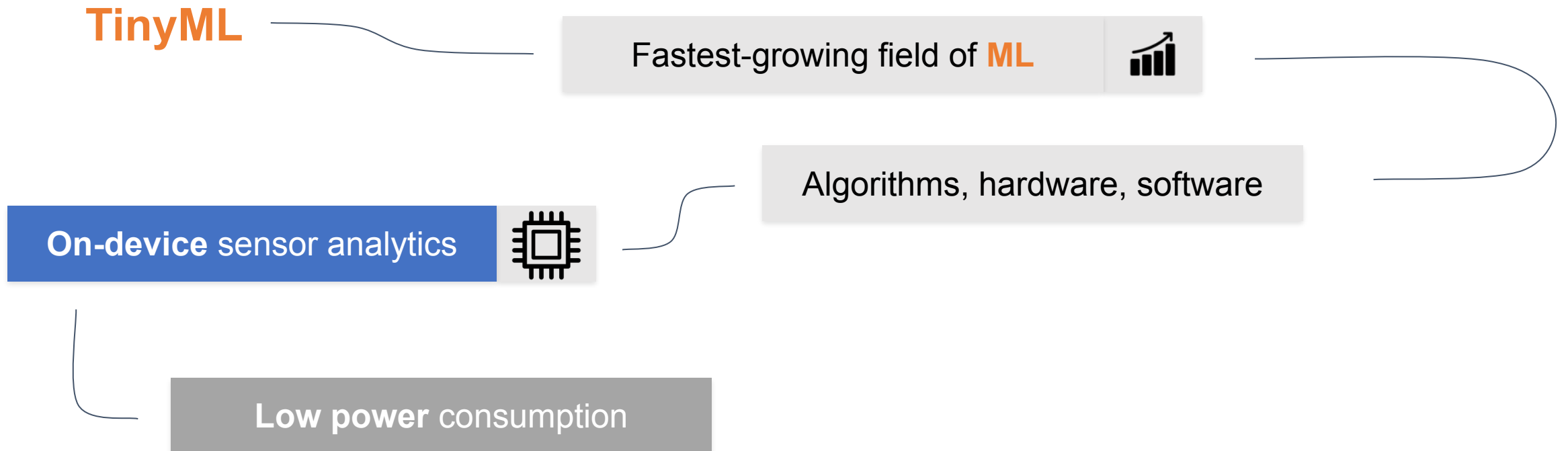
**TinyML**

Fastest-growing field of **ML**
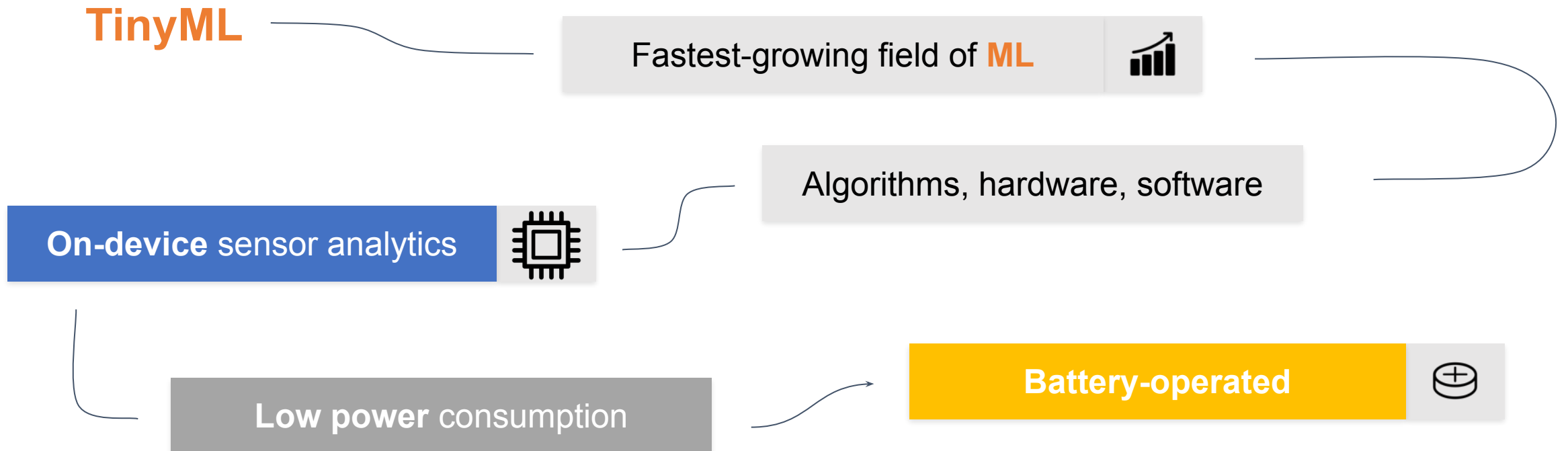
Algorithms, hardware, software

**On-device** sensor analytics

# What is Tiny Machine Learning (TinyML)?

**TinyML**

Fastest-growing field of **ML**

Algorithms, hardware, software

**On-device** sensor analytics

**Low power** consumption

# What is Tiny Machine Learning (TinyML)?

**TinyML**

Fastest-growing field of **ML**

Algorithms, hardware, software

**On-device** sensor analytics

**Low power** consumption

**Battery-operated**

# What is Tiny Machine Learning (TinyML)?

**TinyML**

Fastest-growing field of **ML**

Algorithms, hardware, software

**On-device** sensor analytics

**Low power** consumption

**Battery-operated**

**Always-on ML**

# What Makes **TinyML** ?

Embedded
Systems

Machine
Learning

**TinyML**

# What Makes **TinyML** ?



**TinyML**

# TinyML Challenges
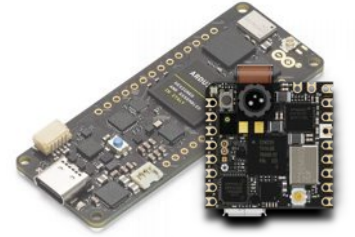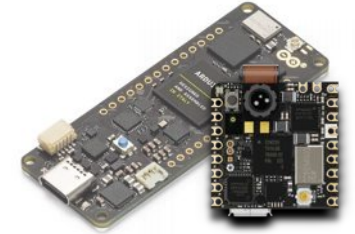
# 250 Billion
*MCUs today*

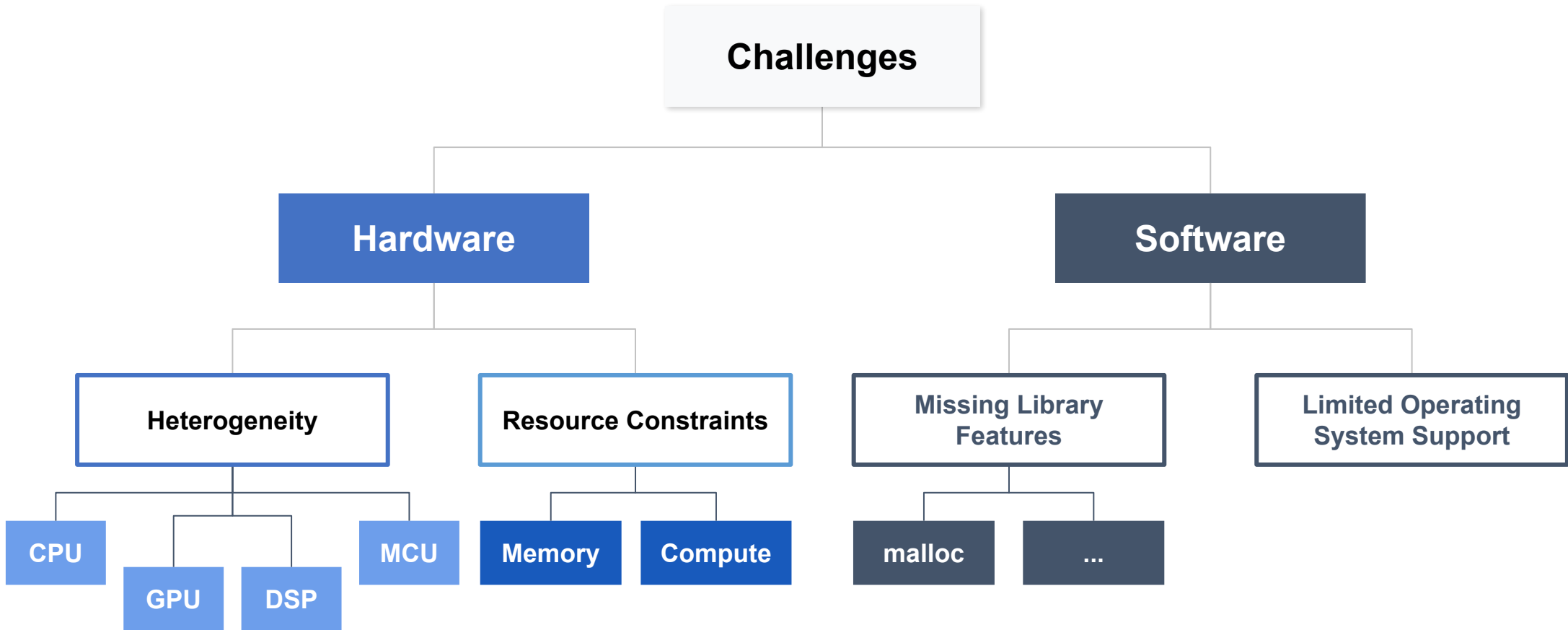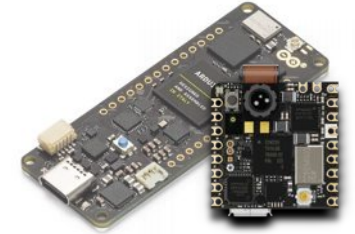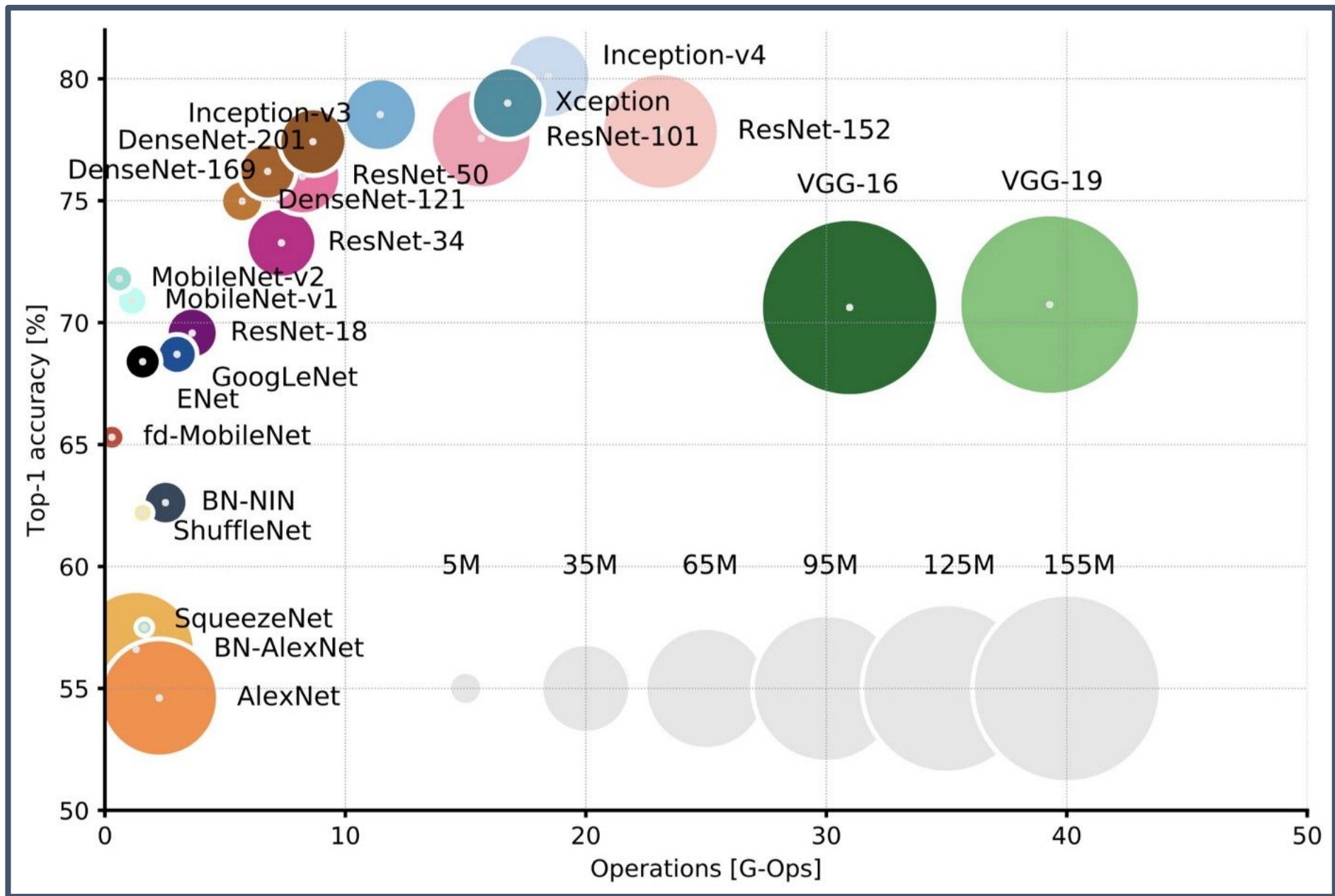# Hardware

