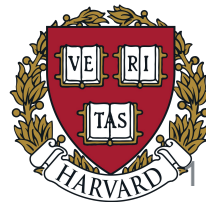# Tiny Machine Learning

*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*

*Chips & Compilers Symposium at MLSys '22, Sep. 1, 2022*

# TinyML

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

**TinyML**

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

**TinyML** → Fast-growing field of **ML**

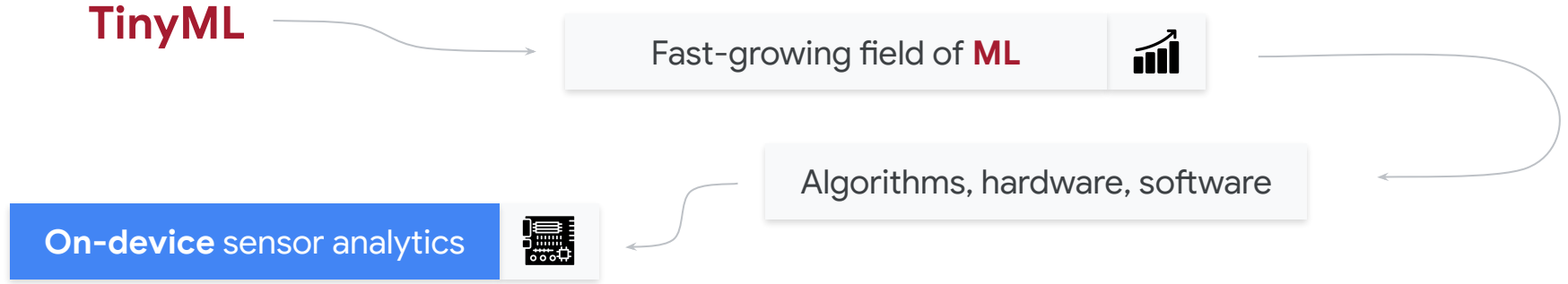# What is Tiny Machine Learning (**TinyML**)?

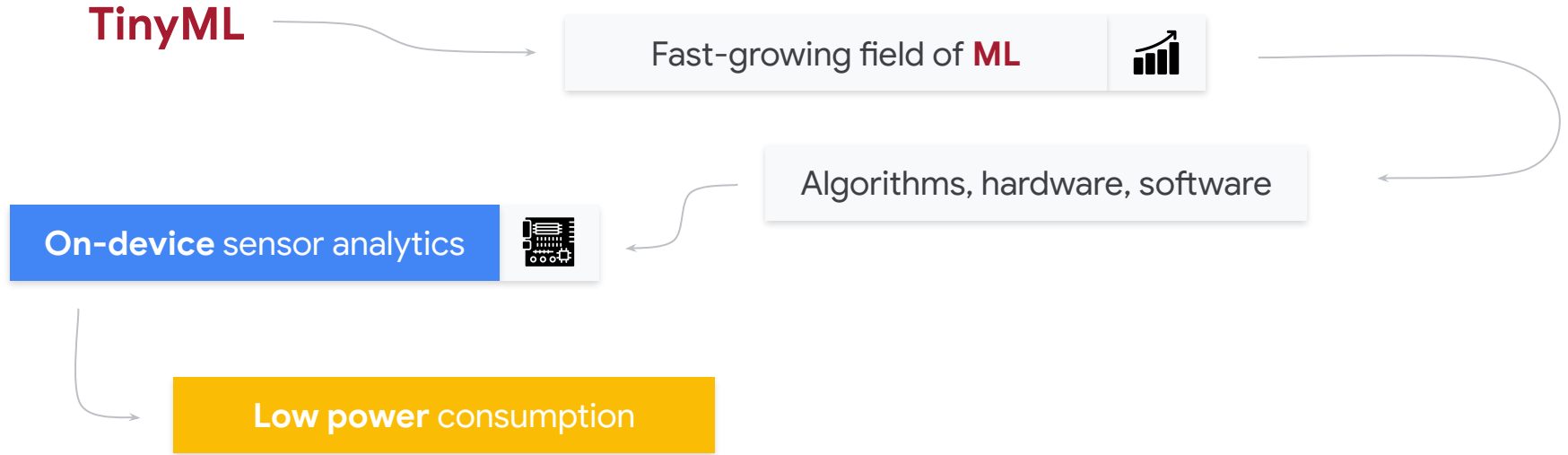**TinyML**

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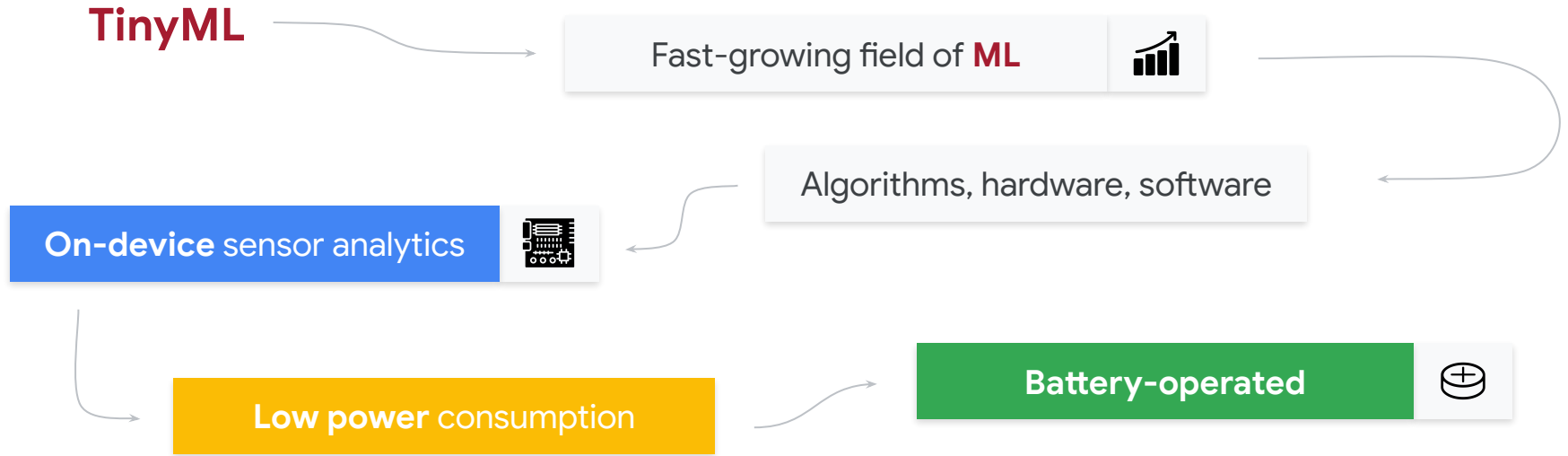
Algorithms, hardware, software

