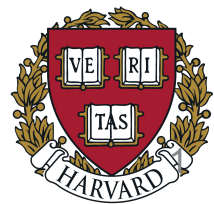# CS249r: Overview and Introduction to TinyML

*Vijay Janapa Reddi, Ph. D. | Associate Professor |*
*John A. Paulson School of Engineering and Applied Sciences | Harvard University |*
*Web: http://scholar.harvard.edu/vijay-janapa-reddi*

# Goals for today

- Introductions
- Course logistics
- Introduction to TinyML

# Introductions

# About Me

- Graduated from Harvard in 2010
- Professor at UT Austin
- Returned to Harvard
- Research
  - Applied machine learning architect
  - Computer architecture and runtime systems
  - Embedded hardware and software

# Teaching Staff

**Matthew Stewart**

Postdoctoral Researcher,
Harvard University

**Ikechukwu Uchendu**

Computer Science, PhD
Student, Harvard University

**Jason Jabbour**

Computer Science, PhD
Student, Harvard University

**Jessica Quaye**

Computer Science, PhD
Student, Harvard University

# Course Logistics

# What you will learn

**Through this course, you would have been exposed to the following:**

- Brief introduction to ML and IoT
- Industry talks about real-world deployments
- Discussion of bleeding edge academic research
- Practical experience through hands-on project assignments

**At the end of the course, you would have achieved the following:**

- Gained familiarity with cutting edge literature in the field of tinyML
- Learnt to train and deploy models on microcontrollers with TF-micro
- Conceived and developed a (novel) TinyML application running on a MCU

# Course Topics

1. **Overview and Introduction to Embedded Machine Learning**
2. Data Engineering
3. Embedded Machine Learning Frameworks
4. Efficient Model Representation and Compression
5. Performance Metrics and Benchmarking of ML Systems
6. Learning on the Edge
7. Hardware Acceleration for Edge ML: GPUs, TPUs and FPGAs
8. Embedded MLOps
9. Secure and Privacy-Preserving On-Device ML
10. Responsible AI
11. Sustainability at the Edge
12. Generative AI at the Edge

# Grading

| | |
|---|---|
| Class Participation | = 10% |
| Paper Reviews | = 10% |
| Paper Presentation | = 10% |
| Programming Assignments | = 25% |
| Final Project | = 45% |

# Paper Readings

- Check to see the **bold** list of papers, these are the ones to read and come prepared for discussion
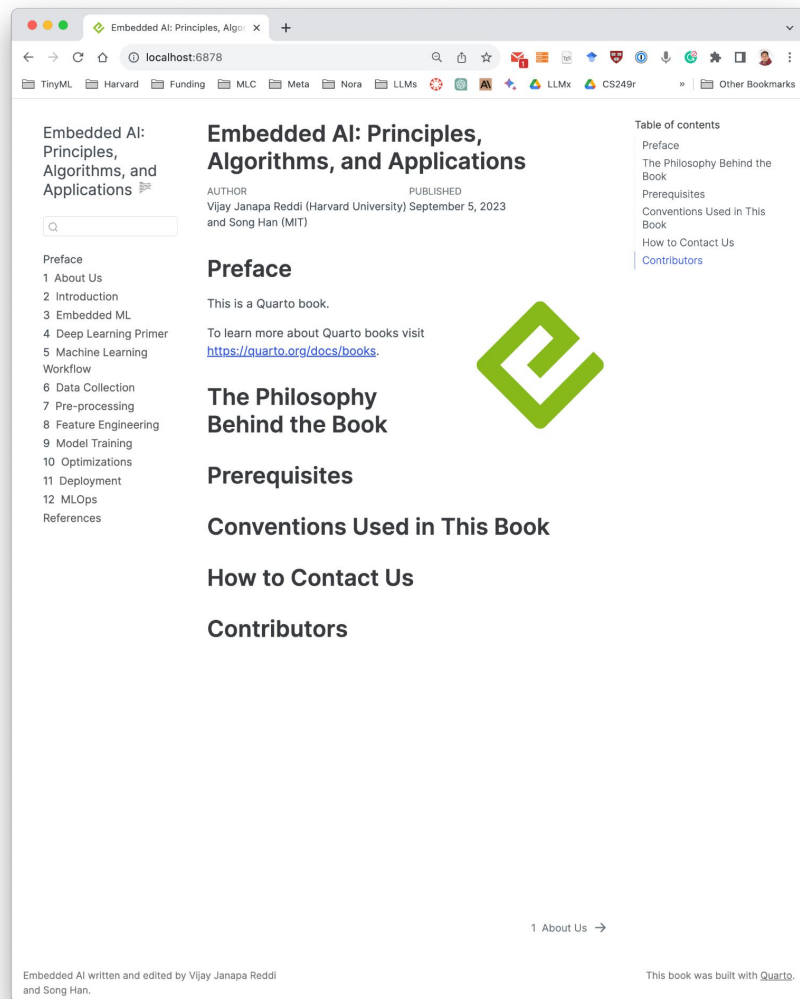- Discuss in groups
- One lead per paper
- Report back to the group

# Paper Reviews & Presentations

Goal: Share the knowledge

- Listen to the lectures
- Read the papers
- Review the material
- Scribe into the open source book

Instructions will be announced soon!

# Class Schedule

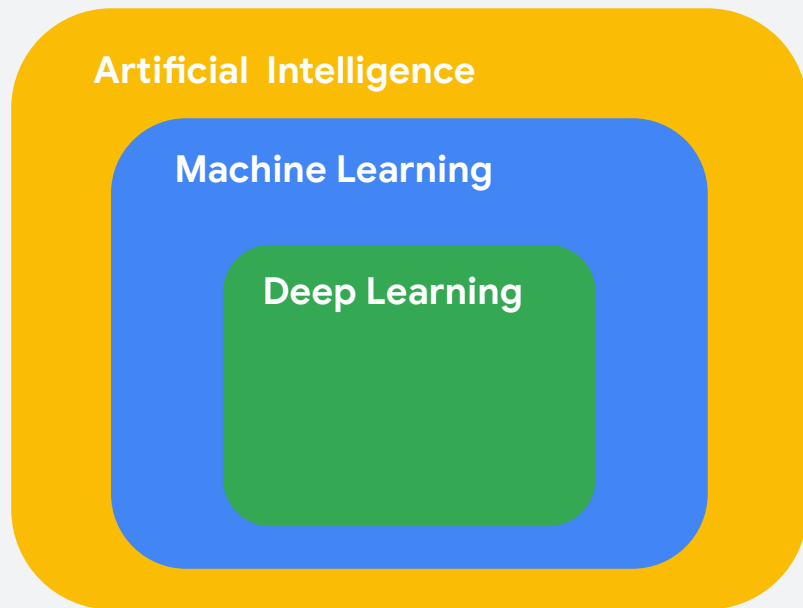| | | Start | ~End |
|---|---|---|---|
| Lecture | | 12:45:00 PM | 1:30:00 PM |
| Break | | 1:40:00 PM | 1:45:00 PM |
| Paper Discussions | Paper 1 | 1:45:00 PM | 2:05:00 PM |
| | Paper 2 | 2:05:00 PM | 2:25:00 PM |
| Break | | 2:25:00 PM | 2:30:00 PM |
| Guest Lecture | | 2:30:00 PM | 3:30:00 PM |

# Introduction to TinyML

# What is Machine Learning?

1. **Machine Learning** is a subfield of **Artificial Intelligence** focused on developing algorithms that learn to **solve problems by analyzing data for patterns**