# Hardware

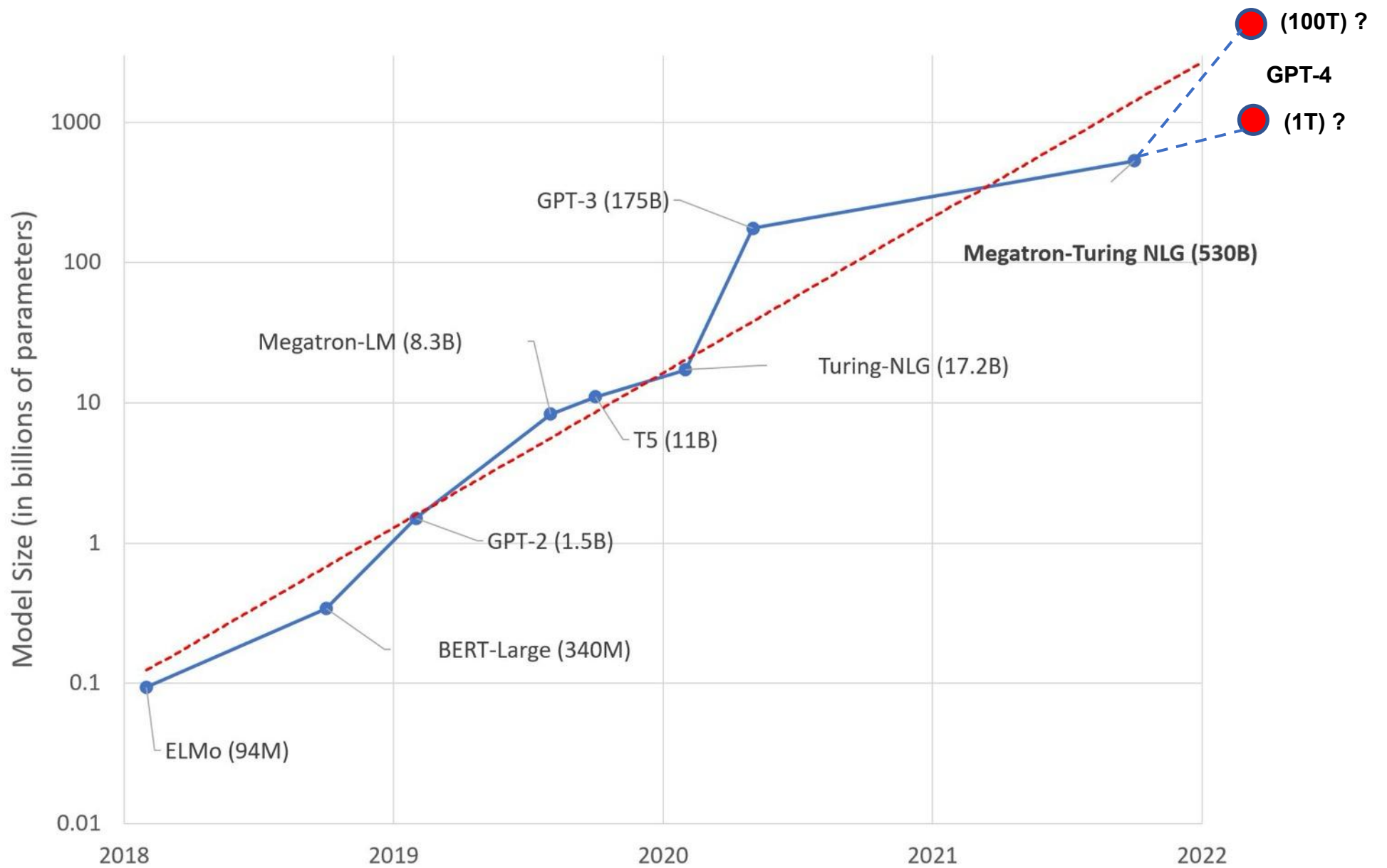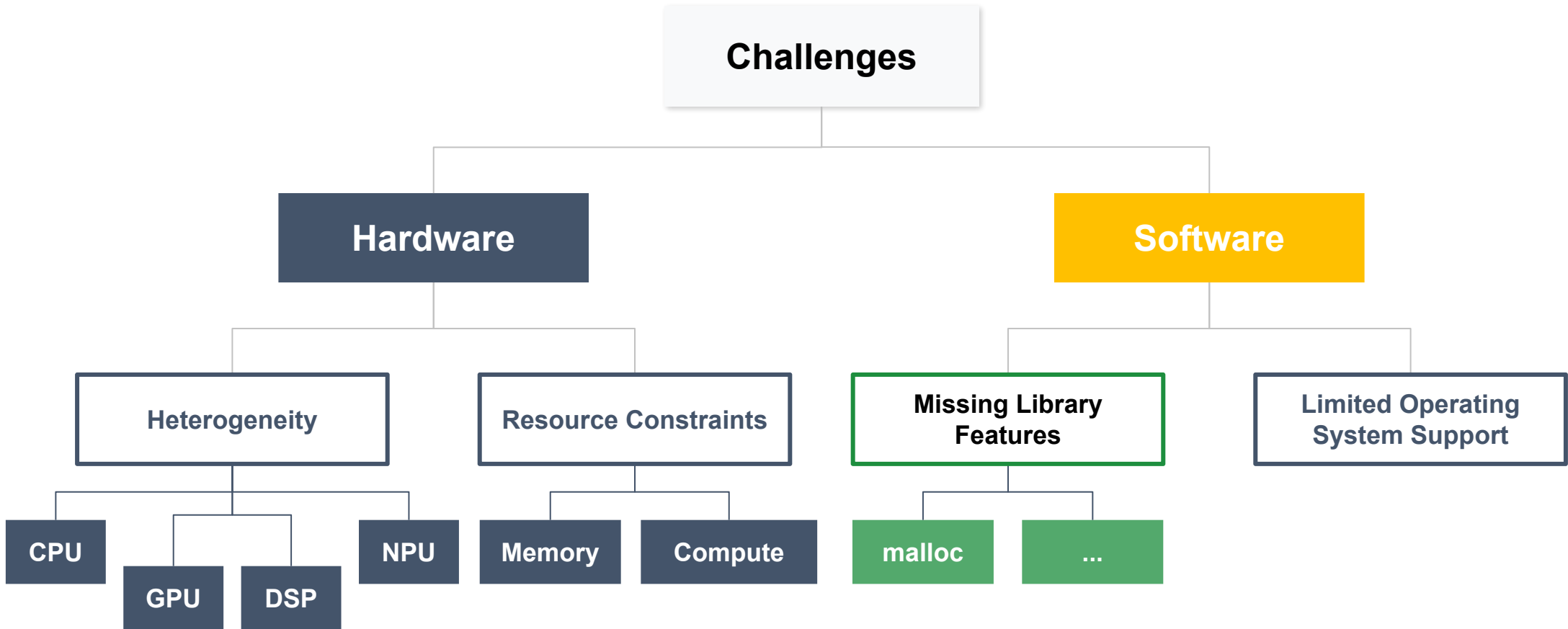| | Raspberry Pico (W) | Arduino Nano Sense | ESP 32 | Seeed XIAO Sense / ESP32S3 | Arduino Pro |
|---|---|---|---|---|---|
| **32Bits CPU** | Dual-core Arm Cortex-M0+ | Arm Cortex-M4F | Xtensa LX6 Dual Core | Arm Cortex-M4F (BLE) Xtensa LX7 Dual Core | Dual Core Arm Cortex M7/M4 |
| **CLOCK** | 133MHz | 64MHz | 240MHz | 64 / 240MHz | 480/240MHz |
| **RAM** | 264KB | 256KB | 520KB (part available) | 256KB / 8MB | 1MB |
| **ROM** | 2MB | 1MB | 2MB | 2MB / 8MB | 2MB |
| **Radio** | (Yes for W) | BLE | BLE/WiFi | BLE / WiFi (ESP32S3) | BLE/WiFi |
| **Sensors** | No | Yes | No | Yes (Sense) | Yes (Nicla) |
| **Bat. Power Manag.** | No | No | No | Yes | Yes |
| **Price** | $ | $$$ | $ | $$ | $$$$ |

https://media.digikey.com/Resources/Maker/the-original-guide-to-boards-2022.pdf

# Hardware

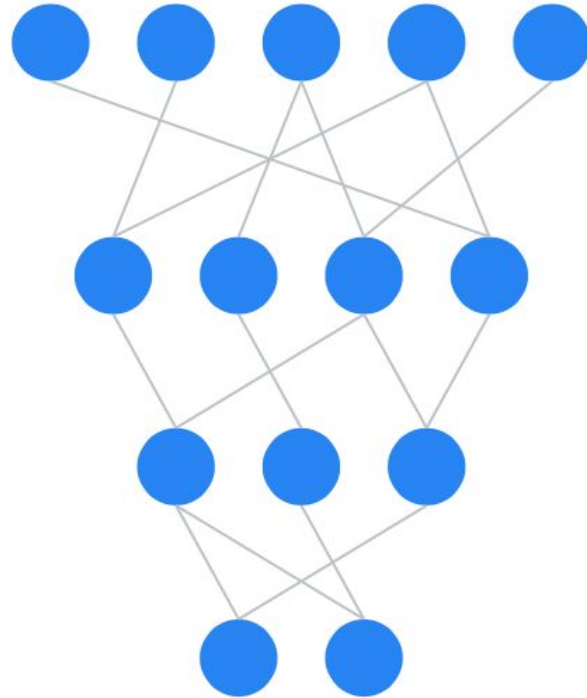| | Raspberry Pico (W) | Arduino Nano Sense | ESP 32 | Seeed XIAO Sense / ESP32S3 | Arduino Pro |
|---|---|---|---|---|---|
| **32Bits CPU** | Dual-core Arm Cortex-M0+ | Arm Cortex-M4F | Xtensa LX6 Dual Core | Arm Cortex-M4F (BLE) Xtensa LX7 Dual Core | Dual Core Arm Cortex M7/M4 |
| **CLOCK** | 133MHz | 64MHz | 240MHz | 64 / 240MHz | 480/240MHz |
| **RAM** | 264KB | 256KB | 520KB (part available) | 256KB / 8MB | 1MB |
| **ROM** | 2MB | 1MB | 2MB | 2MB / 8MB | 2MB |
| **Radio** | (Yes for W) | BLE | BLE/WiFi | BLE / WiFi (ESP32S3) | BLE/WiFi |
| **Sensors** | No | Yes | No | Yes (Sense) | Yes (Nicla) |
| **Bat. Power Manag.** | No | No | No | Yes | Yes |
| **Price** | $ | $$$ | $ | $$ | $$$$ |

https://media.digikey.com/Resources/Maker/the-original-guide-to-boards-2022.pdf

38

Model Size (in billions of parameters)

(100T) ?

GPT-4

(1T) ?