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

**TinyML**

Fast-growing field of **ML**

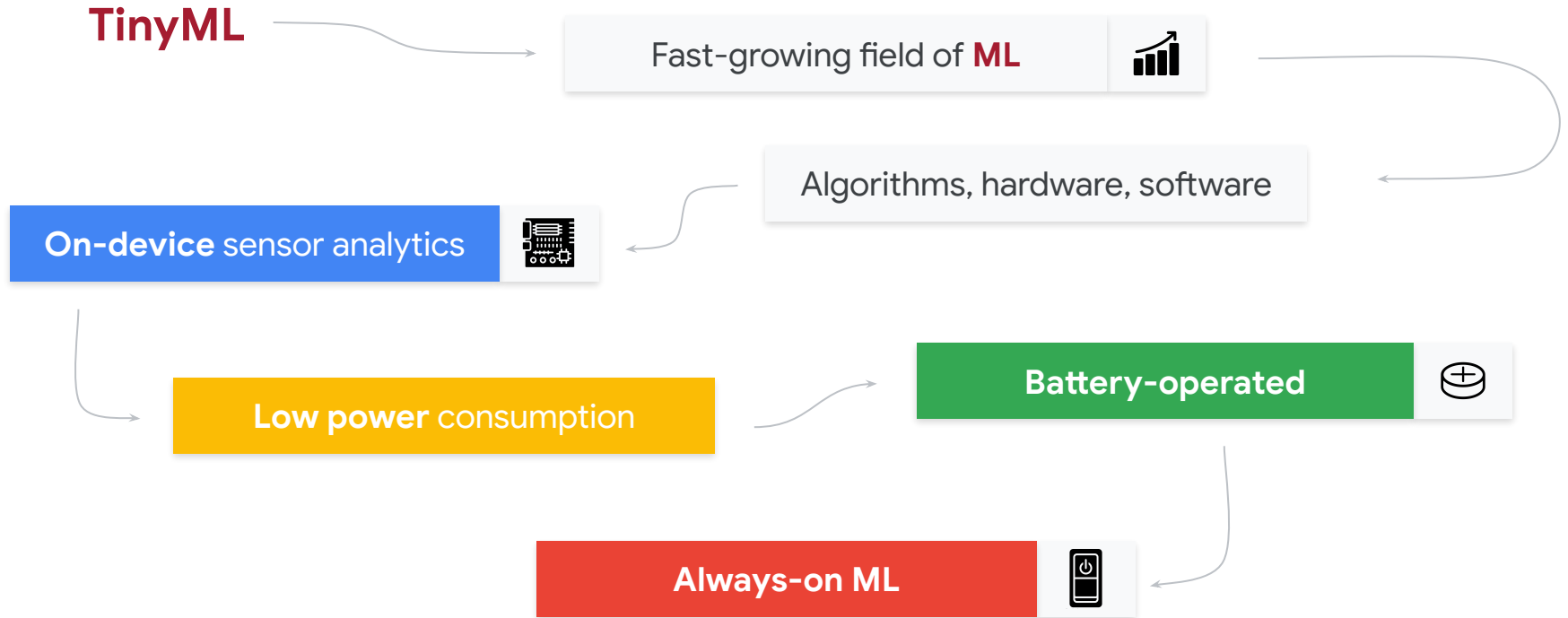Algorithms, hardware, software

**On-device** sensor analytics

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

**TinyML**

Fast-growing field of **ML**

Algorithms, hardware, software

**On-device** sensor analytics

**Low power** consumption

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

**TinyML**

Fast-growing field of **ML**

Algorithms, hardware, software

**On-device** sensor analytics

**Low power** consumption

**Battery-operated**
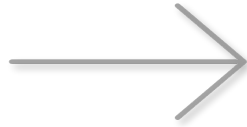
# What is Tiny Machine Learning (**TinyML**)?

**TinyML**

Fast-growing field of **ML**

Algorithms, hardware, software

**On-device** sensor analytics

**Low power** consumption

**Battery-operated**

**Always-on ML**

# Mobile

# IoT 1.0: Internet of Things
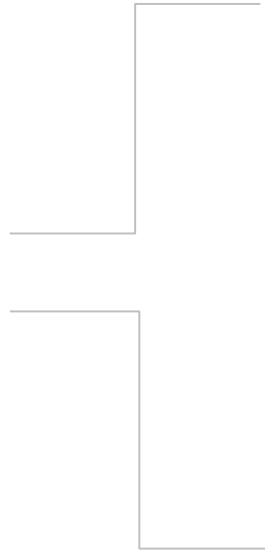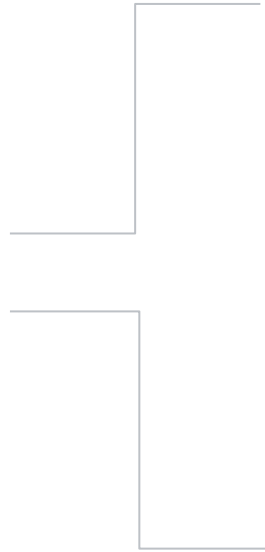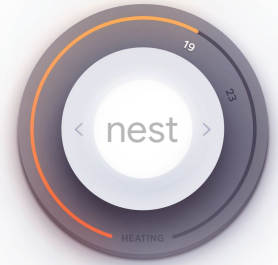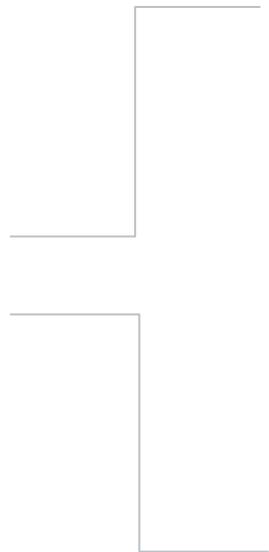
Google Assistant

# IoT 2.0: Intelligence on Things

# IoT 2.0: Intelligence on Things

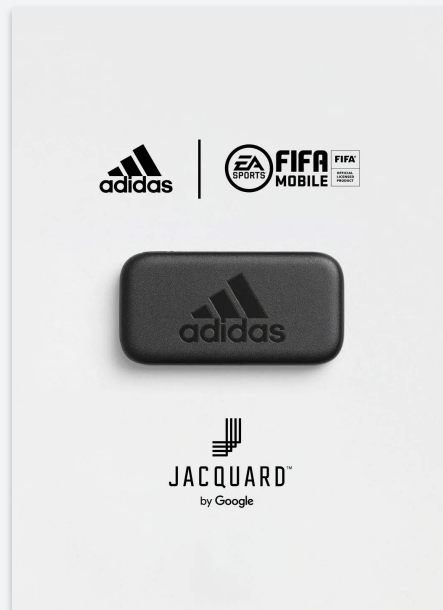**Bandwidth**
**Reliability**
**Latency**
**Privacy**
**Energy**

# Emerging TinyML Use Cases

**Example: Smart shoes**

- Kicking
- Penalty kicking
- Passing
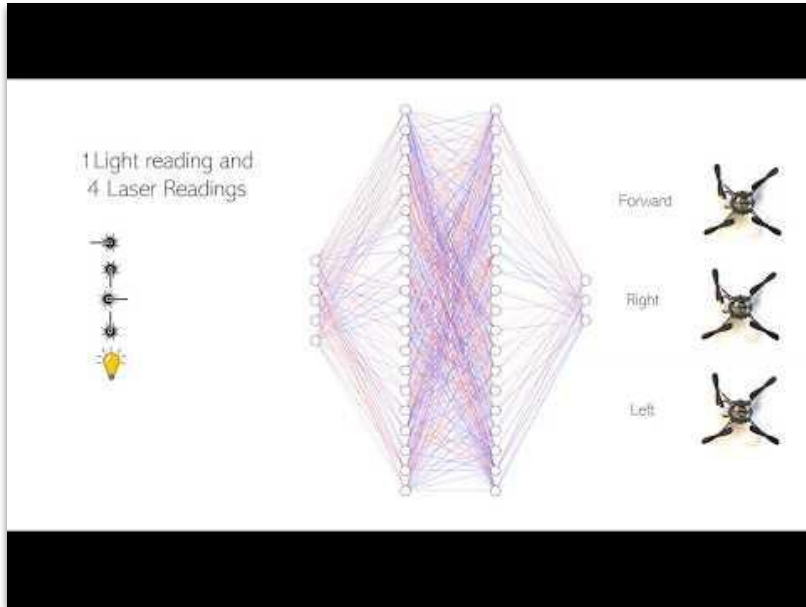- Dribbling
- ...