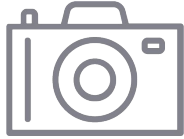**Artificial  Intelligence**

**Machine Learning**

# What is (Deep) Machine Learning?

1. Machine Learning is a subfield of Artificial Intelligence focused on developing algorithms that learn to solve problems by analyzing data for patterns

2. **Deep Learning** is a type of Machine Learning that leverages **Neural Networks** and **Big Data**



Artificial Intelligence

Machine Learning

Deep Learning

15

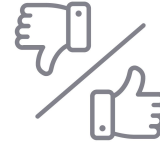# Applications of Machine Learning

# Applications of Machine Learning

Label

# Applications of Machine Learning

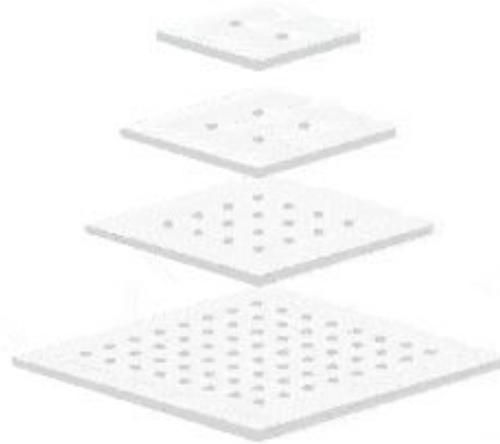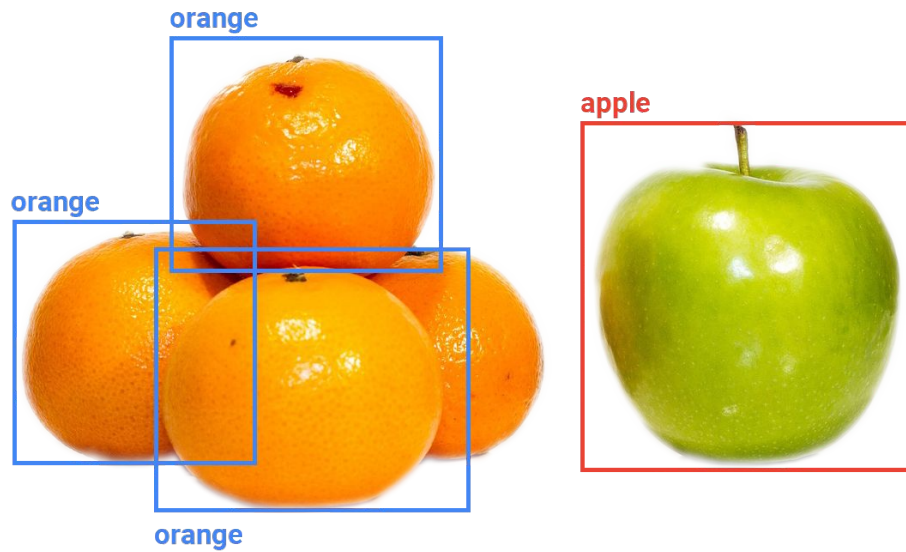# Image Classification

# Object Detection

# Segmentation

# Machine Translation

# Recommendations

# Datacenter

# Datacenter

# AI Compute: 300,000x Increase in Demand

"… **since 2012 the amount of compute used in the largest AI training runs has been increasing exponentially with a 3.5 month-doubling time** (by comparison, Moore's Law had an 18-month doubling period). Since 2012, this metric has **grown by more than 300,000x (an 18-month [Moore's Law] doubling period would yield only a 12x increase)**. Improvements in compute have been a key component of AI progress, so as long as this trend continues, it's worth preparing for the implications of systems far outside today's capabilities."

**AlexNet to AlphaGo Zero: A 300,000x Increase in Compute**



Source: https://blog.openai.com/ai-and-compute/

26

# Two Eras of Computing

"… **since 2012 the amount of compute used in the largest AI training runs has been increasing exponentially with a 3.5 month-doubling time** (by comparison, Moore's Law had an 18-month doubling period). Since 2012, this metric has **grown by more than 300,000x (an 18-month [Moore's Law] doubling period would yield only a 12x increase)**. Improvements in compute have been a key component of AI progress, so as long as this trend continues, it's worth preparing for the implications of systems far outside today's capabilities."



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

Source: https://blog.openai.com/ai-and-compute/

# TPUs/GPUs

But... Bigger Is Not Always Better.

# Common carbon footprint benchmarks

in lbs of CO2 equivalent

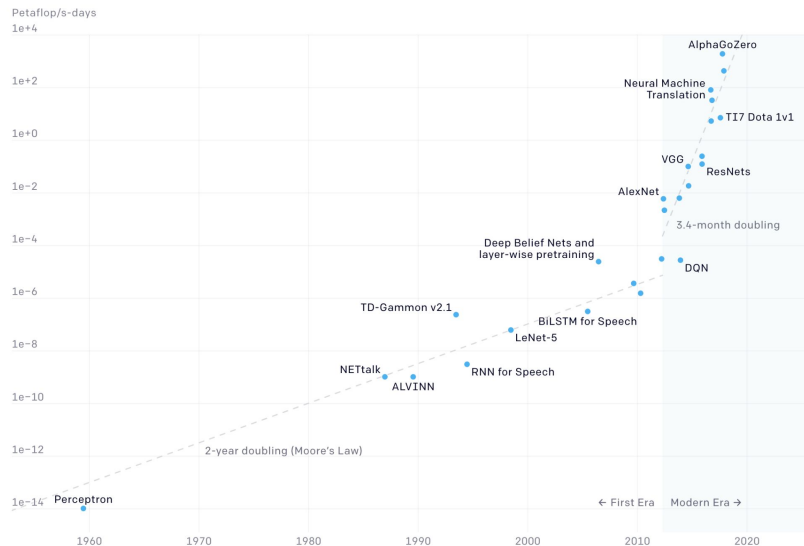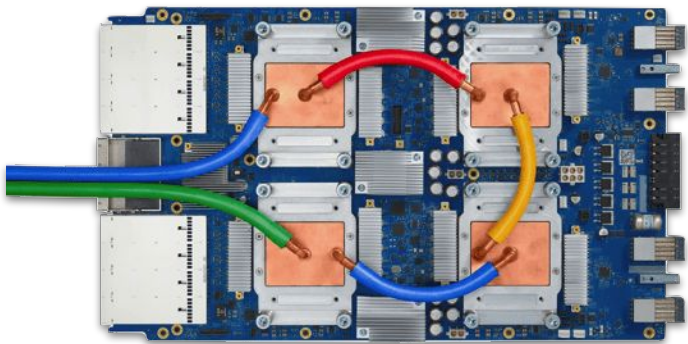| | |
|---|---|
| Roundtrip flight b/w NY and SF (1 passenger) | 1,984 |
| Human life (avg. 1 year) | 11,023 |
| American life (avg. 1 year) | 36,156 |
| US car including fuel (avg. 1 lifetime) | 126,000 |
| Transformer (213M parameters) w/ neural architecture search | 626,155 |

---

## Energy and Policy Considerations for Deep Learning in NLP

Emma Strubell    Ananya Ganesh    Andrew McCallum
College of Information and Computer Sciences
University of Massachusetts Amherst
{strubell, aganesh, mccallum}@cs.umass.edu

arXiv:1906.02243v1 [cs.CL] 5 Jun 2019

### Abstract

Recent progress in hardware and methodology for training neural networks has ushered in a new generation of large networks trained on abundant data. These models have obtained notable gains in accuracy across many NLP tasks. However, these accuracy improvements depend on the availability of exceptionally large computational resources that necessitate similarly substantial energy consumption. As a result these models are costly to train and develop, both financially, due to the cost of hardware and electricity or compute time, and environmentally, due to the carbon footprint required to fuel modern tensor processing hardware. In this paper we bring this issue to the attention of NLP researchers by quantifying the approximate financial and environmental costs of training a variety of recently successful neural network models for NLP. Based on these findings, we propose actionable recommendations to reduce costs and improve equity in NLP research and practice.