GPT-3 (175B)

Megatron-Turing NLG (530B)

Megatron-LM (8.3B)

Turing-NLG (17.2B)

T5 (11B)

GPT-2 (1.5B)

BERT-Large (340M)

ELMo (94M)

**Datasets Preprocessing**

**Quantization Pruning**

**Resource constraints**

Sound

Vision

Vibration

End-to-end **TinyML** application design

# Software

# Application Complexity vs. HW

Power

EdgeML

TinyML

Video Classification 2 MB+

Object Detection Complex Voice Processing 1 MB+

Image Classification 250 KB+

KeyWord Spotting Audio Classification 50 KB

Anomaly Detection Sensor Classification 20 KB

ESP32    XIAO

Application Complexity

CPU Power / Memory

Rpi-Pico (Cortex-M0+)

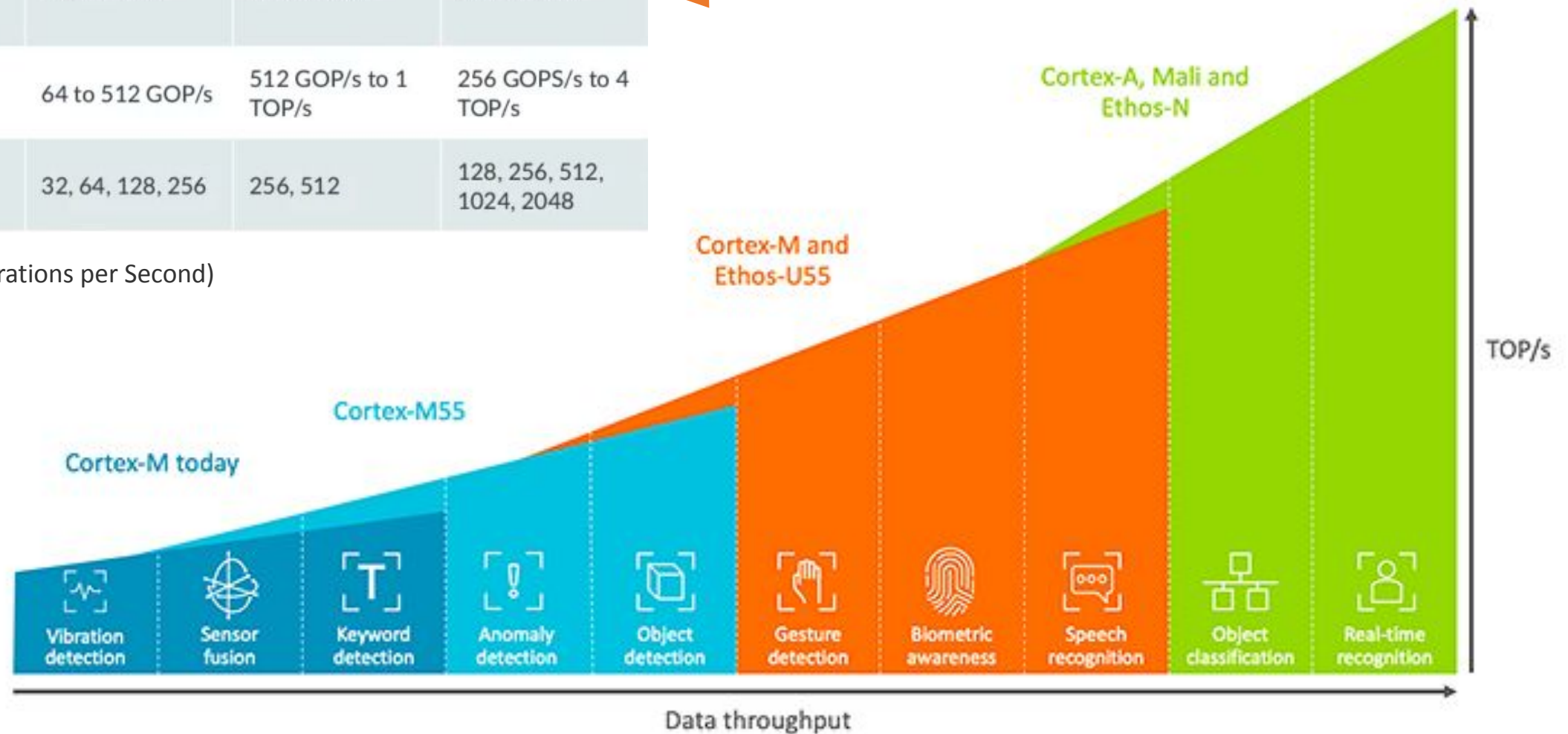Arduino Nano (Cortex-M4)

Arduino Pro (Cortex-M7)

RaspberryPi (Cortex-A)

SmartPhone

Jetson Nano/Orin (Cortex-A + GPU)

# ML- optimized Solutions (w/microNPUs)

| | Ethos-U55 | Ethos-U65 | Ethos-U85 |
|---|---|---|---|
| Performance (At 1 GHz) | 64 to 512 GOP/s | 512 GOP/s to 1 TOP/s | 256 GOPS/s to 4 TOP/s |
| MACs (8x8) | 32, 64, 128, 256 | 256, 512 | 128, 256, 512, 1024, 2048 |

**TOPS** (Tera Operations per Second)
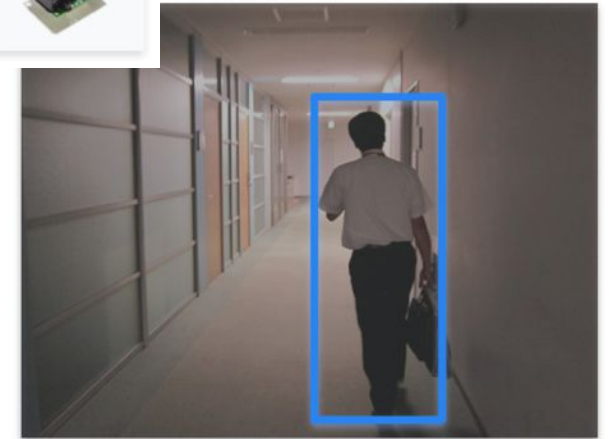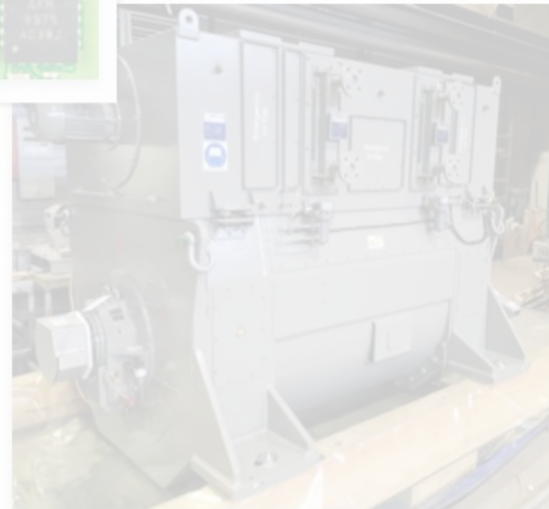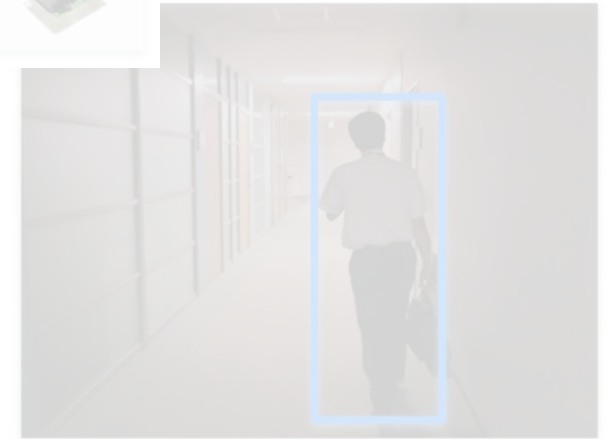
# TinyML Application
Examples

Machine Learning

| Supervised | Unsupervised | Reinforcement |
|---|---|---|
| Task driven | Data driven | Learns to react to an environment |

- Regression
- Classification
- Object Detection
- Time Series Forecasting

- Anomaly Detection
- Clustering
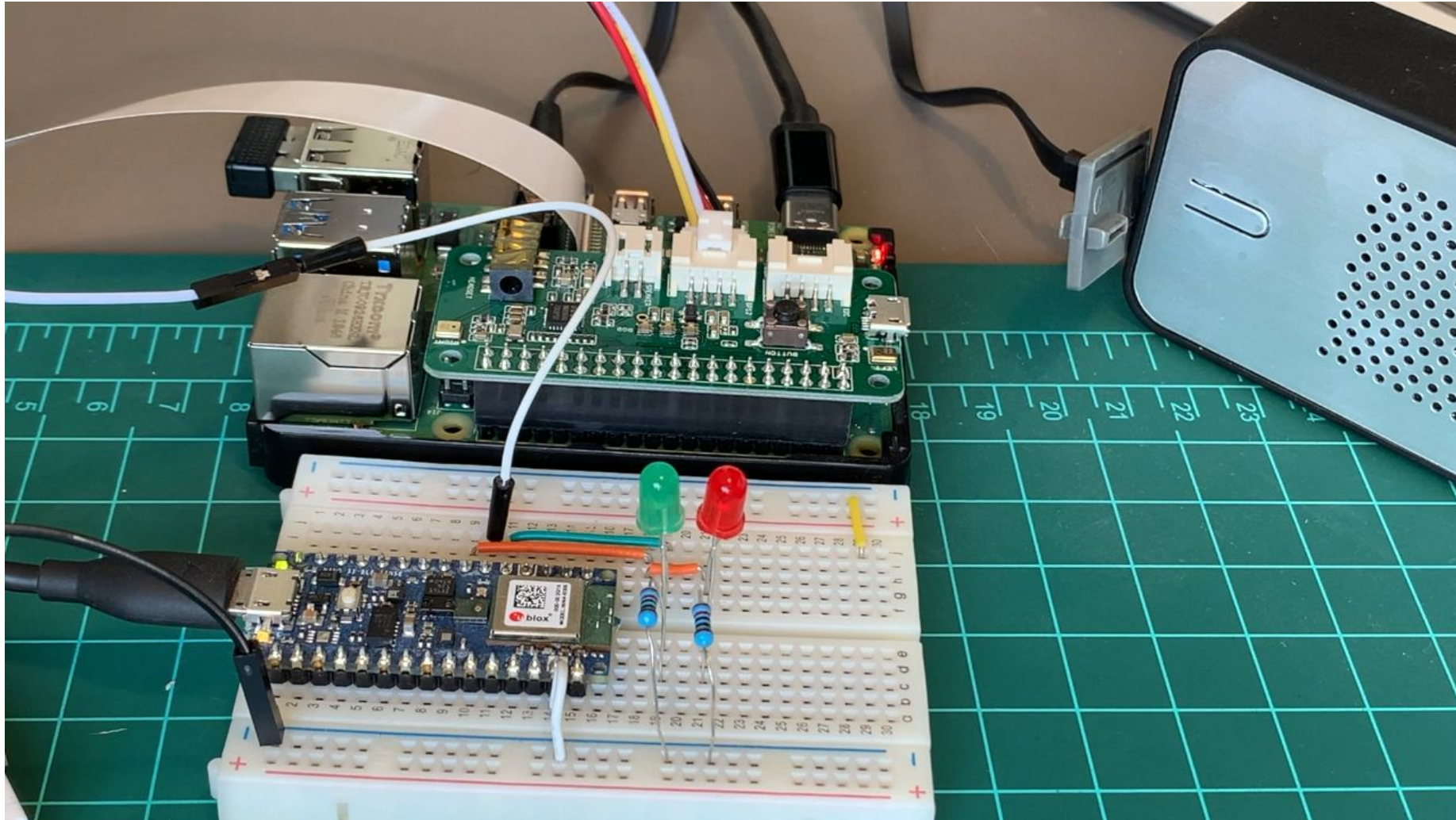- Dimensionality Reduction
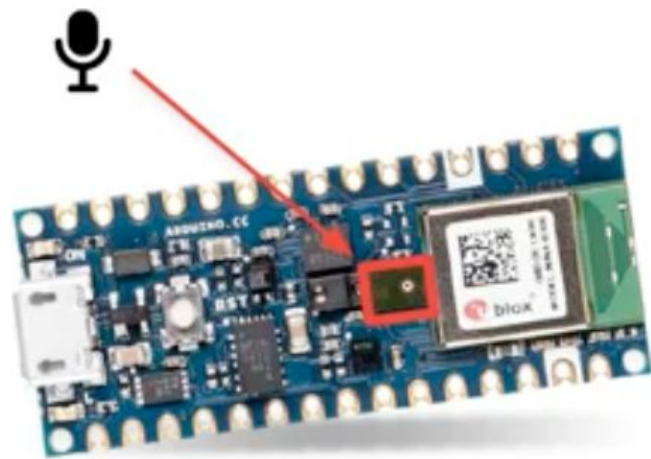
- Autonomous Navigation