# Emerging TinyML Use Cases

**Example: Augmented Reality**

- Eye tracking

- Hand tracking

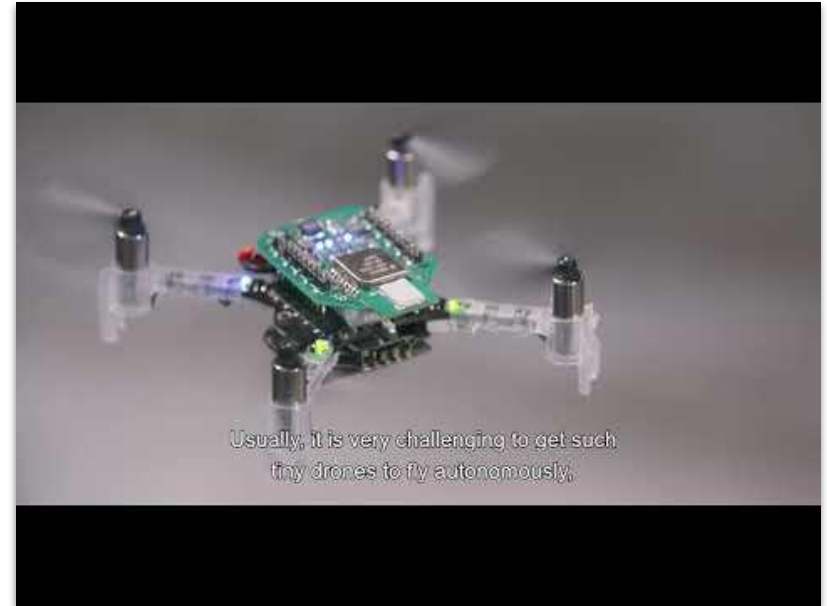- Computer vision

- Superresolution

- …

# Tiny Robot Learning



Duisterhof, B.P., Krishnan, S., Cruz, J.J., Banbury, C.R., Fu, W., Faust, A., de Croon, G.C. and Reddi, V.J., 2021, May. Tiny robot learning (tinyrl) for source seeking on a nano quadcopter. In *2021 IEEE International Conference on Robotics and Automation (ICRA)* (pp. 7242-7248). IEEE.
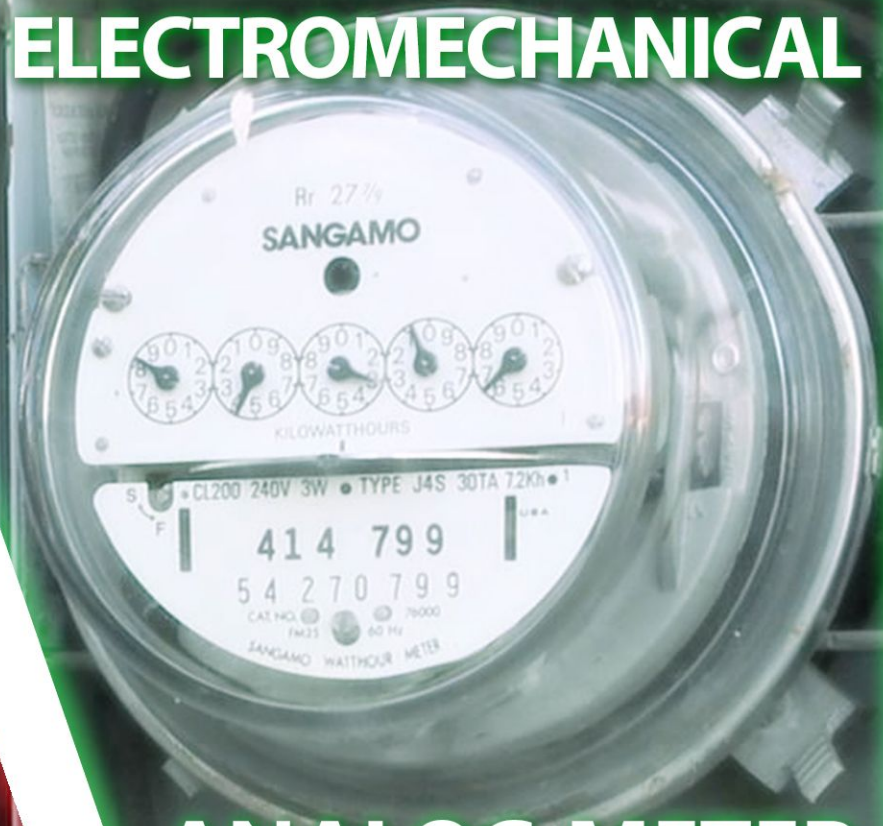


Duisterhof, B.P., Li, S., Burgués, J., Reddi, V.J. and de Croon, G.C., 2021, September. Sniffy bug: A fully autonomous swarm of gas-seeking nano quadcopters in cluttered environments. In 2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) (pp. 9099-9106). IEEE.

A DIGITAL

AN ELECTROMECHANICAL

"SMART" METER

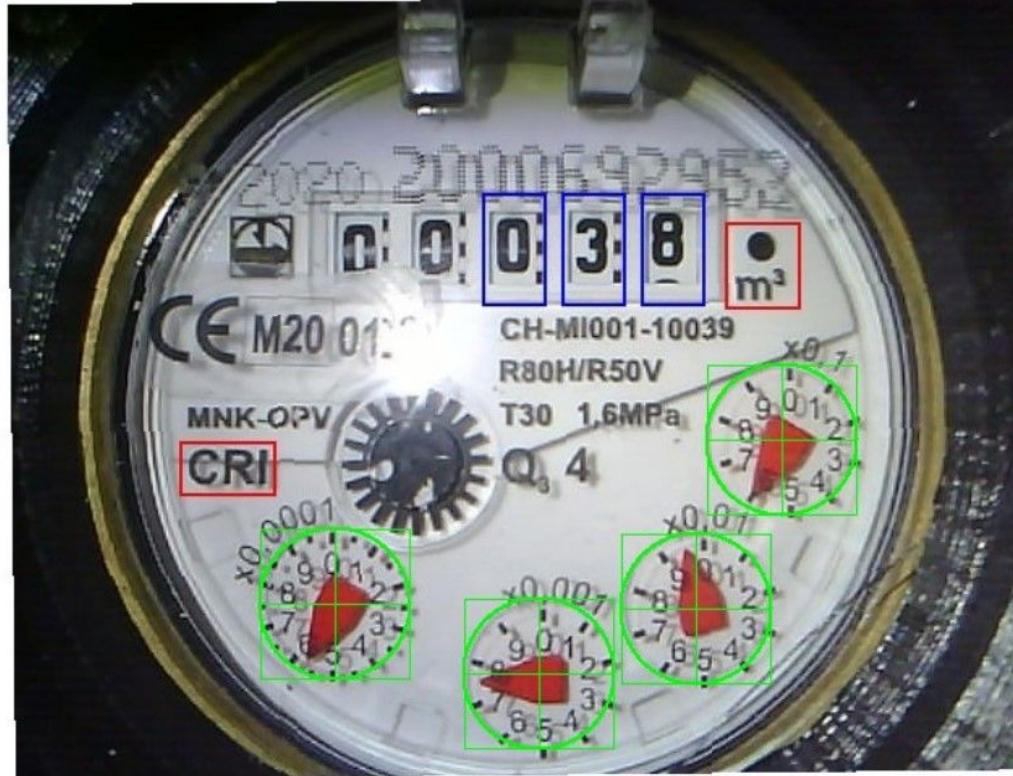ANALOG METER

# 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 | |

# 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

...

# 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

# Forbes

Subscribe   Sign In

Nov 8, 2021, 08:30am EST | 15,696 views

## Meet TinyML: The Latest Machine Learning Tech Having An Outsize Business Impact

**Dr. Nicholas Nicoloudis** Brand Contributor
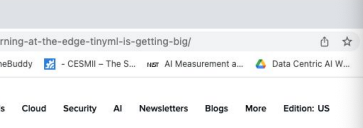SAP BRANDVOICE | Paid Program
Innovation

As device sensors proliferate across
product development through insp
surfacing to provide actionable insi
There are sound economic reasons
researchers predict IoT will have a
trillion by 2025, identifying manufa
trillion).

The rise of tinyML to collect data from edge de
explosion of sensors in pretty much every indu

The tinyML community was establi
learning architectures, techniques,
on-device analytics for a variety of
chemical, and others) at low power
devices. One of the tinyML founde

*"..we are in the midst of the digital
ultimate benefits of extreme energy
intelligence and analytics at low c
features...".*

---

# EETimes

---

# ZDNet

Windows 11   5G   Best VPNs   Cloud   Security   AI   Newsletters   Blogs   More   Edition: US

⚠ MUST READ:   **Log4j flaw:** Now state-backed hackers are using bug as part of attacks,

## Machine learning at the edge: Tin
getting big

Being able to deploy machine learning applications at the edge is the key to unlocking
TinyML is the art and science of producing machine learning models frugal enough to
rapid growth.

Written by **George Anadiotis**, Contributing Writer
Posted in Big on Data on June 7, 2021 | Topic: Big Data

Is it $61 billion and 38.4% CAGR by 2028 or $43 billion and 37.4% CAGR by 2027? Depends on
which report outlining the growth of edge computing you choose to go by, but in the end it's not
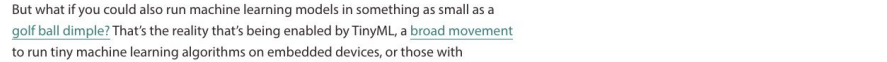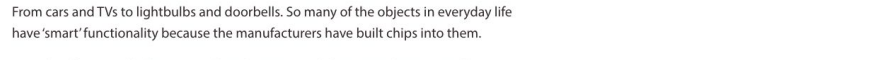that different.

What matters is that edge computing is booming. There is growing interest by vendors, and ample
coverage, for good reason. Although the definition of what constitutes edge computing is a bit
fuzzy, the idea is simple. It's about taking compute out of the data center, and bringing it as close
to where the action is as possible.

Whether it's stand-alone IoT sensors, devices of all kinds, drones,
or autonomous vehicles, there's one thing in common. Increasingly,
data generated at the edge are used to feed applications powered
by machine learning models. There's just one problem: machine
learning models were never designed to be deployed at the edge.
Not until now, at least. Enter TinyML.

EXECUTIVE GUIDE

**What is machine learning?**
Everything you need to

Tiny machine learning (TinyML) is broadly defined as a fast growing

---

# SEMICONDUCTOR ENGINEERING

---

# CIO

DIGITAL ISSUE   AWARDS   EVENTS   CIO THINK TANK   NEWSLETTERS   RESOURCES   INSIDER   SIGN IN   REGISTER

**Hot Topics**   IT Leadership   Digital Transformation   Innovation   Data Analytics & AI   Enterprise Applications   Diversity and Inclusion

Home

SAP

### NEXT EVOLUTION OF MACHINE LEARNING IS UPON US

SPONSORED

## How TinyML is powering big ideas across critical industries

**BrandPost** Sponsored by SAP | Learn More | JUL 18, 2021 4:31 PM PDT

shutterstock

From cars and TVs to lightbulbs and doorbells. So many of the objects in everyday life
have 'smart' functionality because the manufacturers have built chips into them.