### 1 Introduction

Advances in techniques and hardware for training deep neural networks have recently enabled impressive accuracy improvements across many fundamental NLP tasks (Bahdanau et al., 2015; Luong et al., 2015; Dozat and Manning, 2017; Vaswani et al., 2017), with the most computationally-hungry models obtaining the highest scores (Peters et al., 2018; Devlin et al., 2019; Radford et al., 2019; So et al., 2019). As a result, training a state-of-the-art model now requires substantial computational resources which demand considerable energy, along with the associated financial and environmental costs. Research and development of new models multiplies these costs by thousands of times by requiring retraining to experiment with model architectures and hyperparameters. Whereas a decade ago most
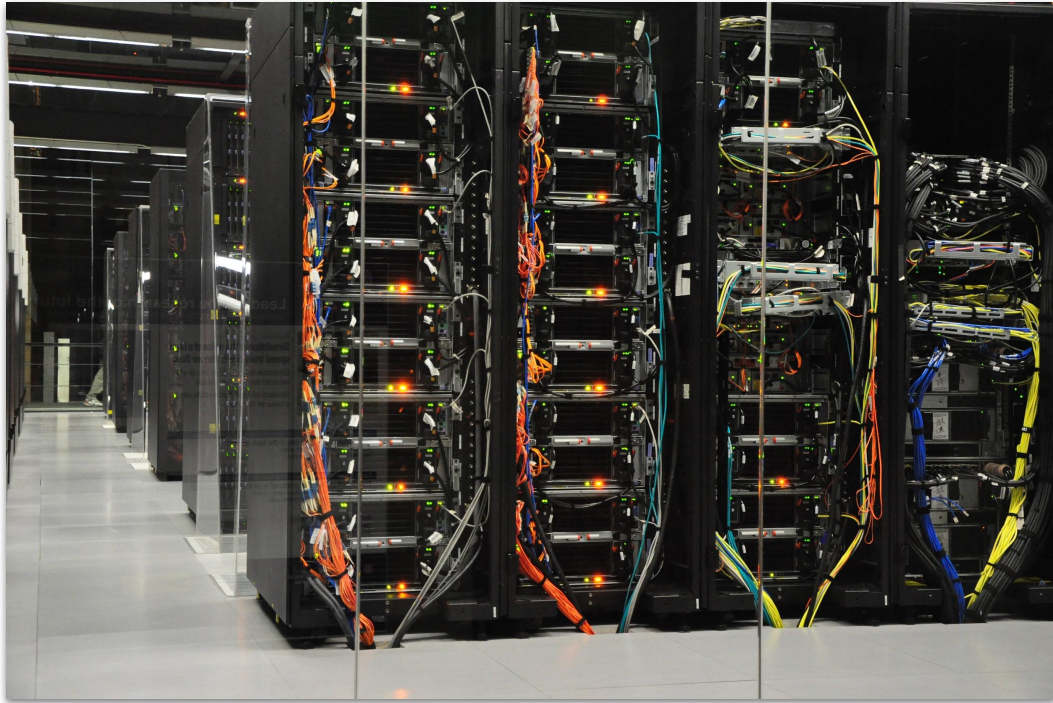
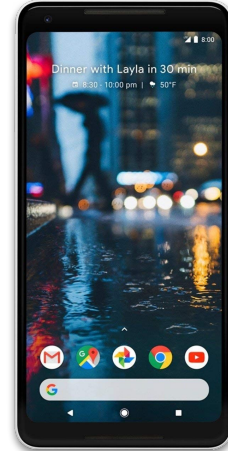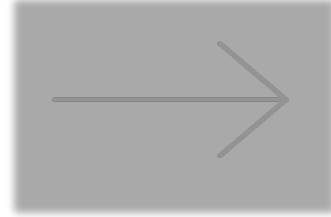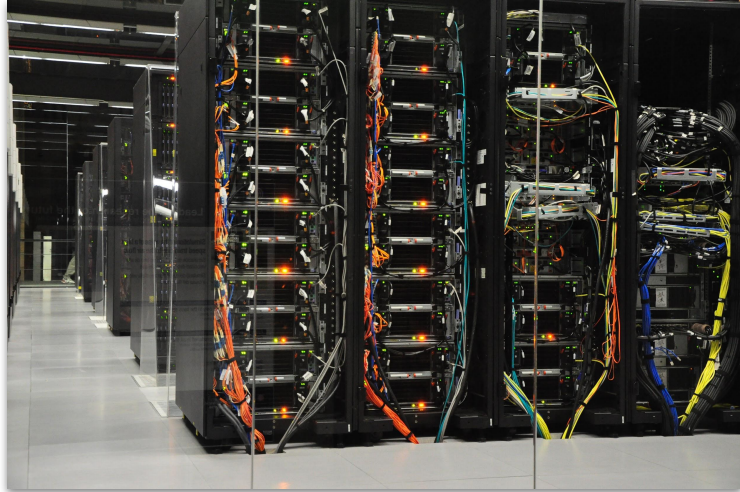| Consumption | CO$_2$e (lbs) |
|---|---|
| Air travel, 1 passenger, NY↔SF | 1984 |
| Human life, avg, 1 year | 11,023 |
| American life, avg, 1 year | 36,156 |
| Car, avg incl. fuel, 1 lifetime | 126,000 |
| **Training one model (GPU)** | |
| NLP pipeline (parsing, SRL) | 39 |
| w/ tuning & experimentation | 78,468 |
| Transformer (big) | 192 |
| w/ neural architecture search | 626,155 |

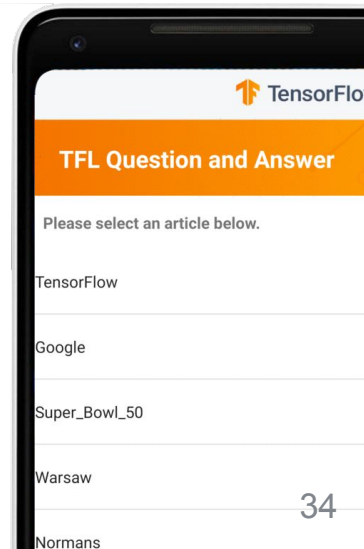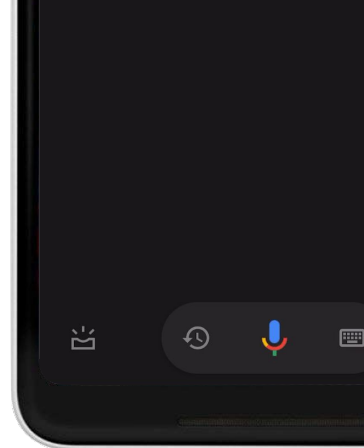Table 1: Estimated CO$_2$ emissions from training common NLP models, compared to familiar consumption.[1]