Sound

Vibration

Vision

Sound

Vibration

Vision

# Personal Assistant

# Personal Assistant

# "Cascade" Detection: multi-stage model



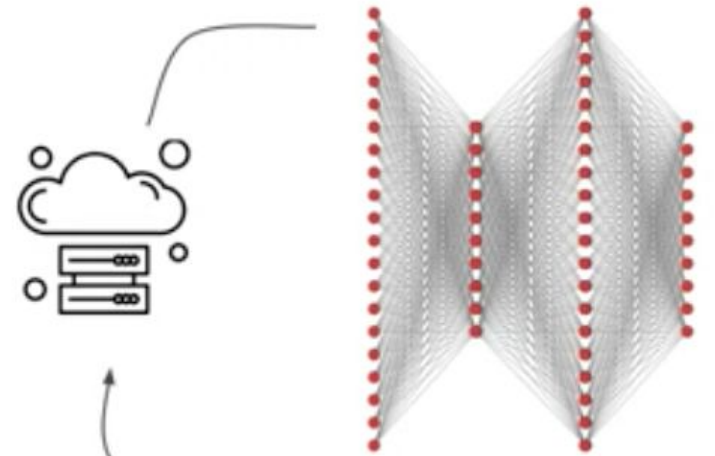5 Process the full speech data with a large model in the cloud

2 Process the data with **TinyML** at the edge

3 Process on a secondary larger model on a larger local device

4 Send the data to the cloud when triggered
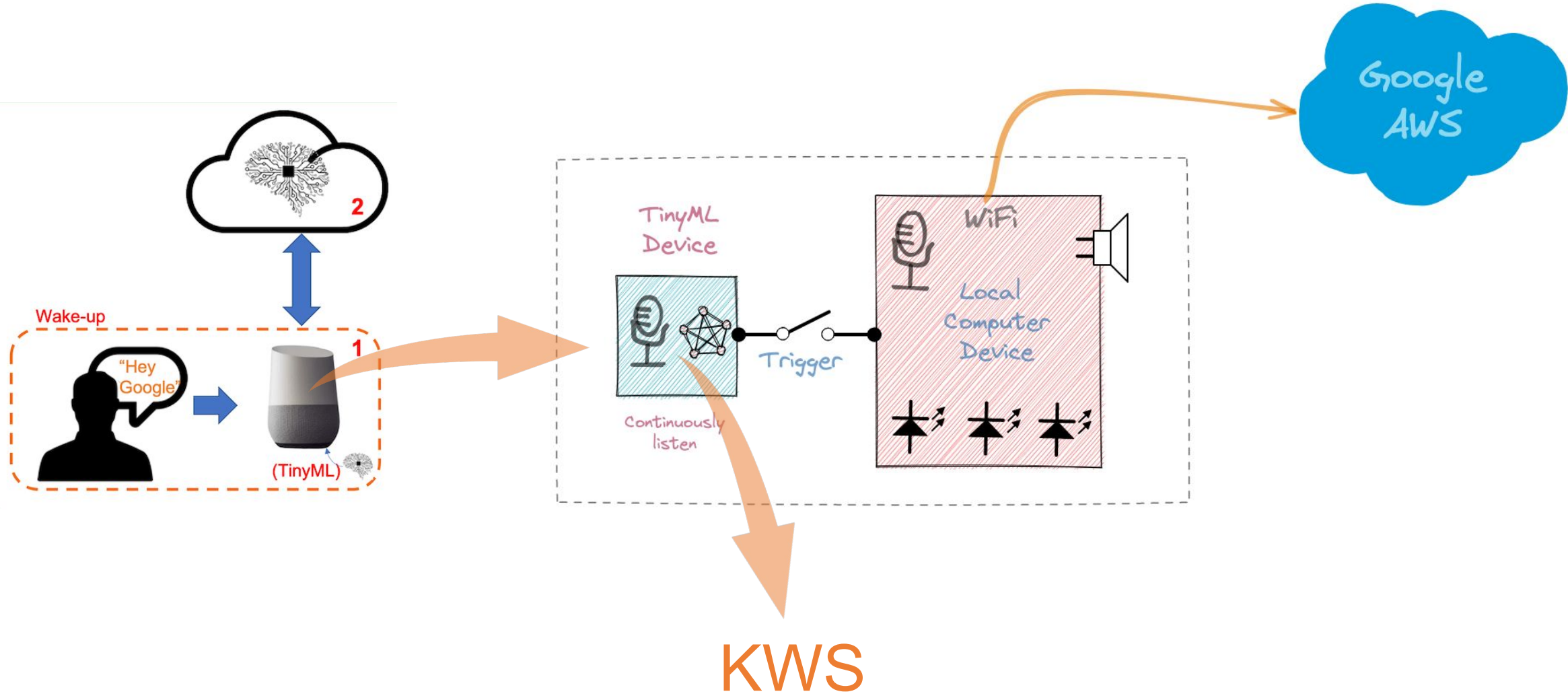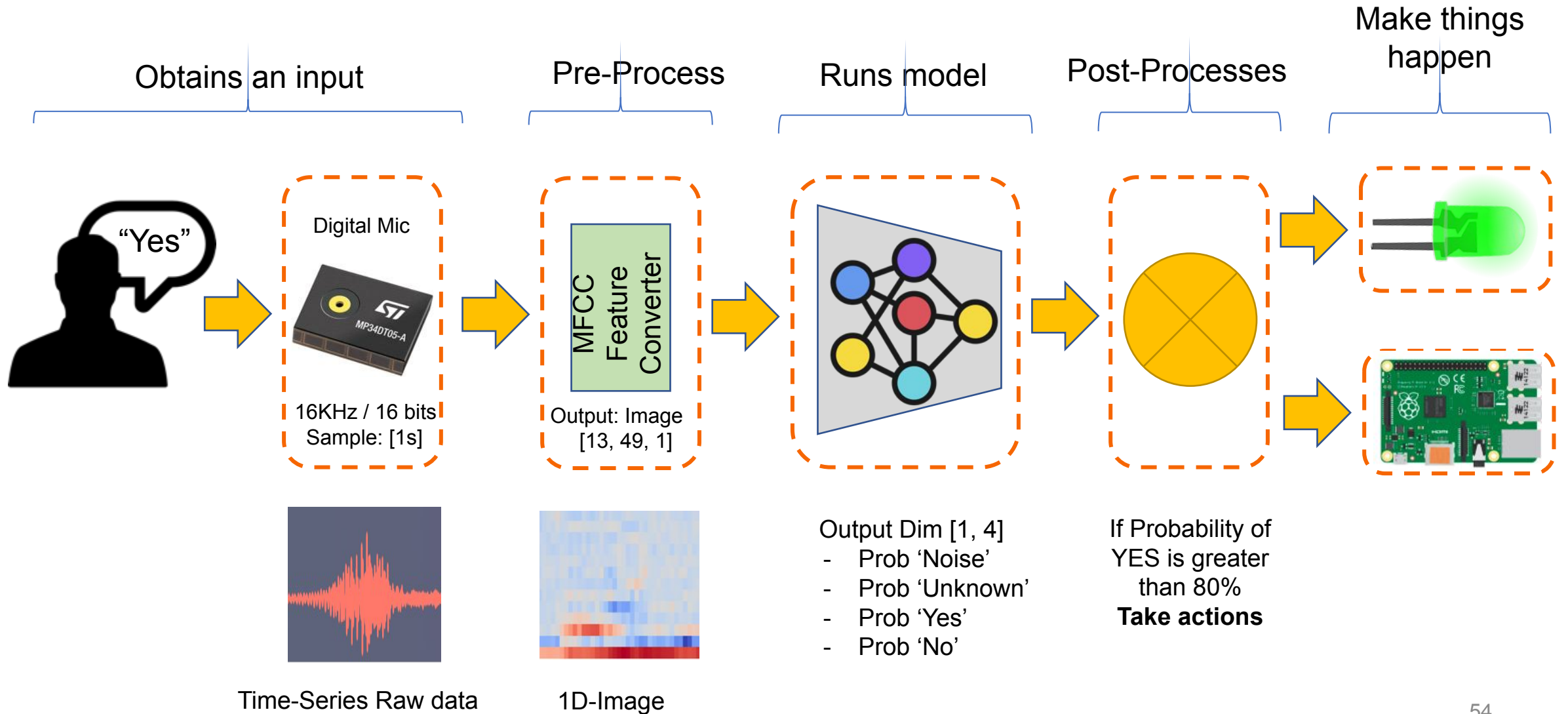
1 Continuously listen on the microcontroller

# Personal Assistant



KWS

# KeyWord Spotting (KWS) - **Inference**



Obtains an input     Pre-Process     Runs model     Post-Processes     Make things happen

"Yes"

**Digital Mic**

16KHz / 16 bits
Sample: [1s]

MFCC Feature Converter

Output: Image
[13, 49, 1]

Output Dim [1, 4]
- Prob 'Noise'
- Prob 'Unknown'
- Prob 'Yes'
- Prob 'No'

If Probability of YES is greater than 80%
**Take actions**

Time-Series Raw data

1D-Image

# Classifying mosquito wingbeat sound using TinyML

Moez Altayeb
University of Khartoum, Sudan
ICTP, Trieste, Italy
mohedahmed@hotmail.com

Marcelo Rovai
Universidade Federal de Itajubá
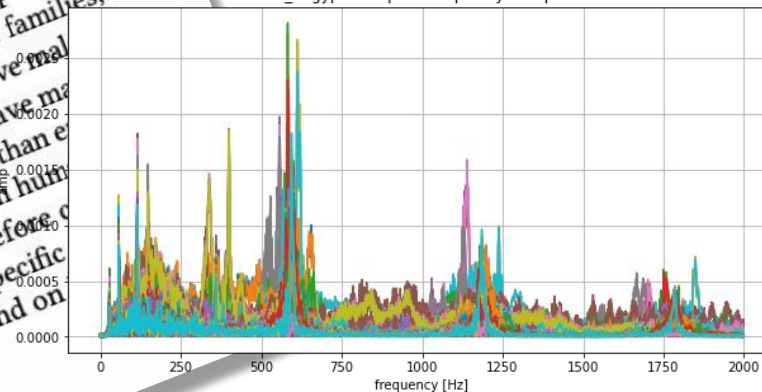Itajubá, Brazil
rovai@unifei.edu.br

Marco Zennaro
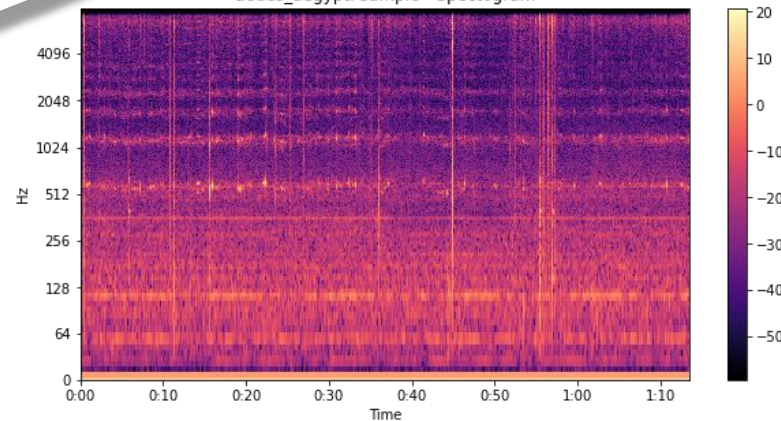ICTP
Trieste, Italy
mzennaro@ictp.it

## ABSTRACT

Every year more than one billion people are infected and more than one million people die from vector-borne diseases including malaria, dengue, zika and chikungunya. Mosquitoes are the best known disease vector and are geographically spread worldwide. It is important to raise awareness of mosquito proliferation by monitoring their incidence, especially in poor regions. Acoustic detection of mosquitoes has been studied for long and ML can be used to automatically identify mosquito species by their wingbeat. We present a prototype solution based on an openly available dataset, on the Edge Impulse platform and on three commercially-available TinyML devices. The proposed solution is low-power, low-cost and can run without human intervention in resource-constrained areas. This insect monitoring system can reach a global scale.