But what if you could also run machine learning models in something as small as a
golf ball dimple? That's the reality that's being enabled by TinyML, a broad movement
to run tiny machine learning algorithms on embedded devices, or those with

# Questions



**01** How do we design an open-source ecosystem to enable TinyML to thrive in the face of heterogeneity?

**02** How do we benchmark the various TinyML solutions to enable "apples to apples" system comparisons?

**03** How do we drive hardware and software co-design in a flexible manner across the complete system stack?

# Questions



How do we design an
**open-source ecosystem to**
**enable TinyML to thrive in**
**the face of heterogeneity?**

**03**

How do we drive hardware
and software co-design in a
flexible manner across the
complete system stack?

**01**

**02**

How do we benchmark the
various TinyML solutions to
enable "apples to apples"
system comparisons?

24

# 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
- Heterogeneity
  - CPU
  - GPU
  - DSP
  - NPU
- Resource Constraints
  - Memory
  - Power

Software
- Missing Library Features
  - malloc
  - ...
- Limited Operating System Support

Challenges

Hardware

Software

Heterogeneity

Resource Constraints

Missing Library Features

Limited Operating System Support

CPU

GPU

DSP

NPU

Memory

Power

malloc

...

Challenges

Hardware
- Heterogeneity
  - CPU
  - GPU
  - DSP
  - NPU
- Resource Constraints
  - Memory
  - Power

Software
- Missing Library Features
  - malloc
  - ...
- Limited Operating System Support

...

TensorFlow Lite Micro

Arduino
BLE Sense 33

Himax
WE-I Plus EVB

SparkFun
Edge 2

Espressif
EYE

...

# **TFLite Micro** Design

- TFLite Micro uses an **interpreter** design

- Store the model as data and loop through its ops at **runtime**



dispatch
**loop**

instruction
**ops**

dispatch
**loop**

instruction
**ops**

```
int main() {
    function_a();
    function_b();

    printf("done!\n");
}

void function_a() {
    doSomething();
    saveTheWorld();
    machineLearning++;

    printf("a is complete\n");
}

void functon_b() {
    x = 50;
    y = 249;
    z = 141;

    int result = run_conv(x,y,z);

    result += 61;

    printf("b is complete\n");
}
```

C/C++
**code**

one time
**compilation**

```
010101010100101010010101010101010
101101010111010101010101010101010
101010111010101001000111001011101
010010101110101110101010101110101
010101010101001010100101010101010
101011011010111010101010101011010
101010101110101010100011100101111
010100010101110101110101010101101
010101010101010010101001010101010
101010110110101110101010101010110
101010101010111010101000011100101
110101000101011110101110101010111
010101010101010100101010101010101
101010101101101011101010101011011
010101010101011101010100011100101
011101010010101110101110101010111
110101010101010101001010100101010
101010101011011010111010101010101
```

compiled
**machine
code**

**Interpreter**
(generally **slower** than compiled code)

**Compiler**
(generally **faster** than interpreted code)

# ML is **Different**

- Each layer like a `Conv` or `softmax` can take tens of thousands or even millions of cycles to complete execution

# ML is **Different**

- Parsing overhead is **relatively small** for the TFMicro interpreter when we consider the **overall network graph**

| Model | Total Cycles | Calculation Cycles | Interpreter Overhead |
|---|---|---|---|
| Visual Wake Words (Ref) | 18,990.8K | 18,987.1K | **< 0.1%** |
| Google Hotword (Ref) | 36.4K | 34.9K | **4.1%** |



Sparkfun Edge 2
(Apollo 3 **Cortex-M4**)

# Interpreter
## **Advantages**

- Change the model
  **without recompiling**
  the code

instruction
**ops**

dispatch
**loop**

dispatch
**loop**

instruction
**ops**

# Interpreter
# **Advantages**

- Change the model
  **without recompiling**
  the code
- **Same operator code**
  can be used across
  multiple **different**
  **models** in the system

| Arduino BLE Sense 33 | Himax WE-I Plus EVB |
|---|---|
| Espressif EYE | SparkFun Edge 2 |

# Interpreter **Advantages**

- Same **portable** model serialization format can be used **across a lots of systems**.

**TensorFlow** Lite Micro

**Core functionality**

< 20KBs

Model loading

Error reporting

Memory planner

...

**Model operators**

conv2D
conv3D
tanh
depthwise_conv2d
l2_normalize
sigmoid
max_pool
...

# Memory
# **Improvements**

- Selective op registration **reduces memory consumption by 30%**

- **Memory reduction varies by model**, depending on the operators used by the model

# 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



David, R., Duke, J., Jain, A., Janapa Reddi, V., Jeffries, N., Li, J., Kreeger, N., Nappier, I., Natraj, M., Wang, T. and Warden, P., 2021. Tensorflow lite micro: Embedded machine learning for tinyml systems. Proceedings of Machine Learning and Systems, 3, pp.800-811.

# Questions



How do we drive hardware and software co-design in a flexible manner across the complete system stack?
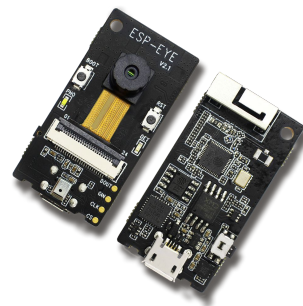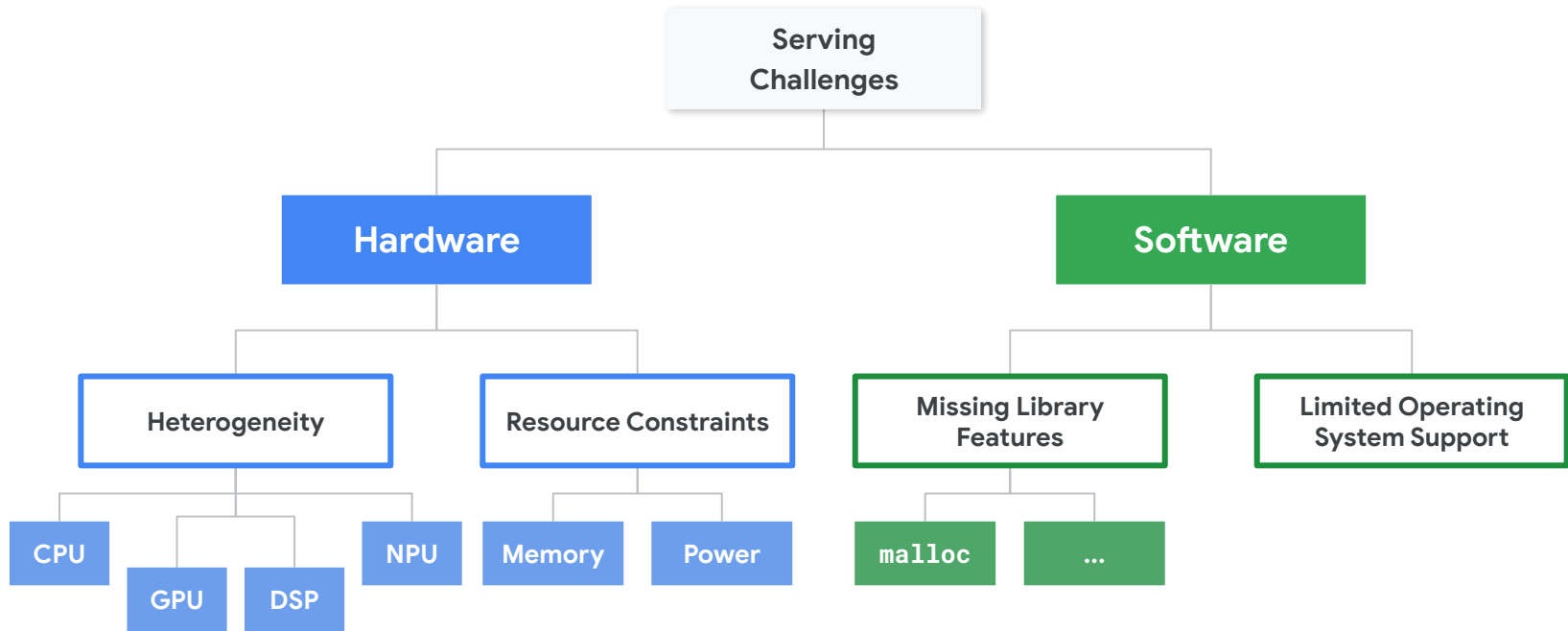
How do we design an open-source ecosystem to enable TinyML to thrive in the face of heterogeneity?