NLP models could be trained and developed on a commodity laptop or server, many now require multiple instances of specialized hardware such as GPUs or TPUs, therefore limiting access to these highly accurate models on the basis of finances.
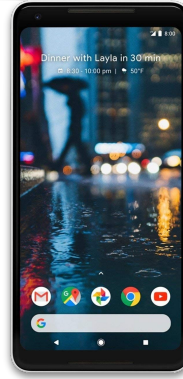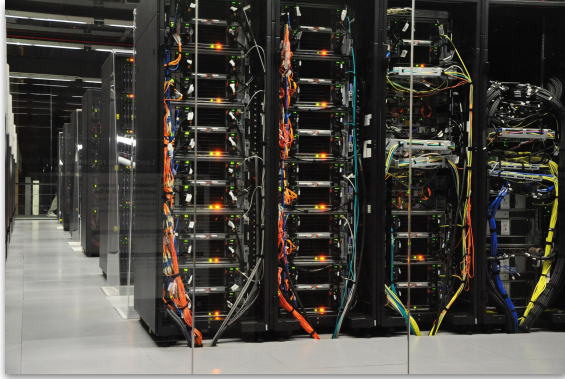
Even when these expensive computational resources are available, model training also incurs a substantial cost to the environment due to the energy required to power this hardware for weeks or months at a time. Though some of this energy may come from renewable or carbon credit-offset resources, the high energy demands of these models are still a concern since (1) energy is not currently derived from carbon-neutral sources in many locations, and (2) when renewable energy is available, it is still limited to the equipment we have to produce and store it, and energy spent training a neural network might better be allocated to heating a family's home. It is estimated that we must cut carbon emissions by half over the next decade to deter escalating rates of natural disaster, and based on the estimated CO$_2$ emissions listed in Table 1,

[1]Sources: (1) Air travel and per-capita consumption: https://bit.ly/2Hw0xWc; (2) car lifetime: https://bit.ly/2Qbr0w1.

1

# TensorFlow

**TFL Question and Answer**

Please select an article below.

TensorFlow

Google

Super_Bowl_50

Warsaw

Normans

Kicking

Penalty kicking

Passing

Dribbling

...

# No Good Data Left Behind

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

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

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

# What is TinyML?

# TinyML

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

**TinyML**

# What is Tiny Machine Learning (**TinyML**)?
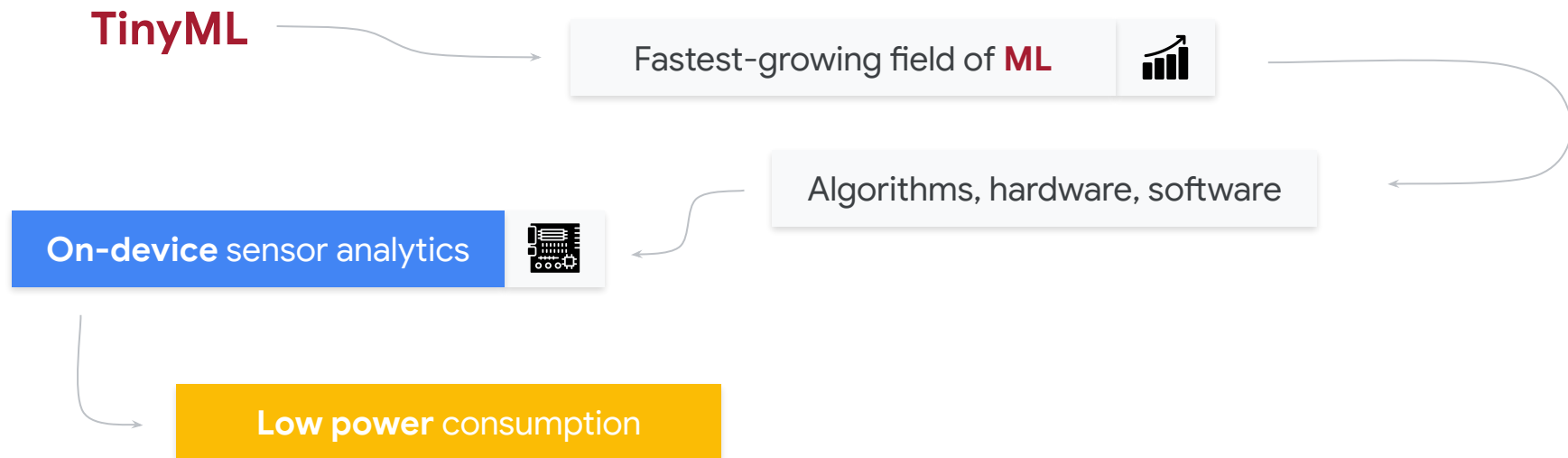
**TinyML**  Fastest-growing field of **ML** 📊
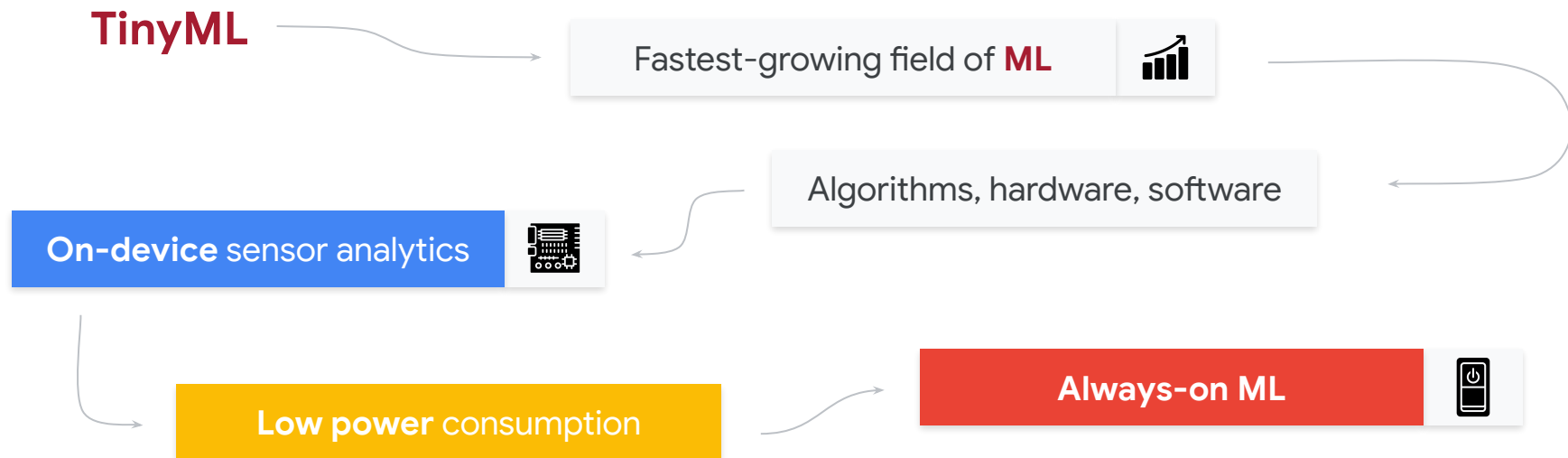
# What is Tiny Machine Learning (**TinyML**)?

**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

**Always-on ML**

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

**TinyML**

Fastest-growing field of **ML**

Algorithms, hardware, software

**On-device** sensor analytics

**Low power** consumption

**Always-on ML**

**Battery-operated**

# Tiny Robot Learning





Usually, it is very challenging to get such tiny drones to fly autonomously.

# Tiny Robot Learning





Usually, it is very challenging to get such tiny drones to fly autonomously.

# Wildlife Conservation



© Elephant Listening Project

# ElephantEdge

**Risk Monitoring**

"Know when an elephant is moving into a high-risk area and send real-time notifications to park rangers."

**Conflict Monitoring**

"Sense and alert when an elephant is heading into an area where farmers live."

**Activity Monitoring**

"Classify the general behavior of the elephant, such as when it is drinking, eating, sleeping, etc."

**Communication Monitoring**

"Listen for vocal communications between elephants via the onboard microphone."

# Rich Array of Sensors

**Motion Sensors**
Gyroscope, radar, magnetometer, accelerator

**Acoustic Sensors**
Ultrasonic, Microphones, Geophones, Vibrometers

**Environmental Sensors**
Temperature, Humidity, Pressure, IR, etc.

**Touchscreen Sensors**
Capacitive, IR

**Image Sensors**
Thermal, Image

**Biometric Sensors**
Fingerprint, Heart rate, etc.

**Force Sensors**
Pressure, Strain

**Rotation Sensors**
Encoders

...

