tinyML EMEA Technical Forum 2021 Proceedings

June 7 – 10, 2021

Virtual Event

www.tinyML.org

# IMPORTANCE OF BENCHMARKING

GIVEN A TASK AT HAND HOW DO WE KNOW WHICH IS THE **RIGHT SYSTEM**?
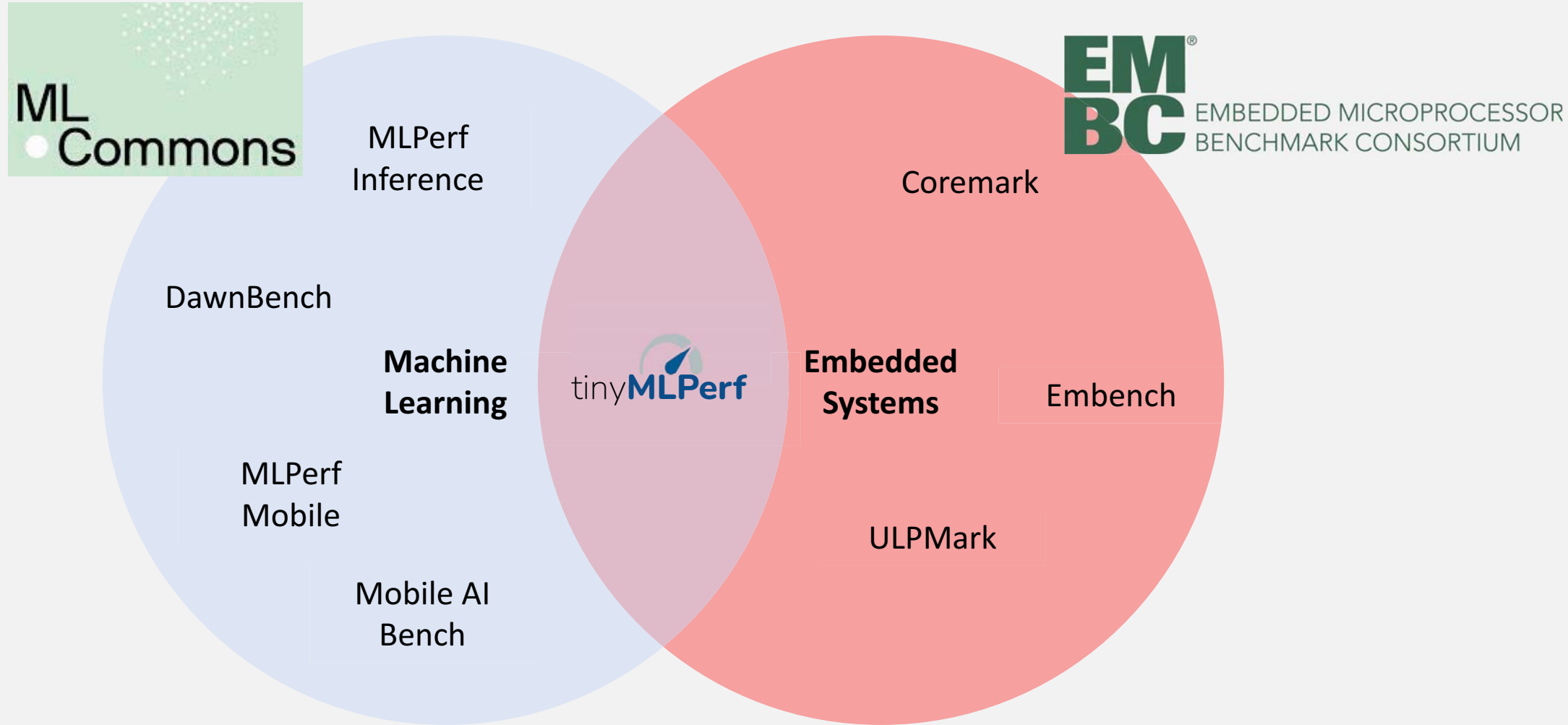
HOW DO WE **COMPARE** THE DIFFERENT SOLUTIONS?