affected. People from poor communities with little access to health care and clean water sources are also at risk. Although anti-malarial drugs exist, there's currently no malaria vaccine. Vector-borne diseases also exacerbate poverty. Illness prevent people from working and supporting themselves and their families, impeding economic development. Countries with intensive mal[...] have much lower income levels than those that don't have ma[...] Countries affected by malaria turn to control rather than e[...] tion. Vector control means decreasing contact between hum[...] disease carriers on an area-by-area basis. It is therefore [...] be able to detect the presence of mosquitoes in a specific [...] paper presents an approach based on TinyML and on [...] embedded devices.
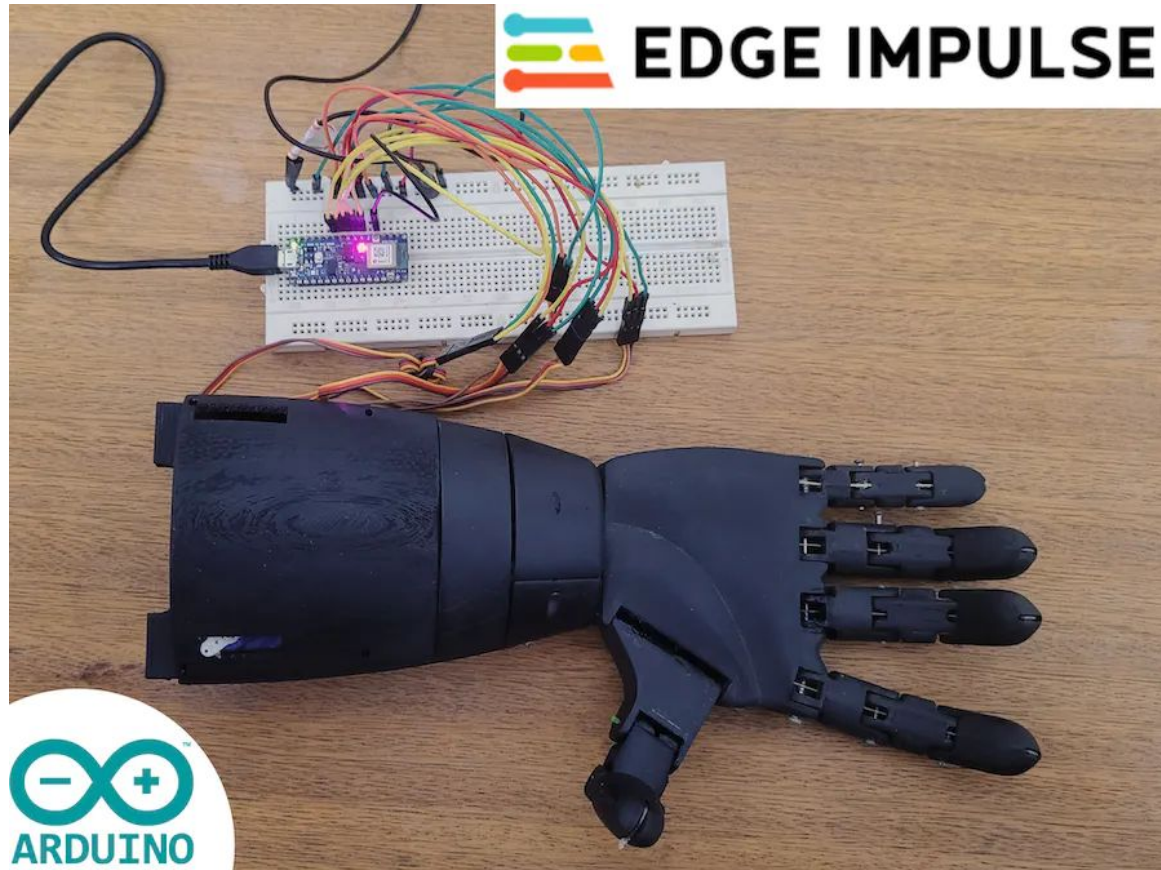
aedes_aegypti sample - Frequency Components
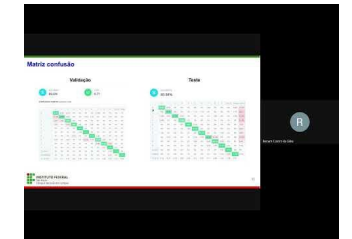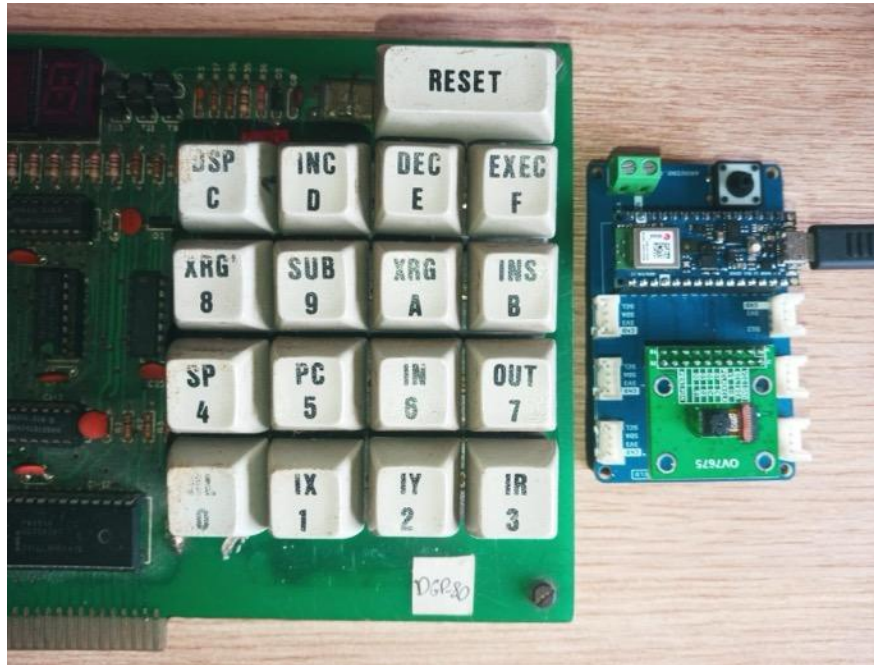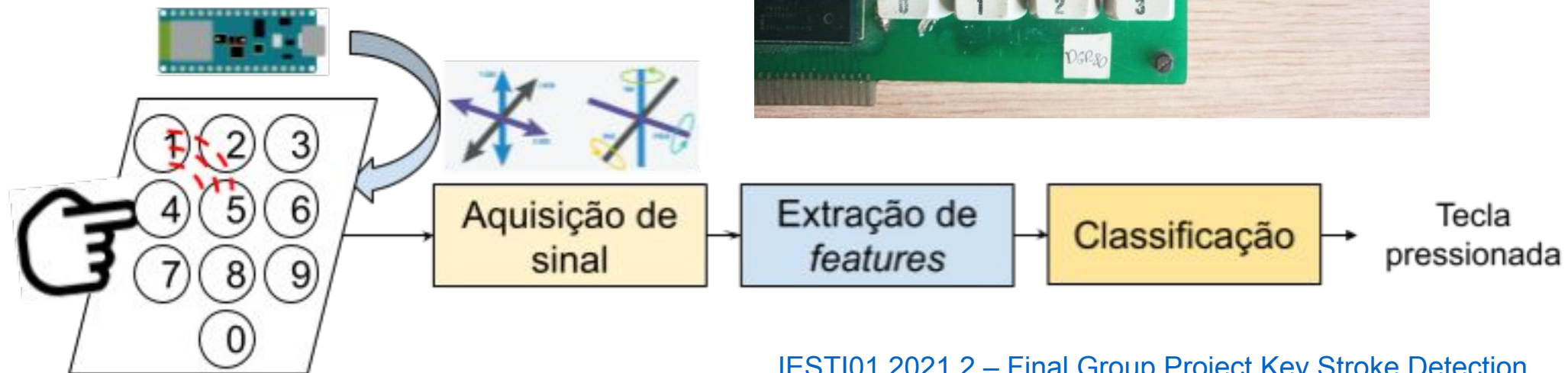


aedes_aegypti sample - Spectogram

55

# Bionic Hand Voice Commands Module



VIDEO

# Keystroke **Sound** Detection



**Renam Castro**
Professor IFESP

Aquisição de sinal → Extração de *features* → Classificação → Tecla pressionada
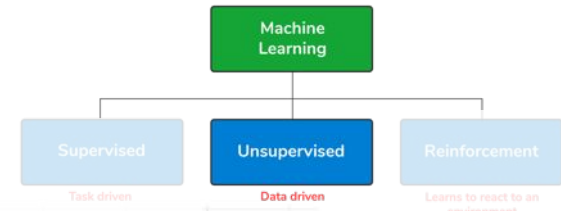
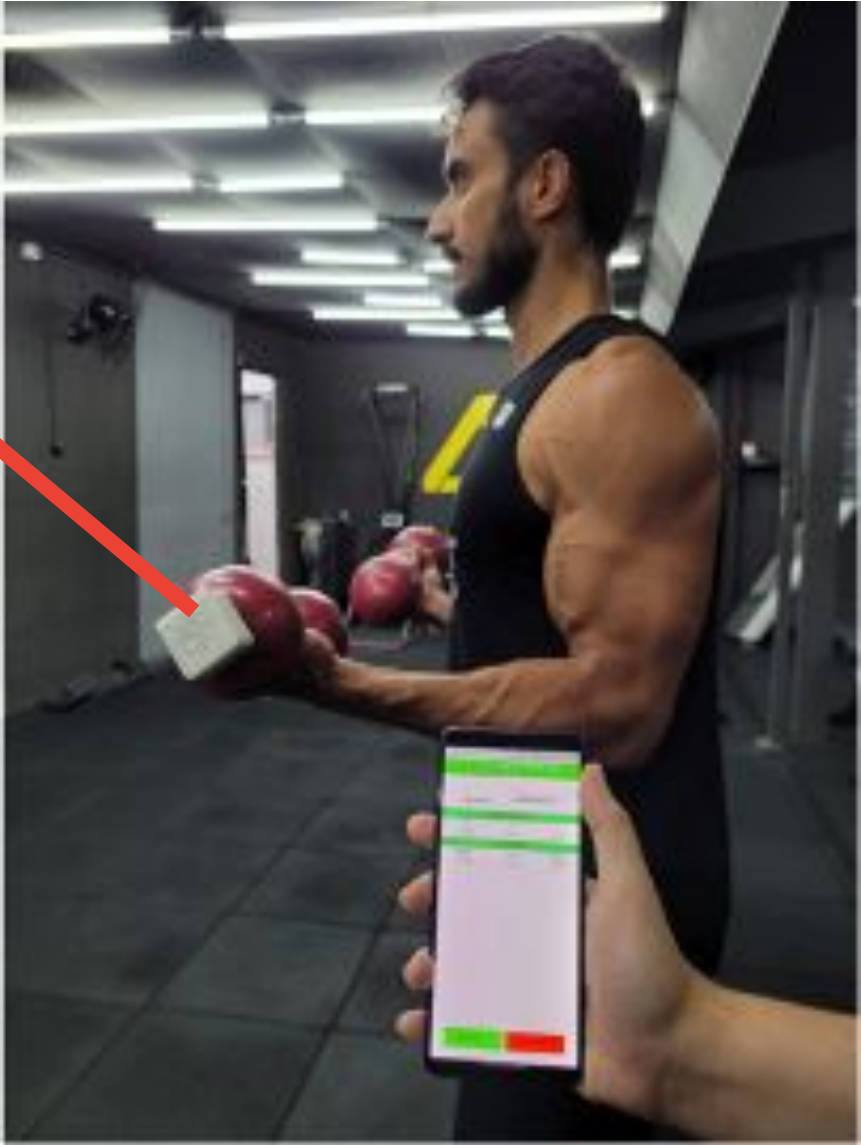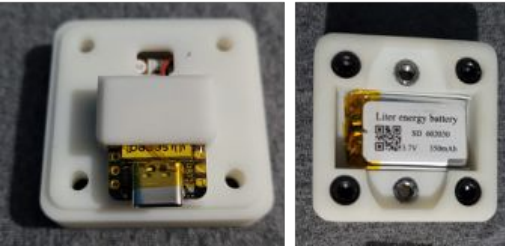IESTI01 2021.2 – Final Group Project Key Stroke Detection
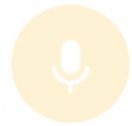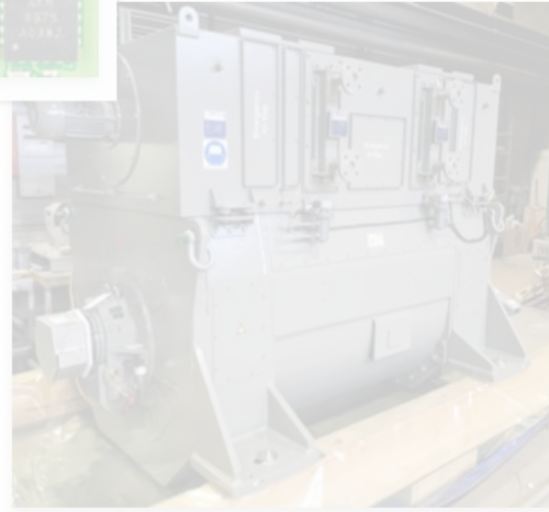
Sound