A DIGITAL

AN ELECTROMECHANICAL

"SMART" METER

ANALOG METER

# Digitizer - AI on the edge

## An ESP32 all inclusive neural network recognition system for meter digitalization

| Raw Value: |
|---|
| 038.5975 |
| **Corrected Value:** |
| 38.5975 |
| **Checked Value:** |
| 38.5975 |
| **Start Time:** |
| 20201118-075416 |
| Last Page Refresh:06:57:39 |

# Challenges

# 250 Billion
*MCUs today*

| Board | MCU / ASIC | Clock | Memory | Sensors | Radio |
|-------|-----------|-------|--------|---------|-------|
| Himax<br>WE-I Plus EVB | HX6537-A<br>32-bit EM9D DSP | 400 MHz | 2MB flash<br>2MB RAM | Accelerometer, Mic, Camera | None |
| Arduino<br>Nano 33 BLE Sense | 32-bit<br>nRF52840 | 64 MHz | 1MB flash<br>256kB RAM | Mic, IMU, Temp, Humidity, Gesture, Pressure, Proximity, Brightness, Color | BLE |
| SparkFun<br>Edge 2 | 32-bit<br>ArtemisV1 | 48 MHz | 1MB flash<br>384kB RAM | Accelerometer, Mic, Camera | BLE |
| Espressif<br>EYE | 32-bit<br>ESP32-D0WD | 240 MHz | 4MB flash<br>520kB RAM | Mic, Camera | WiFi, BLE |

**Challenges**

```
                        ┌──────────────┐
                        │  Challenges  │
                        └──────────────┘

        ┌──────────────┐              ┌──────────────┐
        │   Hardware   │              │   Software   │
        └──────────────┘              └──────────────┘

 ┌──────────────┐ ┌────────────────┐  ┌──────────────┐ ┌──────────────────┐
 │Heterogeneity │ │Resource        │  │Missing Library│ │Limited Operating │
 │              │ │Constraints     │  │Features       │ │System Support    │
 └──────────────┘ └────────────────┘  └──────────────┘ └──────────────────┘
```
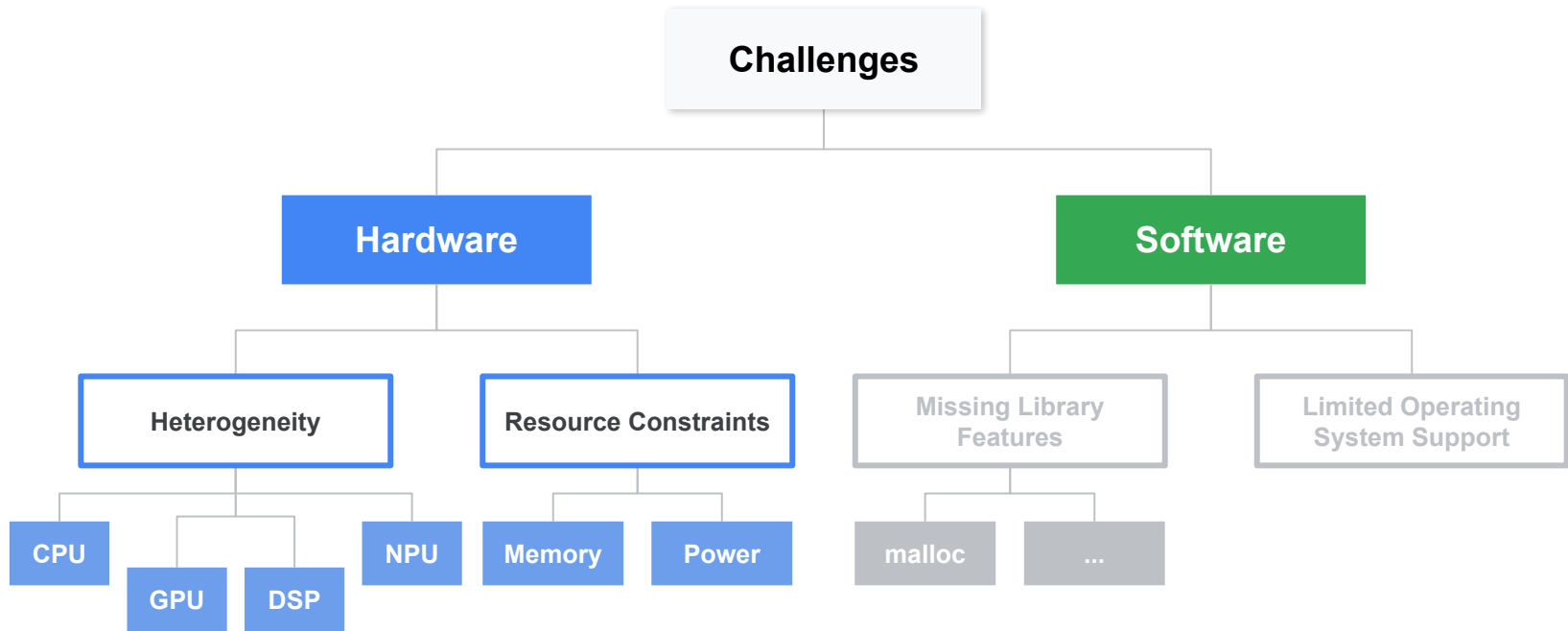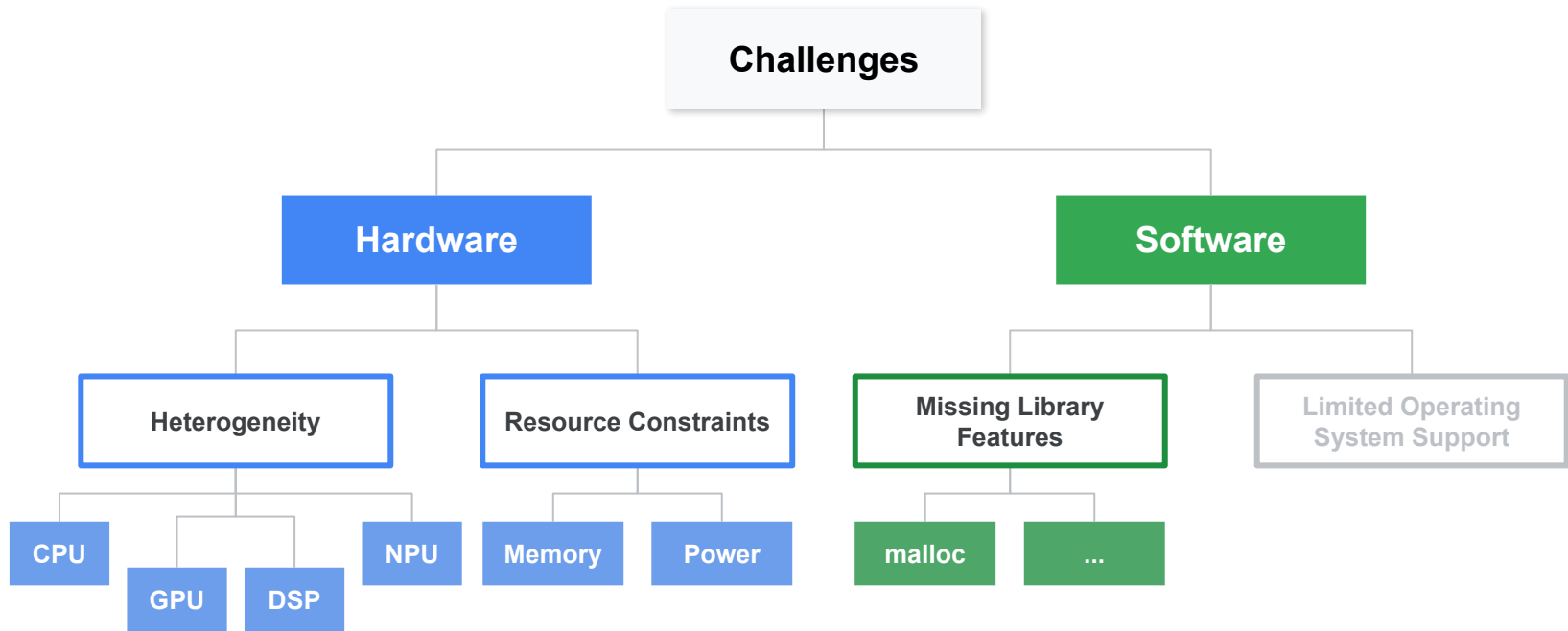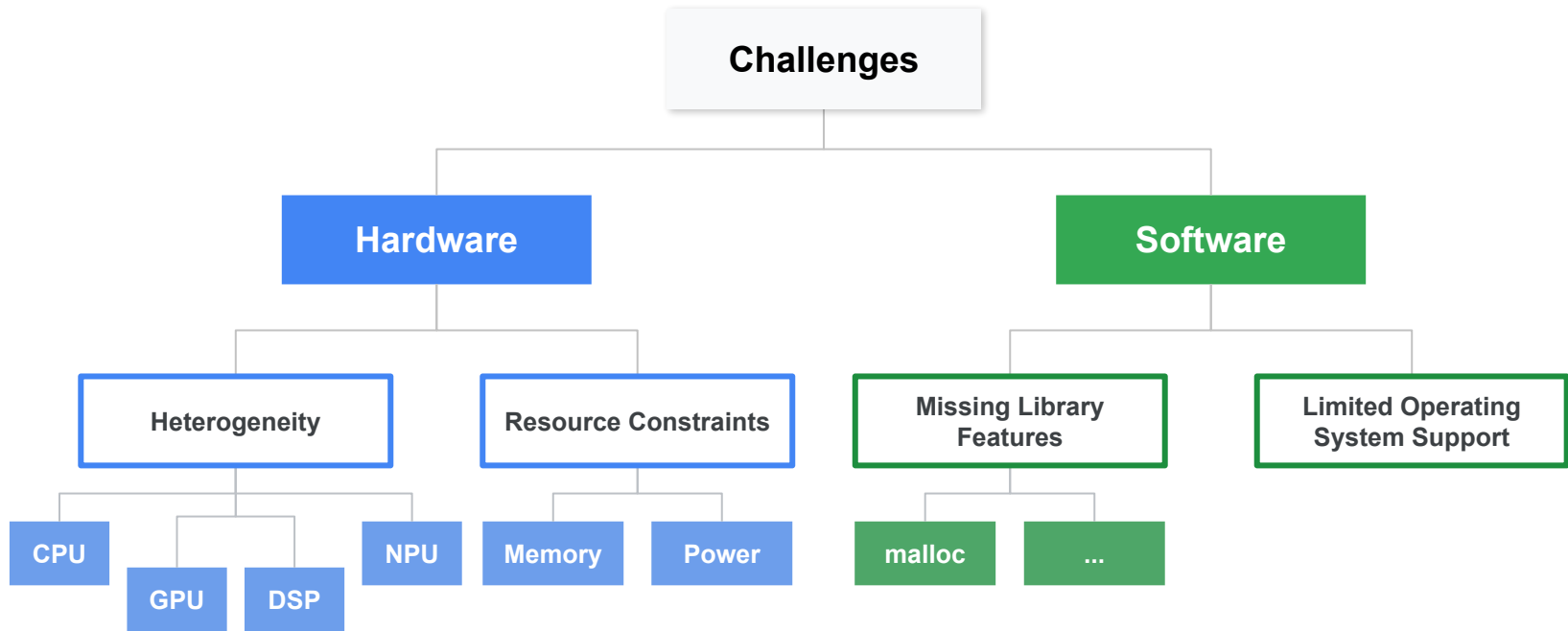
**Challenges**