**03**

**01**

**02**

**How do we benchmark the various TinyML solutions to enable "apples to apples" system comparisons?**

| 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 |

Serving Challenges

- Hardware
  - Heterogeneity
    - CPU
    - GPU
    - DSP
    - NPU
  - Resource Constraints
    - Memory
    - Power
- Software
  - Missing Library Features
    - malloc
    - ...
  - Limited Operating System Support

# TinyML System Stack is Complicated

- Machine learning system stack is **complicated**

- Many **different** models, datasets, models, frameworks, formats, compilers, libraries, operating systems, targets

- The **cross-product** makes it challenging to decipher system performance

| | | | | |
|---|---|---|---|---|
| Sensors | Camera | Microphone | IMU | |
| ML Applications | Person Detection | Keyword Spotting | Anomaly Detection | |
| ML Datasets | Visual Wake Words | Google Speech Commands | ToyADMOS | |
| ML Models | MobileNet | MicroNets | RNN | AutoEncoder |
| Training Framework | | TensorFlow | PyTorch | |
| Graph Formats | | TFLite | ONNX | |
| Inference Framework | TensorFlow Lite for Microcontrollers | uTVM | STM Cube.AI | TinyEngine |
| Optimized Libraries | CMSIS-NN | | embARC | CEVA |
| Operating Systems | MBED OS | RTOS | Zephyr | VxWorks |
| Hardware Targets | MCU | DSP | uNPU | Accelerators |

# Apples-to-apples comparison



**ML System X**

**ML System Y**

What task?
What model?
What dataset?
What batch size?
What quantization?
What software libraries?
...

## bench·mark

/ˈben(t)SHmärk/

See definitions in:

All    Technology    Surveying

*noun*

1. a standard or point of reference against which things may be compared or assessed.
   "a benchmark case"

   Similar:   standard   point of reference   basis   gauge   criterion   specification

2. a surveyor's mark cut in a wall, pillar, or building and used as a reference point in measuring altitudes.

*verb*

evaluate or check (something) by comparison with a standard.
"we are **benchmarking** our performance **against** external criteria"

Definitions from Oxford Languages      Feedback

# Benchmarking

**Use to**

- **Compare** solutions
- **Inform** selection
- **Measure** and track progress
- **Raise** the bar, **advance** the field

**Requires**

- **Methodology** that is both fair and rigorous
- **Community** support and consensus



**Provides**

- **Standardization** of use cases and workloads
- **Comparability** across heterogeneous HW/SW systems
- **Complex characterization** of system compromises
- **Verifiable and Reproducible** results

# Wide Array of ML Tasks

| Task Category | Use Case | Model Type | Datasets |
|---|---|---|---|
| Audio | Audio Wake Words<br>Context Recognition<br>Control Words<br>Keyword Detection | DNN<br>CNN<br>RNN<br>LSTM | Speech Commands<br>Audioset<br>ExtraSensory<br>Freesound<br>DCASE |
| Image | Visual Wake Words<br>Object Detection<br>Gesture Recognition<br>Object Counting<br>Text Recognition | DNN<br>CNN<br>SVM<br>Decision Tree<br>KNN<br>Linear | Visual Wake Words<br>CIFAR10<br>MNIST<br>ImageNet<br>DVS128 Gesture |
| Physiological /<br>Behavioral Metrics | Segmentation<br>Anomaly Detection<br>Forecasting<br>Activity Detection | DNN<br>Decision Tree<br>SVM<br>Linear | Physionet<br>HAR<br>DSA<br>Opportunity |
| Industry Telemetry | Sensing<br>Predictive Maintenance<br>Motor Control | DNN<br>Decision Tree<br>SVM<br>Linear<br>Naive Bayes | UCI Air Quality<br>UCI Gas<br>UCI EMG<br>NASA's PCoE |

MLPerf

# Goals

**MLPerf**

**Enforce performance result replicability** to ensure reliable results

Use **representative workloads**, reflecting production use-cases

**Encourage innovation** to improve the state-of-the-art of ML

Accelerate progress in ML via **fair and useful measurement**

Serve both the **commercial and research communities**

Keep **benchmarking affordable** so that all can participate

# MLPerf "Tiny" Tasks



### Keyword Spotting

Audio Sample Data

< 30ms >
< 20ms >

FFT

< 256 values >

Average

< 43 values >

Spectrogram

Warden, Pete. "Speech commands: A dataset for limited-vocabulary speech recognition." *arXiv preprint arXiv:1804.03209* (2018).

### Visual Wake Words

(a) 'Person'

(b) 'Not-person'

Chowdhery, Aakanksha, et al. "Visual wake words dataset." *arXiv preprint arXiv:1906.05721* (2019).

### Anomaly Detection

Pump

Microphone array

Fan          90 deg.          Valve
180 deg.          0 deg.
500 mm          100 mm
270 deg.

500 mm

Slide rail

Purohit, Harsh, et al. "MIMII dataset: Sound dataset for malfunctioning industrial machine investigation and inspection." *arXiv preprint arXiv:1909.09347* (2019).

### Tiny Image Classification

airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

Krizhevsky, Alex, and Geoffrey Hinton. "Learning multiple layers of features from tiny images." (2009): 7.

| Problem definition | Dataset selection (public domain) | Model selection | Model training code | Derive "Tiny" version: Quantization | Embedded implementation | Benchmarking harness integration | Deploy on device | Example benchmark run |
|---|---|---|---|---|---|---|---|---|

**Anomalous Sound Detection System**

Normal

**Anomaly**

FP32 → INT8 → **ARM** mbed OS

Training Code

| Problem | AD |
|---|---|
| Model | FC-AE |
| Size | 270 Kpar |
| Latency | 10.4 ms/inf. |
| Accuracy | .86 AUC |
| Energy | 516 µJ/inf. |

# Metrics

## Latency

Small fast dataset

Loop of inferences

No data-dependent execution

Host

USB Hub

API

DUT

a.

```
Runtime requirements have been met.
Performance results for window 10:
  # Inferences :      1000
  Runtime      :    10.524 sec.
  Throughput   :    95.020 inf./sec.
Runtime requirements have been met.
-------------------------------------
Median throughput is 95.019 inf./sec.
```

## Accuracy

Evaluate on larger dataset

Top-1 accuracy & AUC

**CLOSED**: meet threshold
v.
**OPEN**: part of the metrics

## Energy

No "cherry-picking"

Power Monitor setup

Median result

Host

LEVEL SHIFTER

USB Hub

IO Manager

Energy Monitor

TRIG   VCC   GND

GPIO/TRIG

UART

API

DUT

ENERGY CHART
/Users/ptorelli/eembc/runner/sessions/20210426101328/trace1-energy.bin

Energy (uJ)

Time (s)

Sample Rate: 1.000 kHz          197.5 sec., 63.8 uJ          [15.1 J, 70.8 mW]

# MLPerf Tiny
# in a Nutshell

Built to benchmark **embedded ML systems**:

- **Standardize best practices** in TinyML benchmarking
- **Measure both ML performance and power** consumption
- Designed to be **portable across a wide variety of systems**

| Division | Dataset | Training | Model | Numerics | Framework | Hardware | Demonstrates |
|---|---|---|---|---|---|---|---|
| Closed | X | X | X | INT-8 PTQ | TensorFlow Lite Micro | ARM MCU | Baseline performance results on the reference platform. |
| Closed | X | X | X | INT-8 PTQ | TensorFlow Lite Micro | RISC-V MCU | Performance of a RISC-V microcontroller customized for neural network inference. |
| Closed | X | X | X | FP-32 & INT-8 PTQ | LEIP Framework | RasPi 4 | Capabilities of a software-only optimization toolchain that is agnostic of the hardware. |
| Closed | X | X | X | INT-8 PTQ | Syntiant TDK | Neural Network Accelerator | Ultra-low power hardware efficiency for running deep neural networks. |
| Open | X | QKeras | ✓ | Int-6/8 QAT | HLS4ML | FPGA | Rapid end-to-end development of machine learning accelerators on reconfigurable fabrics. |



Banbury, C., Reddi, V.J., Torelli, P., Holleman, J., Jeffries, N., Kiraly, C., Montino, P., Kanter, D., Ahmed, S., Pau, D. and Thakker, U., 2021. Mlperf tiny benchmark. NeurIPS'21

# Toward Emerging Multi-DNN Models

**Pipelined DNNs**



**Keyword Spotting** **Speech Processing**

- **Back-to-back execution**
- **Execution dependency**

**Concurrent DNNs**



**Eye Tracking** **Obstacle Detection** **Video Processing**

- **Concurrent execution**
- **Execution deadline**

**Concurrent & Pipelined DNNs**

**Obstacle Detection**

**Eye Tracking** → **Foveated Rendering**

- **Challenges from both pipelined and concurrent**

64

# MetaBench
# in a Nutshell (Stay Tuned!)

- We **demystify the unique features and challenges of MMMT workloads** for Metaverse applications
- We **provide a taxonomy of MMMT workloads** to understand new classes of deep learning inference workloads and discuss their feature and challenges
- Based on realistic applications, we propose a **real-time MMMT benchmark suite** that models the different Metaverse end-user usage scenarios.
- We also **discuss the need for new scoring metrics** that reflect ML system performance in a useful manner.