Vibration

Vision
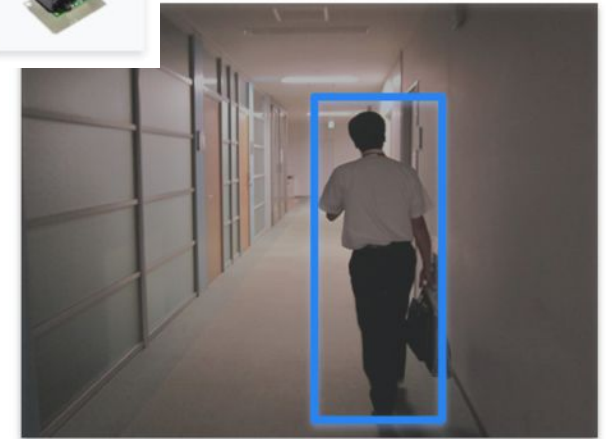
# Industrial – Anomaly Detection

# Movement Classification

Sound

Vibration

Vision

# Computer Vision Main Types
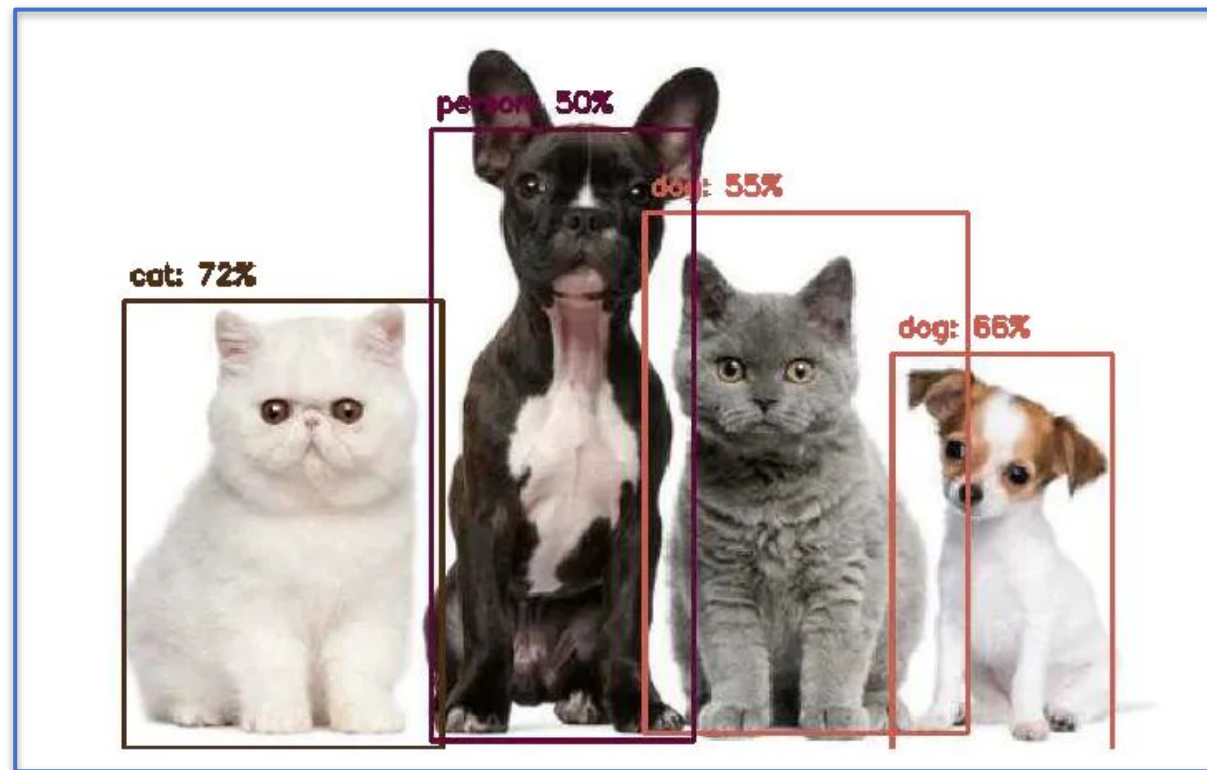
**Image Classification**
(**Multi-Class** Classification)

**Object Detection**
**Multi-Label** Classification + Object **Localization**



Cat: 70%

Dog: 80%



cat: 72%

person: 50%

dog: 55%

dog: 66%

# Computer Vision Main Types

**Image Classification**
(**Multi-Class** Classification)

**Object Detection**
**Multi-Label** Classification + Object **Localization**



Cat: 70%

Dog: 80%



cat: 72%

person 50%

dog: 55%

dog: 66%

# Forest Fire Detection



TinyML Aerial Forest Fire Detection



IESTI01 - Forest Fire Detection – Proof of Concept

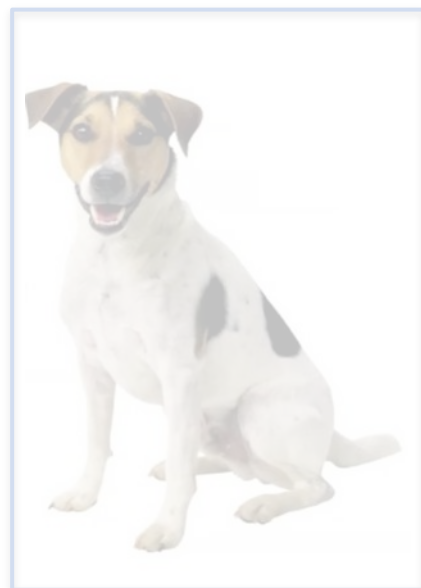# Coffee Disease Classification

**João Vitor Yukio Bordin Yamashita**
Graduando em Engenharia Eletrônica pela UNIFEI

https://www.hackster.io/Yukio/coffee-disease-classification-with-ml-b0a3fc

# Computer Vision Main Types

**Image Classification**
(**Multi-Class** Classification)

**Object Detection**
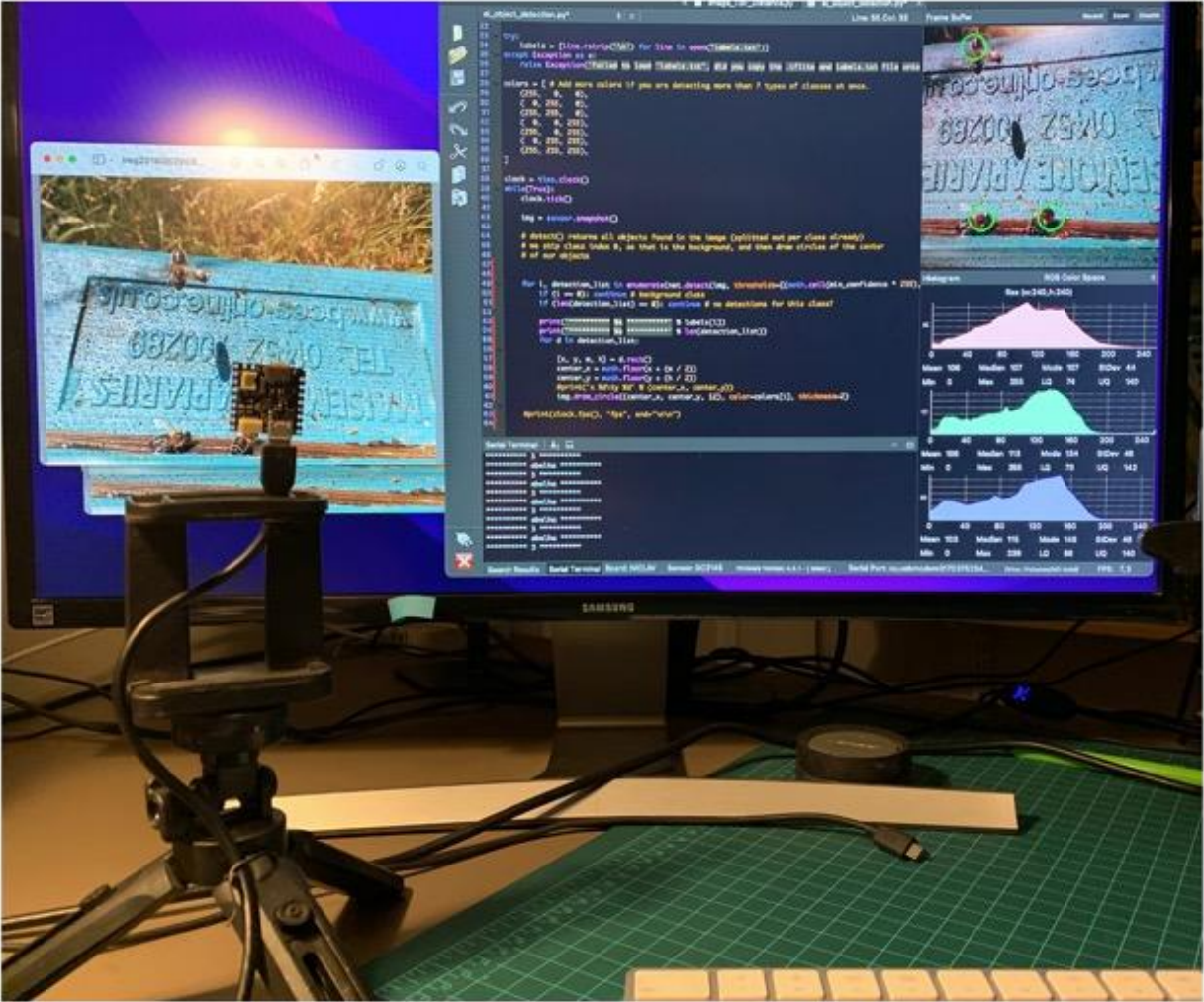**Multi-Label** Classification + Object **Localization**



Cat: 70%

Dog: 80%

# Detecting Objects using TinyML (FOMO)



[EdgeAI made simple - Exploring Image Processing (Object Detection) on microcontrollers with Arduino Portenta, Edge Impulse FOMO, and OpenMV](#)
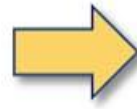
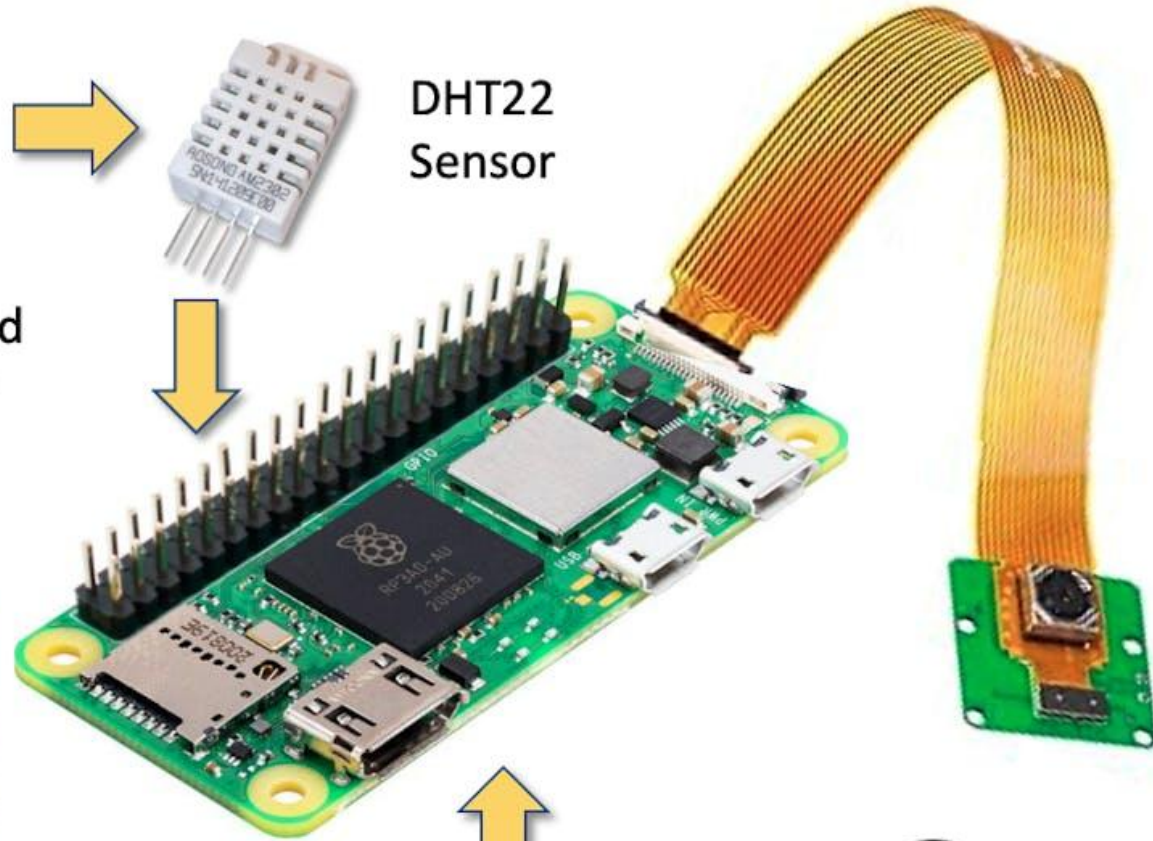# Detecting Objects using TinyML (FOMO)