**Hardware**  **Software**

**Heterogeneity**  **Resource Constraints**  **Missing Library Features**  **Limited Operating System Support**

**CPU**  **NPU**  **Memory**  **Power**  **malloc**  **...**

**GPU**  **DSP**

```
                            Challenges
                    ┌───────────┴───────────┐
                Hardware                  Software
            ┌───────┴───────┐         ┌───────┴───────┐
      Heterogeneity   Resource      Missing Library   Limited Operating
                      Constraints   Features          System Support
    ┌───┬──┴──┬───┐    ┌───┴───┐      ┌───┴───┐
   CPU GPU DSP NPU  Memory Power   malloc  ...
```

**Challenges**

**Hardware**  **Software**

Heterogeneity  Resource Constraints  Missing Library Features  Limited Operating System Support

CPU  GPU  DSP  NPU  Memory  Power  malloc  ...

Challenges

Hardware

Software

Heterogeneity

Resource Constraints

Missing Library Features

Limited Operating System Support

CPU

GPU

DSP

NPU

Memory

Power

malloc

...

TensorFlow Lite Micro

...

Arduino
BLE Sense 33

Himax
WE-I Plus EVB

SparkFun
Edge 2

Espressif
EYE

...

Training

TensorFlow Lite

.tflite

Training

Conversion

Training

Conversion

TensorFlow Lite

.tflite

Array modeling

C array models

Training

Conversion

Inference
Learning

Real Time Data

Array modeling

TensorFlow Lite

.tflite

C array models

# TensorFlow Lite Micro in a Nutshell

Built to fit on **embedded systems**:

- Very **small binary footprint**
- **No** dynamic memory allocation
- **No** dependencies on complex parts of the standard C/C++ libraries
- **No** operating system dependencies, **can run on bare metal**
- Designed to be **portable** across a wide variety of systems



**[MLSys'21]**

# A Greener Tomorrow with TinyML

# Sustainable TinyML

Climate change

Water demand

Freshwater eutrophication

Protochemical oxidant formation

# Sustainable TinyML



| Climate change | Water demand | Freshwater eutrophication | Protochemical oxidant formation |
| --- | --- | --- | --- |
| Total impact 390 g CO2-eq. | Total impact 23 L | Total impact 120 mg P-eq. | Total impact 823 mg NMVOC |
| (a) | (b) | (c) | (d) |

Raw materials · ST production site · Transport · Use · End-of-life

# Tiny Footprint of a Microcontroller



**Total Impact**

| | 390g $CO_2$-eq 1.6km by car | 23L 23 bottles of water | 120mg P-eq 0.2 washing cycles | 823mg NMVOC 2.7km by car |
|---|---|---|---|---|
| | Climate Change | Water Demand | Freshwater Eutrophication | Protochemical Oxidant Formation |
| End of Life | <1% | <1% | <1% | <1% |
| Logistics | 1% | <1% | <1% | 1% |
| Use | 8% | 6% | 28% | 8% |
| Raw Materials | 10% | 41% | 27% | 10% |
| Production: Other | 24% | 15% | 18% | 2% |
| Production: Energy Consumption | 56% | 39% | 27% | 71% |

# Common Carbon Footprint Benchmarks

in lbs of CO2 equivalent

| Benchmark | Value |
|---|---|
| Roundtrip flight b/w NY and SF (1 passenger) | 1,984 |
| Human life (avg. 1 year) | 11,023 |
| American life (avg. 1 year) | 36,156 |
| US car including fuel (avg. 1 lifetime) | 126,000 |
| Transformer (213M parameters) w/ neural architecture search | 626,155 |