65

# Questions

**03**

**How do we drive hardware and software co-design in a flexible manner across the complete system stack?**

How do we design an open-source ecosystem to enable TinyML to thrive in the face of heterogeneity?

**01**

**02**

How do we benchmark the various TinyML solutions to enable "apples to apples" system comparisons?

# The Hardware Lottery



- Sara Hooker's observation that the success of new ML approaches depends on their compatibility with downstream software and hardware. Here you can "make your own luck"!

**MCUs: KBs** of RAM, **Fixed/slow** processor



**Specialized Hardware Customization** (on FPGAs)

# CFU Playground

- **ML library**
  - TensorFlow Lite     -- *open source*
- **CPU ISA**
  - RISC-V              -- *open source*
- **CPU design**
  - VexRiscv            -- *open source*
- **FPGA SoC/IP**
  - LiteX               -- *open source*
- **FPGA synth/PnR**
  - SymbiFlow,Yosys, -- *open source*
    Nextpnr, VPR       -- *open source*

*FPGA vendor tools can be used if you wish*

- **Python HW gen**
  - Migen, nMigen      -- *open source*
- **Simulation**
  - Renode, Verilator  -- *open source*

The only proprietary component is the FPGA itself



**Full-Stack Open-Source Framework**
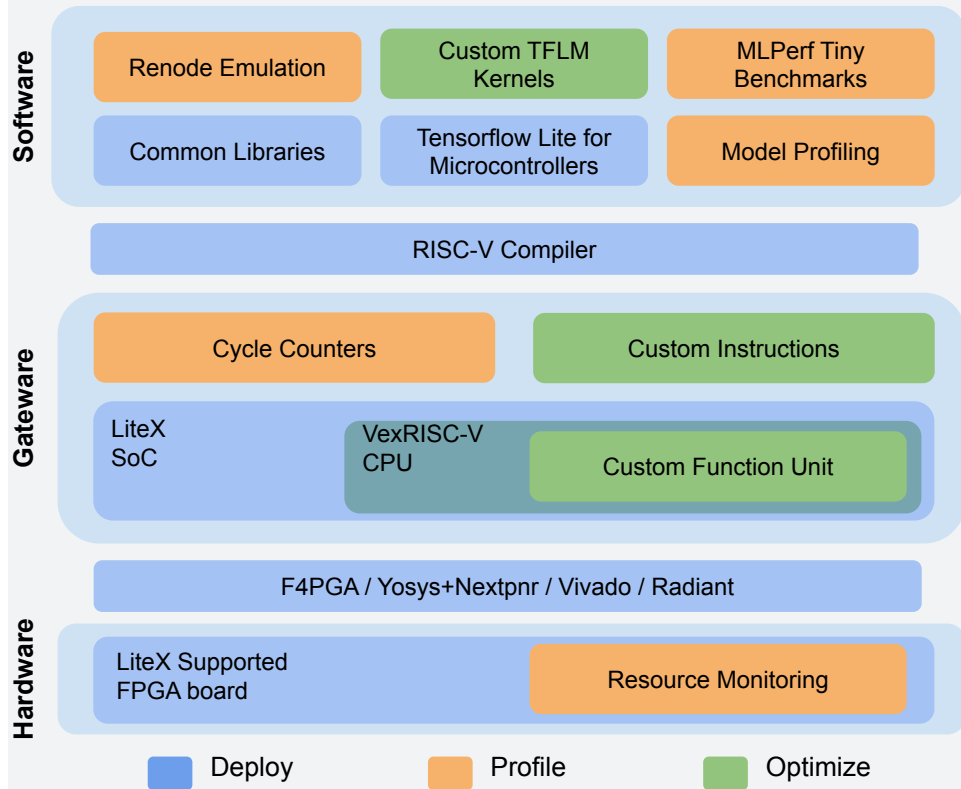
69

# CFU Playground

- **ML library**
  - TensorFlow Lite    *-- open source*
- **CPU ISA**
  - RISC-V             *-- open source*
- **CPU design**
  - VexRiscv           *-- open source*
- **FPGA SoC/IP**
  - LiteX              *-- open source*
- **FPGA synth/PnR**
  - SymbiFlow,Yosys,   *-- open source*
    Nextpnr, VPR       *-- open source*

*FPGA vendor tools can be used if you wish*

- **Python HW gen**
  - Migen, nMigen      *-- open source*
- **Simulation**
  - Renode, Verilator  *-- open source*
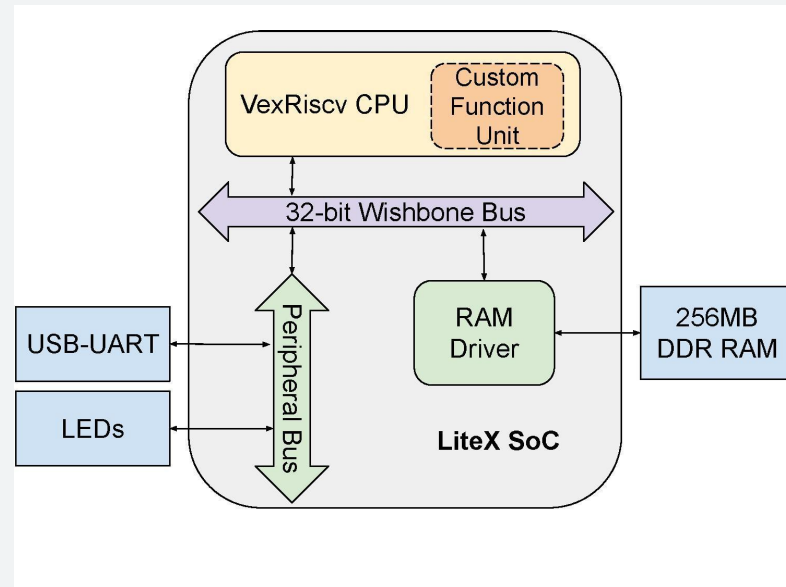
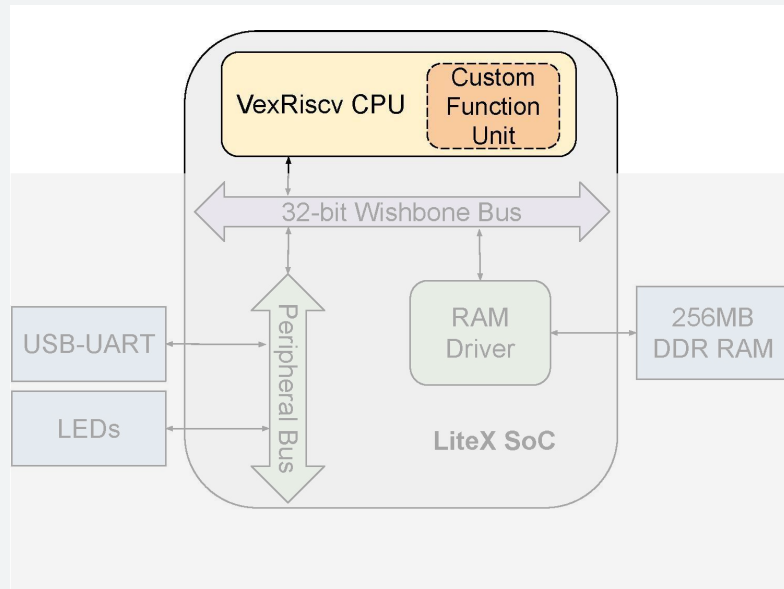The only proprietary component is the FPGA itself

# CFU Playground

- **ML library**
  - TensorFlow Lite    -- *open source*
- **CPU ISA**
  - RISC-V              -- *open source*
- **CPU design**
  - VexRiscv            -- *open source*
- **FPGA SoC/IP**
  - LiteX               -- *open source*
- **FPGA synth/PnR**
  - SymbiFlow,Yosys, -- *open source*
    Nextpnr, VPR      -- *open source*