YOLO

Air Temperature and Relative humidity

DHT22 Sensor

Number of objects: 36 bees

Local Database

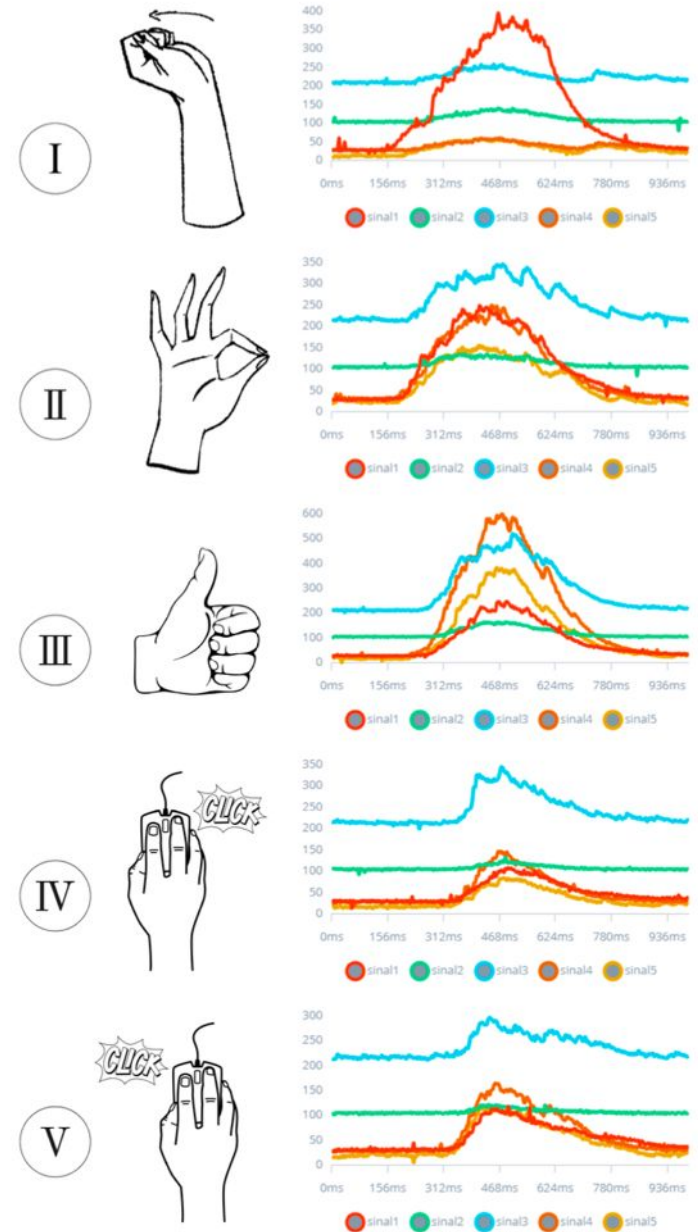sampleFreq ➔ 10 s
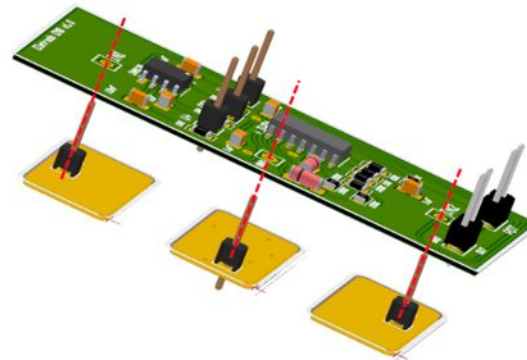
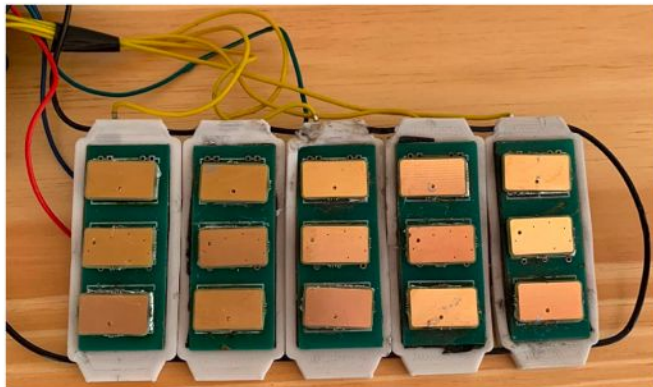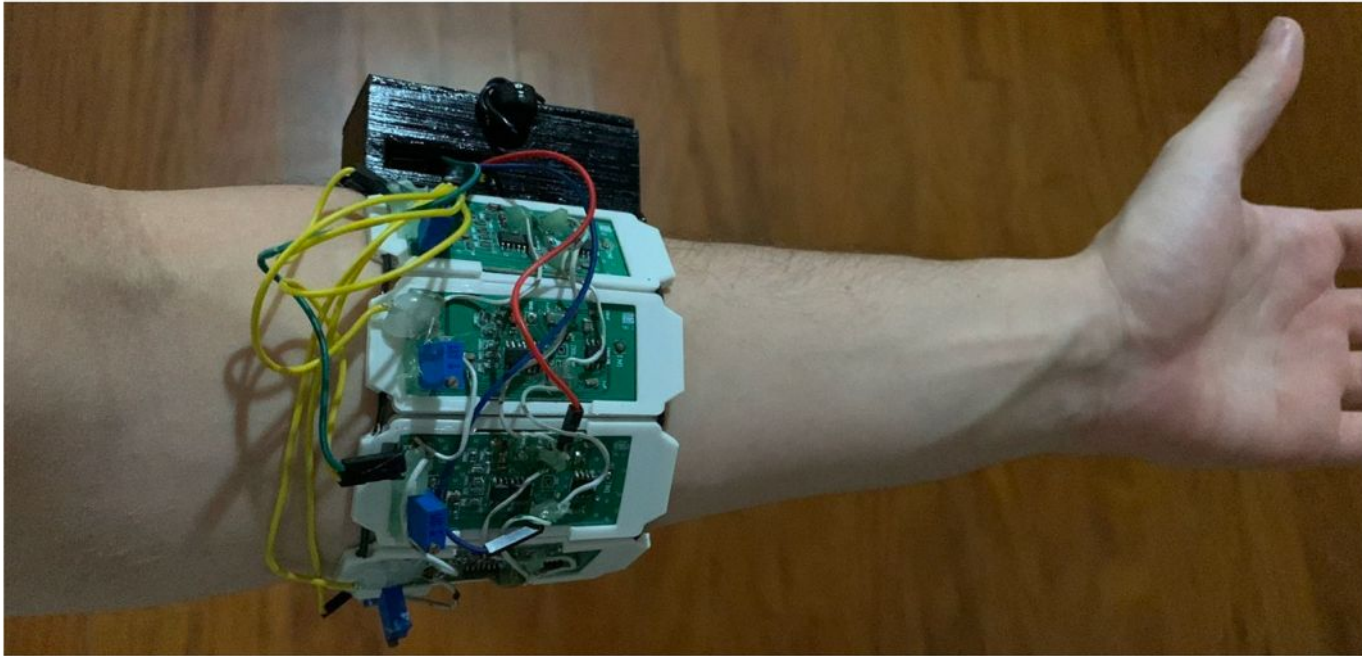BuzzTech: Machine Learning at the Edge
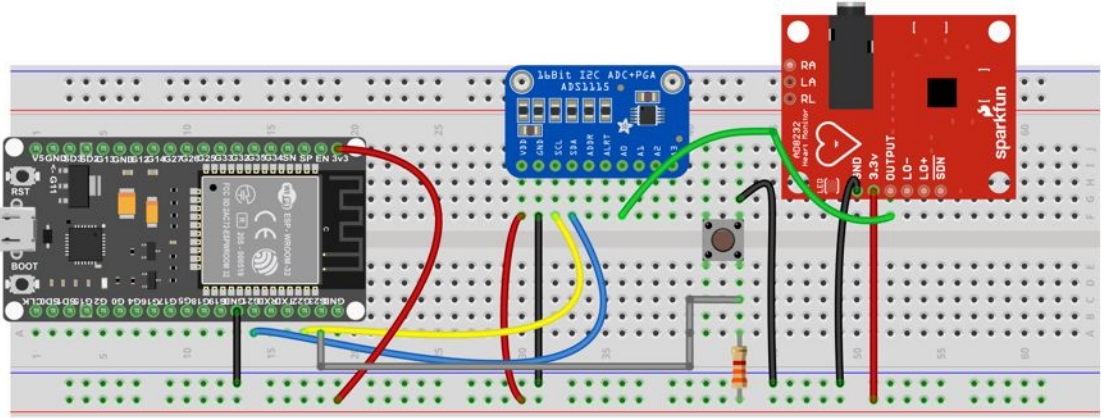
# YOLO