Chart: MIT Technology Review • Source: Strubell et al. • Created with Datawrapper

Logistics Costs

Energy    Raw Materials

Microcontrollers ← Manufacturing

Design

Sustainable Development Goals

TinyML Cycle

Deploy    Maintain

Sustainable Development Goals icons:
2 ZERO HUNGER
6 CLEAN WATER AND SANITATION
12 RESPONSIBLE CONSUMPTION AND PRODUCTION
13 CLIMATE ACTION
14 LIFE BELOW WATER
15 LIFE ON LAND

Sustainable Development Goals

Logistics Costs

Energy — Raw Materials

Microcontrollers ← Manufacturing

Design

TinyML Cycle

Deploy

Maintain

End-of-life

Energy Consumption

E-waste

2 ZERO HUNGER
6 CLEAN WATER AND SANITATION
12 RESPONSIBLE CONSUMPTION AND PRODUCTION
13 CLIMATE ACTION
14 LIFE BELOW WATER
15 LIFE ON LAND

Sustainable Development Goals / TinyML Cycle

Logistics Costs

Microcontrollers

Design

TinyML Cycle

Energy Consumption

Deploy

Maintain

Energy    Raw Materials

Manufacturing

Upcycling

Recycling

End-of-life

E-waste

2 ZERO HUNGER

6 CLEAN WATER AND SANITATION

12 RESPONSIBLE CONSUMPTION AND PRODUCTION

13 CLIMATE ACTION

14 LIFE BELOW WATER

15 LIFE ON LAND

Energy

Raw Materials

Manufacturing

Cloud Hardware Microcontrollers

Logistics Costs

Recycling

Upcycling

End-of-life

E-waste

Applications & Consumer Use

TEST
DESIGN
ML
MAINTAIN
DEPLOY
Ops
TRAIN

Sustainable Agriculture        Conservation

Disease Prevention        **Benefits**        Clean Air & Water

Low Cost        Low Power

Storage        Communication

**Constraints**

Compute        Connectivity

Energy Consumption

# The MLOps Process

# ML Development

ML development entails experimenting with and establishing a dependable and repeatable model training procedure.



Data & model management

**ML development**

Training operationalization

Continuous training

Model deployment

Prediction Serving

Continuous monitoring

# Training Operationalization

Training operationalization is all about automating the packaging, testing, and deployment of repeatable and dependable training pipelines.

| Data & model management | ML development |
| | **Training operationalization** |
| | Continuous training |
| | Model deployment |
| | Prediction Serving |
| | Continuous monitoring |

# Continuous Training

Continuous training entails running the training pipeline on a regular basis, maybe with fresh training settings, in response to new data or code modifications.

**Data & model management**

ML development

Training operationalization

**Continuous training**

Model deployment

Prediction Serving

Continuous monitoring

# Model Deployment

Packaging, testing, and deploying a model to a serving environment for online experimentation and production serving is what model deployment is all about.

# Prediction Serving

Serving the model that is deployed in production for inference is known as prediction serving.



Data & model management

ML development

Training operationalization

Continuous training

Model deployment

**Prediction Serving**

Continuous monitoring

# Continuous Monitoring

Continuous monitoring refers to keeping track of a deployed model's effectiveness and efficiency.



Data & model management

ML development

Training operationalization

Continuous training

Model deployment

Prediction Serving

**Continuous monitoring**

# Data & Model Management

Data and model management is a central, cross-cutting function for governing ML artifacts to support ability, traceability, and compliance. Data and model management can also promote shareability, reusability, and discoverability of ML assets.

**Data & model management**

- ML development
- Training operationalization
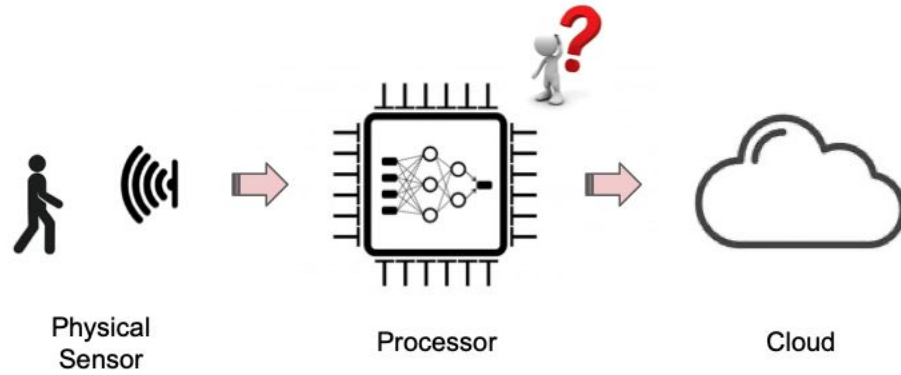- Continuous training
- Model deployment
- Prediction Serving
- Continuous monitoring

# ML Sensors

Physical Sensor — Processor — Cloud

**Sensor 1.0**

Sensor 1.0

Physical Sensor — Processor — Cloud

Sensor 2.0

Machine Learning (ML) Sensor — Processor (MCU) — Cloud

Sensor 1.0

Physical Sensor — Processor — Cloud

Sensor 2.0

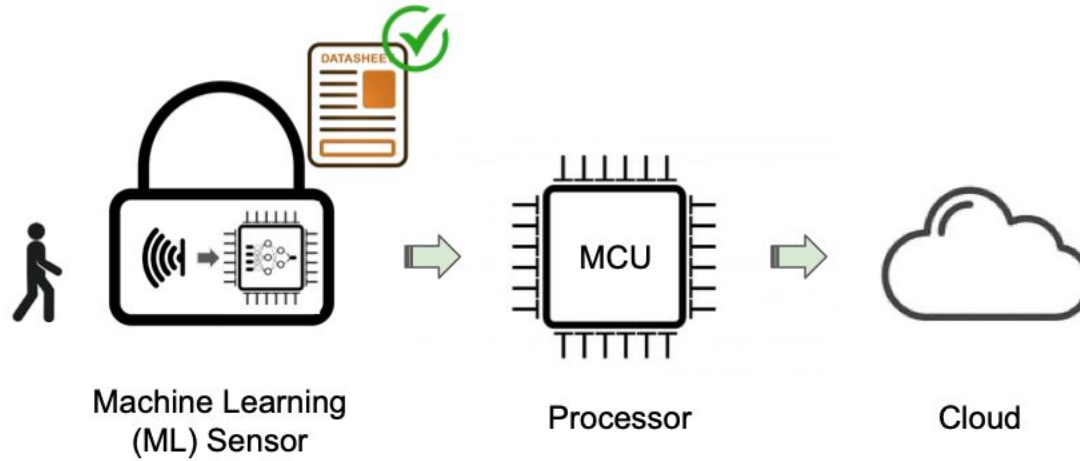Machine Learning (ML) Sensor — Processor — Cloud

97

# Datasheets for ML Sensors

ML sensors must be **transparent**, indicating in a **publicly and freely accessible** ML sensor **datasheet** all the relevant information such as **fact sheets**, **model cards**, and **dataset nutrition labels** to supplement the traditional EE hardware information typically available for sensors.