*FPGA vendor tools can be used if you wish*

- **Python HW gen**
  - Migen, nMigen      -- *open source*
- **Simulation**
  - Renode, Verilator -- *open source*
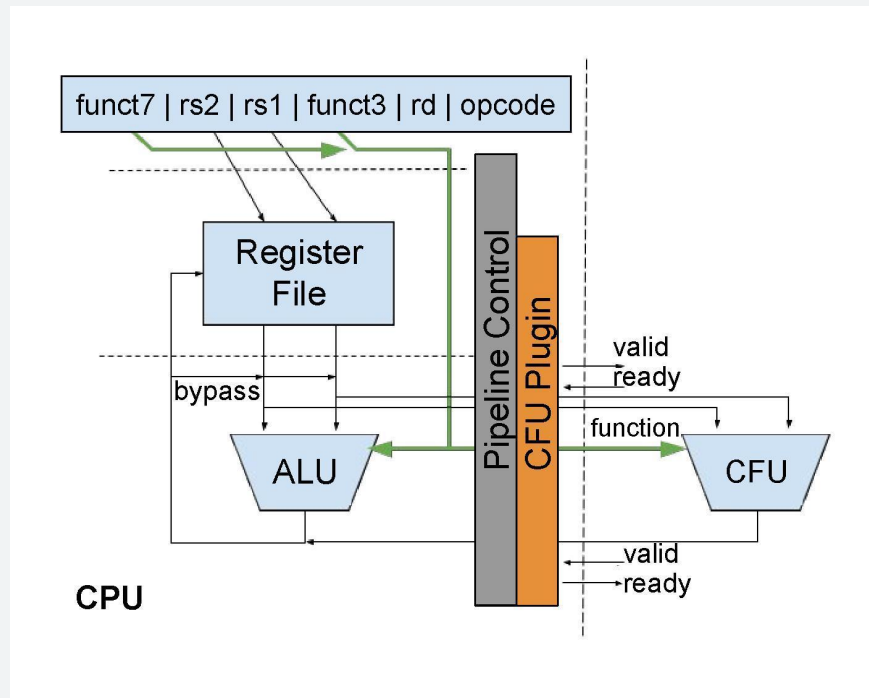
The only proprietary component is the FPGA itself

# CFU Playground

- **ML library**
  - TensorFlow Lite   -- *open source*
- **CPU ISA**
  - RISC-V           -- *open source*
- **CPU design**
  - VexRiscv         -- *open source*
- **FPGA SoC/IP**
  - LiteX            -- *open source*
- **FPGA synth/PnR**
  - SymbiFlow,Yosys, -- *open source*
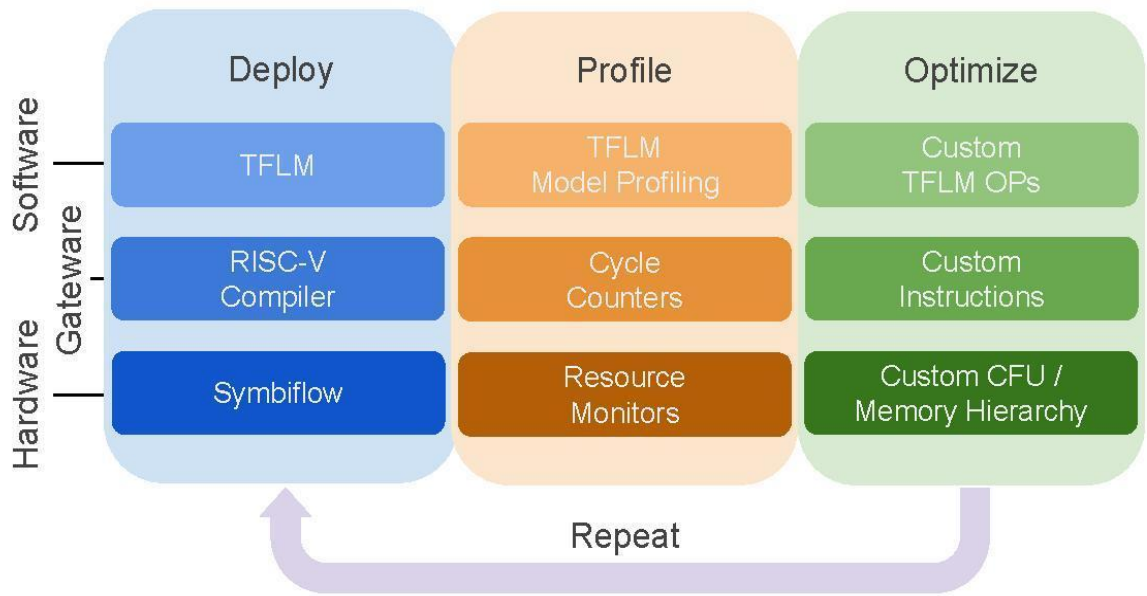    Nextpnr, VPR     -- *open source*

*FPGA vendor tools can be used if you wish*

- **Python HW gen**
  - Migen, nMigen    -- *open source*
- **Simulation**
  - Renode, Verilator -- *open source*

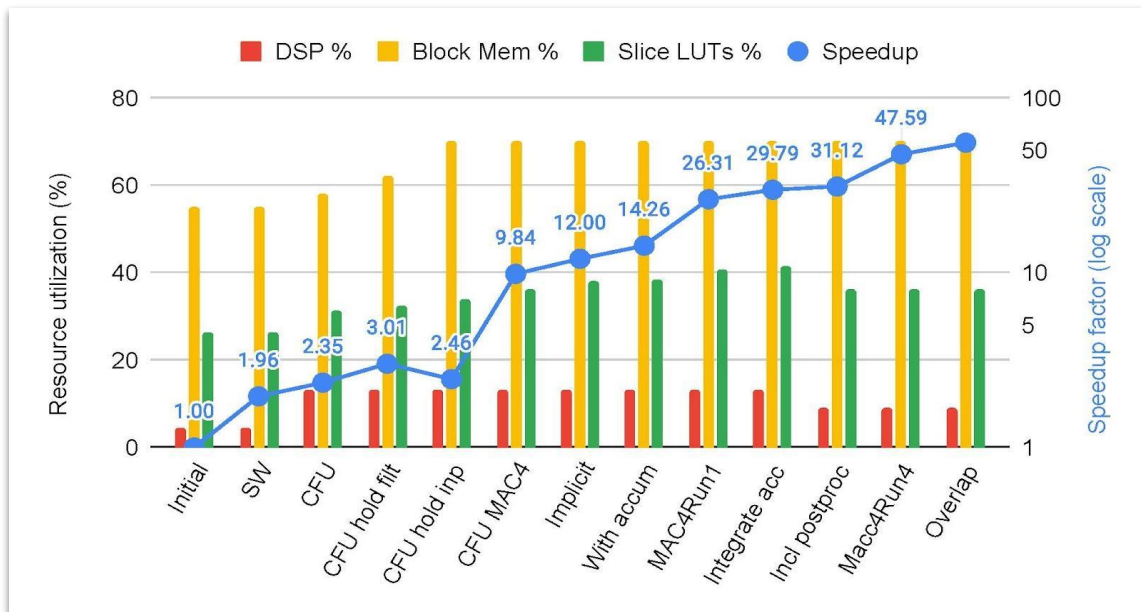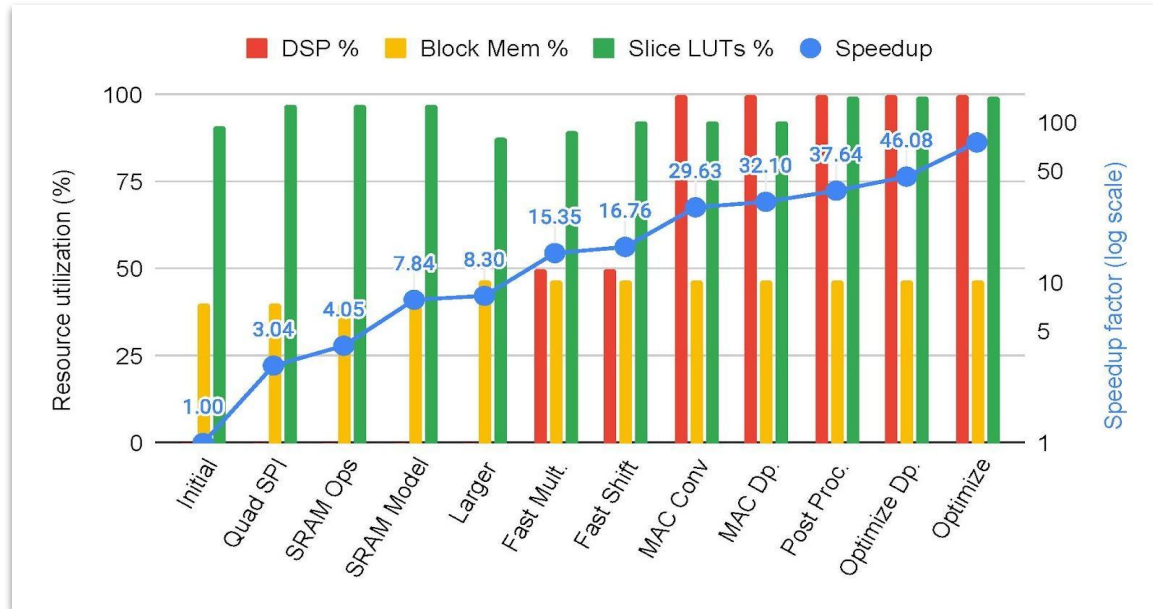The only proprietary component is the FPGA itself