## Ant Detection

# Other Sensors / MCUs / Models

Examples

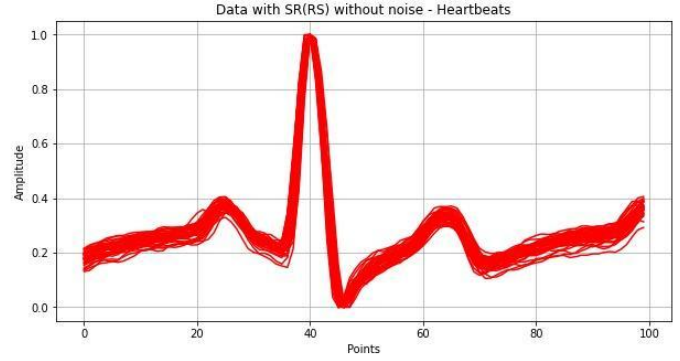# Surface electromyography

# AD8232 - Single Lead Heart Rate Monitor



Atrial Fibrillation Detection on ECG using TinyML
Silva et al. UNIFEI 2021

**Guilherme Silva**
Engenheiro - UNIFEI

# Regression on TinyML



## Sensor fusion

**On-Device IoT-Based Predictive Maintenance Analytics Model: Comparing TinyLSTM and TinyModel from Edge Impulse**

**TinyML Made Easy: Exploring Regression - White Wine Quality**

# LSTM











**ESP32 LSTM Phenolic Sponge Moisture**

# Reinforcement on TinyML

Machine Learning

Supervised — Task driven

Unsupervised — Data driven

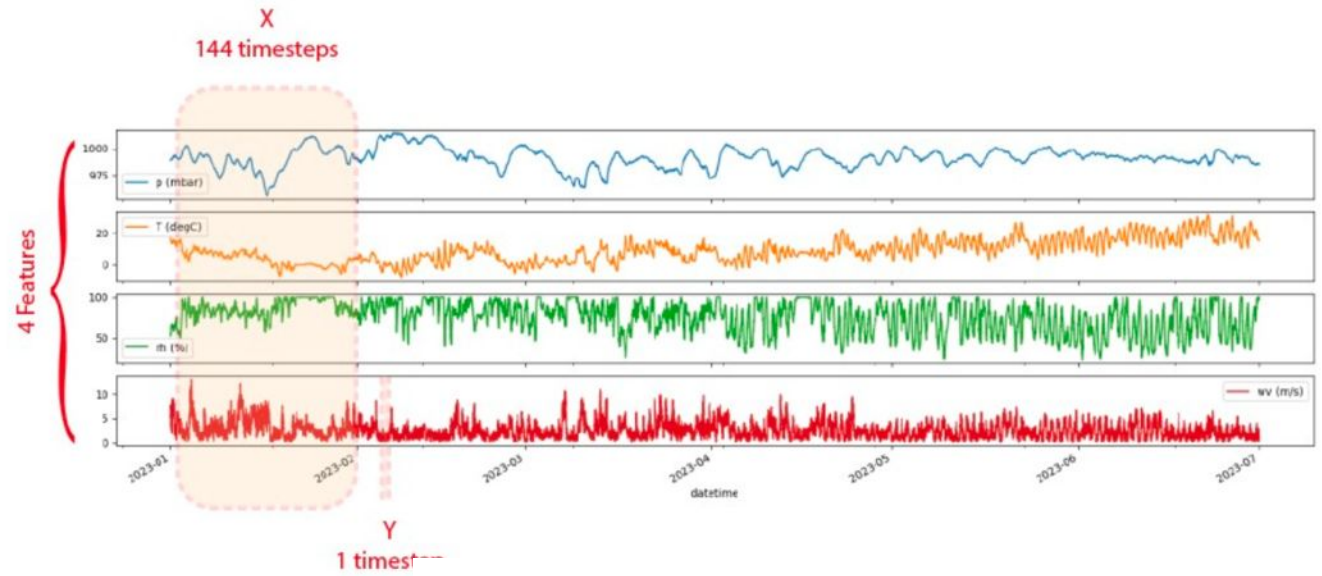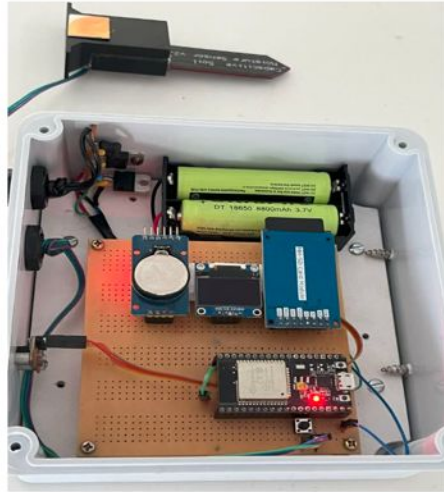Reinforcement — Learns to react to an environment

## Deep Reinforcement Learning for Autonomous Source Seeking on a Nano Drone

Bardienus P. Duisterhof[1,3]  Srivatsan Krishnan[1]  Jonathan J. Cruz[1]  Colby R. Banbury[1]  William Fu[1]

Aleksandra Faust[2]  Guido C. H. E. de Croon[3]  Vijay Janapa Reddi[1,4]

[1]Harvard University, [2]Robotics at Google, [3]Delft University of Technology, [4]The University of Texas at Austin

https://arxiv.org/abs/1909.11236

https://youtu.be/wmVKbX7MOnU

# More MCUs…

## ESP32-TinyML

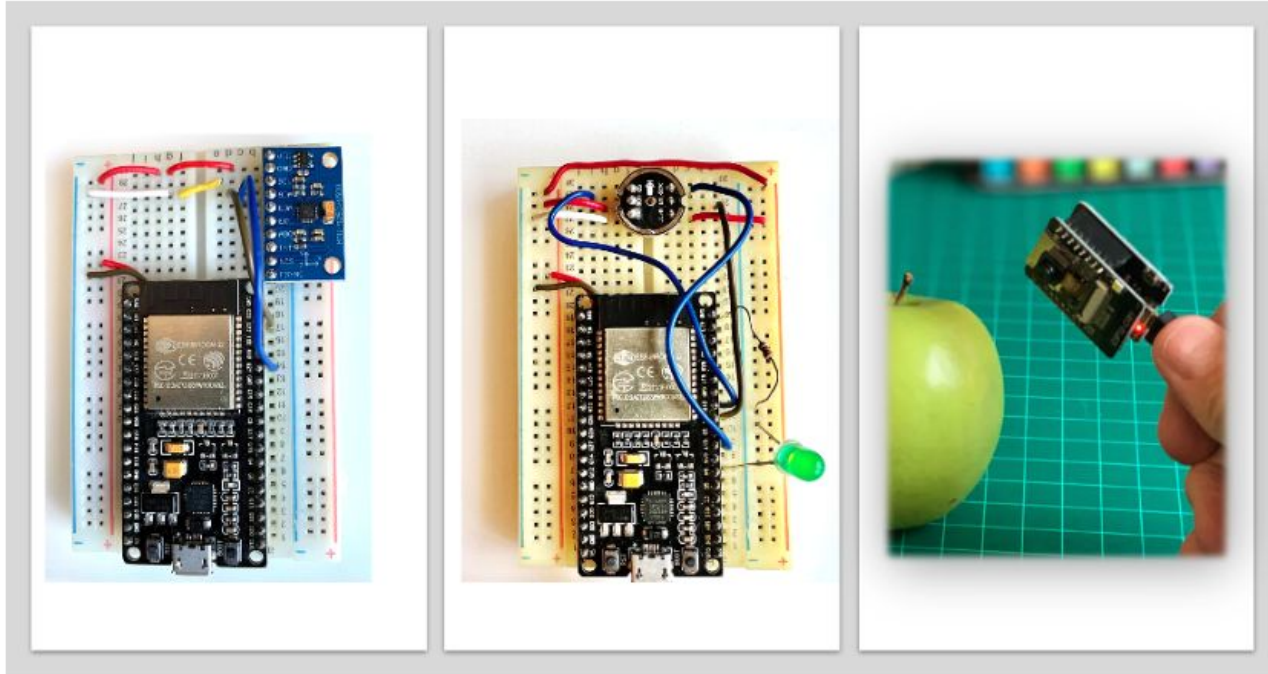Exploring TinyML with ESP32 MCUs.



## Seeed-XIAO-BLE-Sense

KWS, Anomaly Detection & Motion Classification and Micropython - Exploring the Seeed XIAO BLE Sense.



**Programming Tiny devices with MicroPython. The easiest way!**
MJRoBot (Marcelo Rovai)

**Sensor DataLogger**
MJRoBot (Marcelo Rovai)

**TinyML Made Easy: Anomaly Detection & Motion Classification**
MJRoBot (Marcelo Rovai)

**TinyML Made Easy: Sound Classification (KWS)**
MJRoBot (Marcelo Rovai)

## XIAO-ESP32S3-Sense



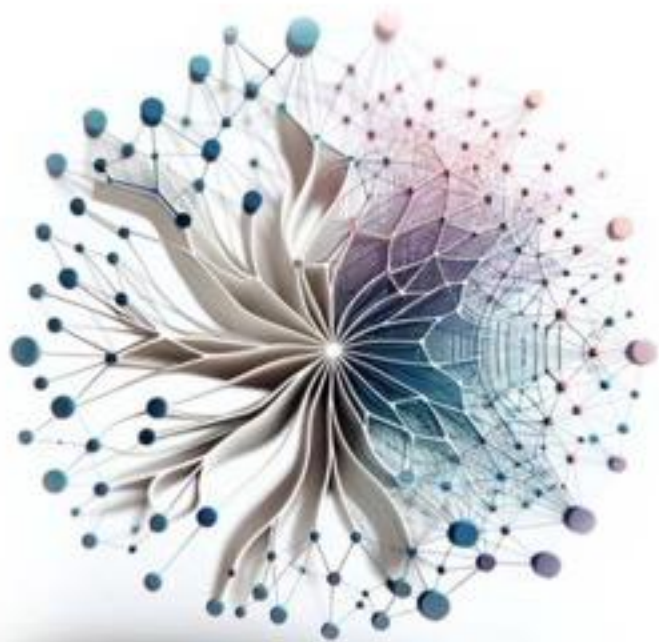**TinyML Made Easy: KeyWord Spotting (KWS)**
MJRoBot (Marcelo Rovai)

**Exploring Machine Learning with the new XIAO ESP32S3**
MJRoBot (Marcelo Rovai)

**TinyML Made Easy: Image Classification**
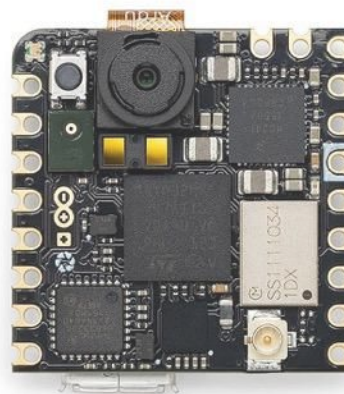MJRoBot (Marcelo Rovai)

Nicla Vision

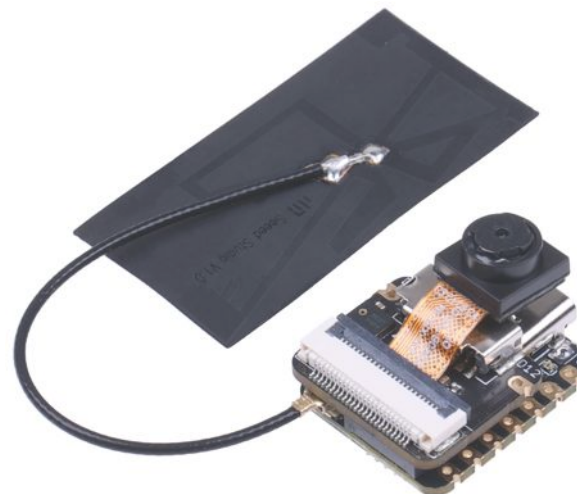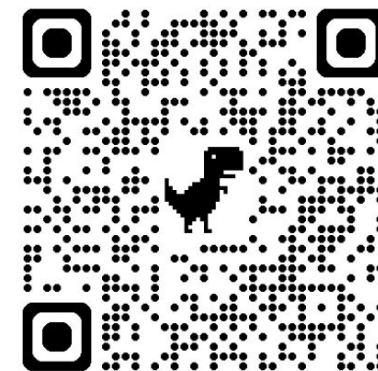XIAO ESP32S3

# Seeed Studio **XIAO**

# To learn more …

**Online Courses**

[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](#)

[UNIFEI-IESTI01 TinyML: "Machine Learning for Embedding Devices"](#)

**Books**

["Python for Data Analysis" by Wes McKinney](#)

["Deep Learning with Python" by François Chollet](#) - [GitHub Notebooks](#)

["TinyML" by Pete Warden and Daniel Situnayake](#)

["TinyML Cookbook 2nd Edition" by Gian Marco Iodice](#)

["Technical Strategy for AI Engineers, In the Era of Deep Learning" by Andrew Ng](#)

["AI at the Edge" book by Daniel Situnayake and Jenny Plunkett](#)

["XIAO: Big Power, Small Board" by Lei Feng and Marcelo Rovai](#)

["MACHINE LEARNING SYSTEMS for TinyML" by a collaborative effort](#)

**Projects Repository**

[Edge Impulse Expert Network](#)

On the **TinyML4D website**, You can find lots of educational materials on TinyML. They are all free and open-source for educational uses – we ask that if you use the material, please cite them! TinyML4D is an initiative to make TinyML education available to everyone globally.

# TinyML4D **Show&Tell** Presentations

**TinymML4D Academic Network Show and Tell Main Index.**

The TinyML4D Academic Network Students should use this form to propose presentations.
**https://forms.gle/ic52HZMqVv4pBrkP7** 2

The Show and Tell are typically held at 2 pm UTC on the last Thursday of each month and will take place in this Meet link:
**https://meet.google.com/rns-yyrx-ggw**

# Conclusion

The Future of ML is Tiny and Bright

**Vijay Janapa Reddi**, *Ph. D. | Associate Professor |*
*John A. Paulson School of Engineering and Applied Sciences | Harvard University |*

# Thanks