## PA1 Person Detection Module

**Description:** The PA1 Person Detection Module enables you to quickly and easily add smarts to your IoT deployment to monitor and detect for humans. You can use this module indoors and outdoors to understand where and when humans arrive at your deployment site.
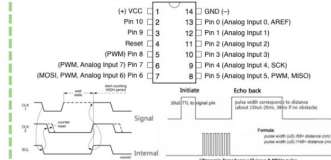
**Features:**
- Real-time Person Detection with On-Device ML
- Indoor and Outdoor use
- Finds a person at a maximum distance of 10 meters to a minimum distance of 5 centimeters
- Operates in low and high light environments (1-20000 Lux) across a wide temperature range (0 to 50 ℃)
- Features Color and Black-and-White Detection Modules

**Use Cases:**
- Smart business and home security systems
- Multi-modal key word spotting for virtual assistants
- Occupancy sensors and other infrastructure sensors
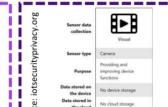
Description, Features, and Use Cases

Compliance

Diagrams and Form Factor

Source: docs.luxonis.com

Hardware Characteristics

Communication Specification and Pinout

Sources: fabacademy.org, electroschematics.com, and nxp.com/docs

Model Characteristics

**Model performance:** Measured with Precision-Recall (PR) and Area Under the PR Curve (PR-AUC). Download raw performance results data here. Disaggregated performance measured with Recall, which captures how often the model misses faces with specific characteristics. Equal recall across subgroups corresponds to the "Equality of Opportunity" fairness criterion.

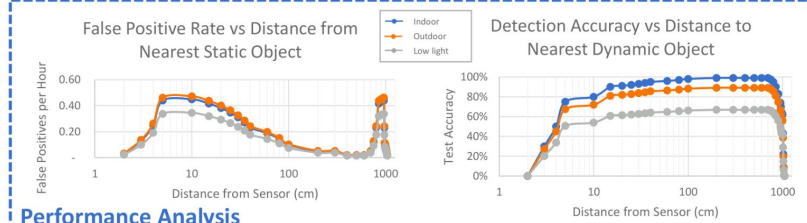**Performance evaluated on:**
- A subset of Open Images
- Face Detection Data Set and Benchmark
- Labeled Faces in the Wild

Source: modelcards.withgoogle.com

Dataset Nutrition Label

Source: datanutrition.org

IoT Security & Privacy Label

Source: iotsecurityprivacy.org

**Environmental Impact:** Full report can be found here. Environmental Impact

Source: st.com

Performance Analysis

Machine Learning Sensors - ML  ×     +

mlsensors.org

TinyML   Harvard   MLC   Research   Seeed   CS141   TimeBuddy   VJs Funding   Enterprise - Suppl...   Geo Chart Exampl...   »   Other Bookmarks

**MLSensors**
Machine Learning Sensors

Home     Whitepaper     GitHub     Email     Team

# Machine Learning Sensors

An ML sensor is a self-contained system that utilizes on-device machine learning to extract useful information by observing some complex set of phenomena in the physical world and reports it through a simple interface to a wider system.

Machine learning sensors represent a paradigm shift for the future of embedded machine learning applications. Current instantiations of embedded ML suffer from complex integration, lack of modularity, and privacy and security concerns from data

## Challenges 🔗

### Interface

What universal interface is needed for ML Sensors?

### Standards

What standards need to be in place for ML Sensors?

### Ethics

What ethical considerations are needed for ML Sensors?

## Call for Working Group Members

**We are actively growing our working group. If you would like to be a part of it please email us at: ml-sensors@googlegroups.com!**

## Example ML Sensor Datasheet

This illustrative example datasheet highlighting the various sections of an ML Sensor datasheet. On the top, we have the items currently found in standard datasheets: the description, features, use cases, diagrams and form factor, hardware characteristics, and communication specification and pinout. On the bottom, we have the new items that need to be included in an ML sensor datasheet: the ML model characteristics, dataset nutrition label, environmental impact analysis, and end-to-end performance analysis. While we compressed this datasheet into a one-page illustrative example by combining features and data from a mixture of sources, on a real datasheet, we assume each of these sections would be longer and include additional explanatory text to increase the transparency of the device to end-users. Interested users can find the most up-to-date version of the datasheet online at https://github.com/harvard-edge/ML-Sensors.

## PA1 Person Detection Module

**Description:** The PA1 Person Detection Module enables you to quickly and easily add smarts to your IoT deployment to monitor and detect for humans. You can use this module indoors and outdoors to understand where and when humans arrive at your deployment site.

**Features:**

Compliance

GDPR          RoHS
              2002/95/EC

# ML Sensors - Guiding Set of Principles

1. We need to **raise the level of abstraction** to enable ease of use for scalable deployment of ML sensors; not everyone should be required to be an systems developer or an engineer to use or leverage ML sensors into their ecosystem.

2. The ML sensor's **design should be inherently data-centric** and defined by its input-output behavior instead of exposing the underlying hardware and software mechanisms that support ML model execution.

3. An ML sensor's **implementation must be clean and complexity-free**. Features such as reusability, software updates, and networking must be thought through to ensure data privacy and secure execution.

4. ML sensors **must be transparent, indicating in a publicly and freely accessible ML sensor datasheet** all the relevant information such as fact sheets, model cards, and dataset nutrition labels to supplement the traditional information available for hardware sensors.

5. We as a community should aim to **foster an open ML sensors ecosystem by maximizing data, model, and hardware transparency** where possible, without necessarily relinquishing any claim to intellectual property.

# Machine Learning Sensors

"An ML sensor is self-contained system that utilizes on-device machine learning extract useful information by observing some complex set of phenomena in the physical world and reports it through a simple interface to a wider system."

**Stage 4: Live Operations**

- Model drift and skew
- Ethical challenges
- Sustainability
- Security and privacy

**Stage 1: Sensor Data**

- Heterogeneous devices
- Multi-modal data
- Sensor drift
- Varying data frequency

**TinyML Lifecycle**

**Stage 3: Model Deployment**

- Model robustness
- Scalable deployment
- Embedded MLOps
- System Integration

**Stage 2: Model Development**

- ML model architecture
- Resource constraints
- Model quality/accuracy
- End-to-end performance

# Course Topics

1. **Overview and Introduction to Embedded Machine Learning**
2. Data Engineering
3. Embedded Machine Learning Frameworks
4. Efficient Model Representation and Compression
5. Performance Metrics and Benchmarking of ML Systems
6. Learning on the Edge
7. Hardware Acceleration for Edge ML: GPUs, TPUs and FPGAs
8. Embedded MLOps
9. Secure and Privacy-Preserving On-Device ML
10. Responsible AI
11. Sustainability at the Edge
12. Generative AI at the Edge

# Conclusion

1. TinyML has the **potential to radically change our future**
2. No free lunch – hardware and software **fragmentation is a serious challenge** to address
3. TinyML **sustainability is crucial** to ensure its broad applicability
4. ML sensors based on TinyML technology must be **transparent**
5. Widening access to applied ML is a must to ensure **equitable access**



*The future of ML is tiny and bright, and its benefits can translate to societal impact.*

105

# CS249r Fall 2022 Projects

# Discussion Topics

# Understanding TinyML

What is tinyML and how does it differ from traditional machine learning approaches?

What are the potential applications of tinyML in everyday life?

How does tinyML align with the current trends in the Internet of Things (IoT)?



**Class Discussion**

# Application and Use Cases

What kind of real-world problems can be solved more efficiently with tinyML compared to traditional ML solutions?

How can tinyML contribute to energy conservation and sustainability?

What industries or sectors could benefit the most from tinyML technologies?

# Security and Privacy

What are the potential security vulnerabilities associated with deploying tinyML in sensitive applications?

How can privacy be maintained when deploying tinyML solutions in personal devices?

How can tinyML aid in the development of secure communication networks?

# Ethics and Society

What are the potential societal impacts of widespread adoption of tinyML technologies?

What ethical considerations should be kept in mind when deploying tinyML in public spaces?

How can tinyML technologies be made accessible and inclusive for different communities?