**Agile Design Methodology**

**Image Classification on Arty**

Legend: DSP %, Block Mem %, Slice LUTs %, Speedup

Speedup values: 1.00, 1.96, 2.35, 3.01, 2.46, 9.84, 12.00, 14.26, 26.31, 29.79, 31.12, 47.59

X-axis categories: Initial, SW, CFU, CFU hold filt, CFU hold inp, CFU MAC4, Implicit, With accum, MAC4Run1, Integrate acc, Incl postproc, Macc4Run4, Overlap

Y-axis (left): Resource utilization (%)
Y-axis (right): Speedup factor (log scale)

55x speedup in 5 weeks (part-time)

**Keyword Spotting on FOMU**

75x speedup in under 4 weeks (intern)

| Software | | | |
|---|---|---|---|
| Renode Emulation | Custom TFLM Kernels | MLPerf Tiny Benchmarks | |
| Common Libraries | Tensorflow Lite for Microcontrollers | Model Profiling | |

RISC-V Compiler

| Gateware | |
|---|---|
| Cycle Counters | Custom Instructions |

LiteX SoC

VexRISC-V CPU — Custom Function Unit

F4PGA / Yosys+Nextpnr / Vivado / Radiant

Hardware

LiteX Supported FPGA board — Resource Monitoring

Deploy    Profile    Optimize

# CFU Playground in a Nutshell

- An **out-of-the-box, full-stack framework that fully integrates open-source tools** across the entire stack to facilitate rich community-driven ecosystem development

- An **agile methodology for developers to progressively and iteratively design bespoke accelerators** for resource-constrained, latency-bound tinyML applications

- Through cross-stack insights, we **demonstrate novel model-specific resource allocation trade-offs between the CFU, CPU, and memory system** that enable optimal ML performance on resource-constrained FPGA platforms

**CFU Playground: Full-Stack Open-Source Framework for Tiny Machine Learning (tinyML) Acceleration on FPGAs**

Shvetank Prakash* Tim Callahan† Joseph Bushagour§ Colby Banbury*
Alan V. Green† Pete Warden† Tim Ansell† Vijay Janapa Reddi*

†*Google* §*Purdue University* *Harvard University*

arXiv:2201.01863v1 [cs.LG] 5 Jan 2022

77

# A Greener Tomorrow with TinyML

# 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% |

# Global CO$_2$ Emissions by Sectors



**Greenhouse Gas Emissions by Sector**
- Agriculture, Forestry, and Other Land Use
- Electricity and Heat Production
- Industry
- Transportation
- Buildings
- Other Energy

24.0%
25.0%
21.0%
14.0%
6.0%
10.0%

# TinyML System - Net Environmental Impact

# ML Sensors

Physical Sensor — Processor — Cloud

**Sensor 1.0**

Physical Sensor → Processor → Cloud

**Sensor 1.0**

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

**Sensor 2.0**

Sensor 1.0

Physical Sensor — Processor — Cloud

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

Sensor 2.0

# 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 °C)
- 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**

**Performance Analysis**

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

Machine Learning Sensors - M ×   +

mlsensors.org

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

## 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.
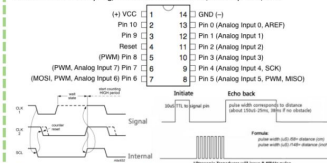
# PA1 Person Detection Module

**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:**

# Machine Learning Sensors

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

Pete Warden [1] Matthew Stewart [2] Brian Plancher [2] Colby Banbury [2] Shvetank Prakash [2] Emma Chen [2] Zain Asgar [1] Sachin Katti [1] Vijay Janapa Reddi [2]

[1]Stanford University [2]Harvard University

### ABSTRACT

Machine learning sensors represent a paradigm shift for the future of embedded machine learning applications. Current instantiations of embedded machine learning (ML) suffer from complex integration, lack of modularity, and privacy and security concerns from data movement. This article proposes a more data-centric paradigm for embedding sensor intelligence on edge devices to combat these challenges. Our vision for "sensor 2.0" entails segregating sensor input data and ML processing from the wider system at the hardware level and providing a thin interface that mimics traditional sensors in functionality. This separation leads to a modular and easy-to-use ML sensor device. We discuss challenges presented by the standard approach of building ML processing into the software stack of the controlling microprocessor on an embedded system and how the modularity of ML sensors alleviates these problems. ML sensors increase privacy and accuracy while making it easier for system builders to integrate ML into their products as a simple component. We provide examples of prospective ML sensors and an illustrative datasheet as a demonstration and hope that this will build a dialogue to progress us towards sensor 2.0.

## 1 INTRODUCTION

Since the advent of AlexNet [43], deep neural networks have proven to be robust solutions to many challenges that involve making sense of data from the physical world. Machine learning (ML) models can now run on low-cost, low-power hardware capable of deployment as part of an embedded device. Processing data close to the sensor on an embedded device allows for an expansive new variety of always-on ML use-cases that preserve bandwidth, latency, and energy while improving responsiveness and maintaining data privacy. This emerging field, commonly referred to as embedded ML or tiny machine learning (TinyML) [73, 18, 39, 59], is paving the way for a prosperous new array of use-cases, from personalized health initiatives to improving manufacturing productivity and everything in-between.

However, the current practice for combining inference and sensing is cumbersome and raises the barrier of entry to embedded ML. At present, the general design practice is to design or leverage a board with decoupled sensors and compute (in the form of a microcontroller or DSP), and for the developer to figure out how to run ML on these embedded platforms. The developer is expected to train and optimize ML models and fit them within the resource constraints of the embedded device. Once an acceptable prototype implementation is developed, the model is integrated with the rest of the software on the device. Finally, the widget is tethered to the device under test to run inference. The current approach is slow, manual, energy-inefficient, and error-prone.
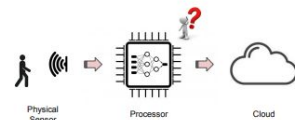
Figure 1. The Sensor 1.0 paradigm tightly couples the ML model with the application processor and logic, making it difficult to provide hard guarantees about the ML sensor's ultimate behavior.
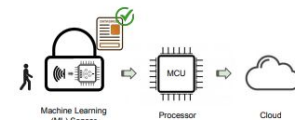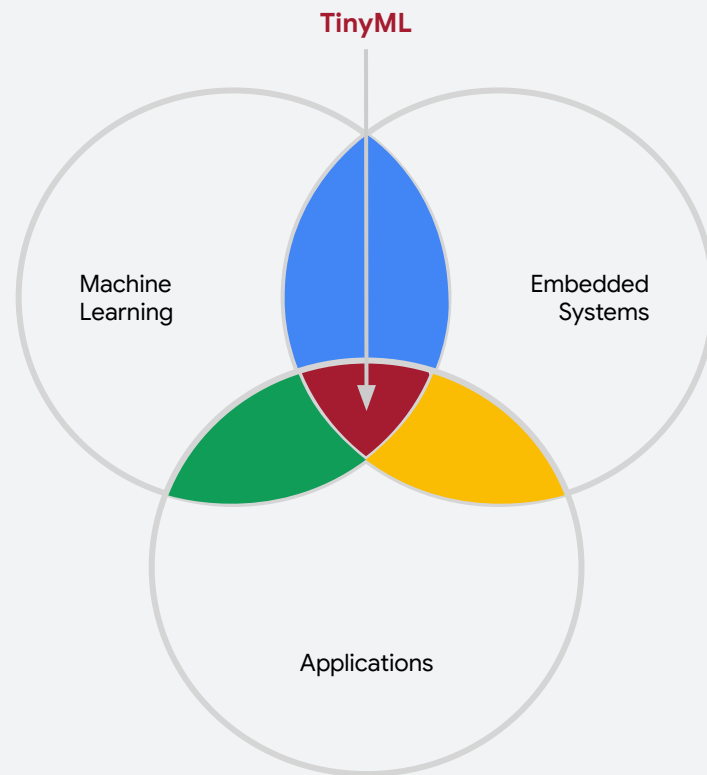
Figure 2. Our proposed Sensor 2.0 paradigm. The ML model is tightly coupled with the physical sensor, separate from the application processor, and comes with an ML sensor datasheet that makes its behavior transparent to the system integrators and developers.

It requires a sophisticated understanding of ML and the intricacies of ML model implementations to optimize and fit a model within the constraints of the embedded device.

# Conclusion

1. TinyML has the **potential to dramatically 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.*

# Conclusion

The Future of ML is Tiny and Bright