# BENCHMARKING IS ALL ABOUT …

- **Comparability** across hardware and software

- **Standardization** of use cases and workloads

- **Measuring** progress via rigorous methodology

- **Community** building and consensus generation

# NEED FOR A TINYML SPECIFIC BENCHMARK

# BENCHMARK DESIGN CHALLENGES

- Number and diversity of use cases for tinyML

  - Which use cases should we focus on?

- TinyML innovation is in HW, in SW, in tooling, in algorithms, etc.

  - How do we support fair comparison and innovation?

- Embedded system design is always a compromise

  - What metrics to choose to fairly evaluate systems?

# BENCHMARK DESIGN CHALLENGES

- **Number and diversity of use cases for tinyML**
  - **Which use cases should we focus on?**

- TinyML innovation is in HW, in SW, in tooling, in algorithms, etc.
  - How do we support fair comparison and innovation?

- Embedded system design is always a compromise
  - What metrics to choose to fairly evaluate systems?

# RELEVANT USE CASES

| Task Category | Use Case | Model Type |
|---|---|---|
| Audio | Audio Wake Words<br>Context Recognition<br>Control Words<br>Keyword Detection | DNN<br>CNN<br>RNN<br>LSTM |
| 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 |
| Industry / Telemetry | Segmentation<br>Anomaly Detection<br>Forecasting<br>Activity Detection | DNN<br>Decision Tree<br>SVM<br>Linear |

# RELEVANT USE CASES & COMMUNITY

| Task Category | Use Case | Model Type |
|---|---|---|
| Audio | **Audio Wake Words** | DNN |
| | Context Recognition | CNN |
| | Control Words | RNN |
| | Keyword Detection | LSTM |
| Image | **Visual Wake Words** | DNN |
| | **Object Detection** | CNN |
| | Gesture Recognition | SVM |
| | Object Counting | Decision Tree |
| | Text Recognition | KNN |
| | | Linear |
| Industry / Telemetry | Segmentation | DNN |
| | **Anomaly Detection** | Decision Tree |
| | Forecasting | SVM |
| | Activity Detection | Linear |

WG lead: Vijay Janapa Reddi + Colby Banbury @ Harvard

Benchmark harness: Peter Torelly @ EEMBC + Nat Jeffries @ Google

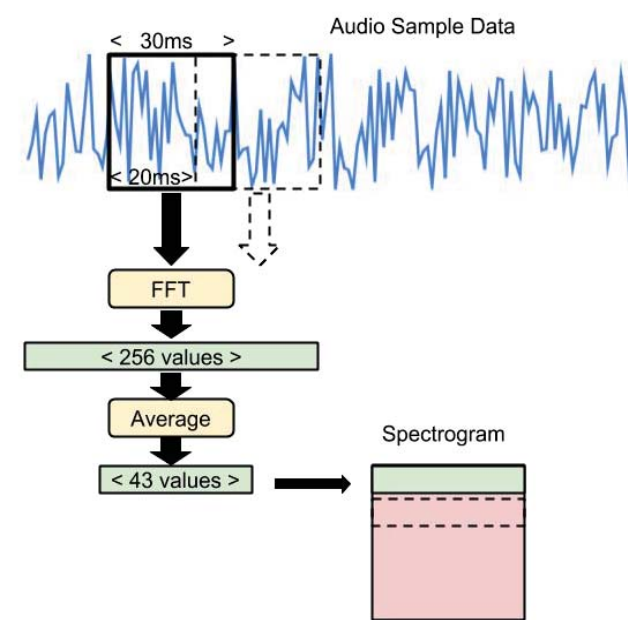Jeremy Holleman @ Syntiant

Nat Jeffries @ Google

Pietro Montino

Csaba Kiraly @ DC

# SELECTED BENCHMARKS

## Keyword Spotting

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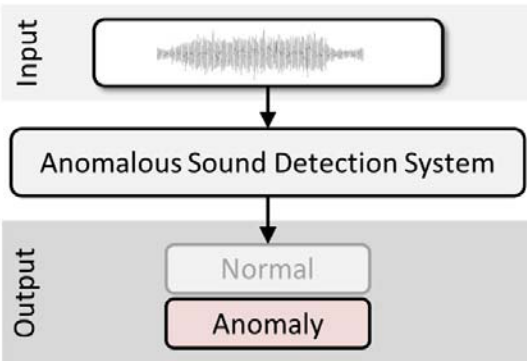
## Visual Wake Words
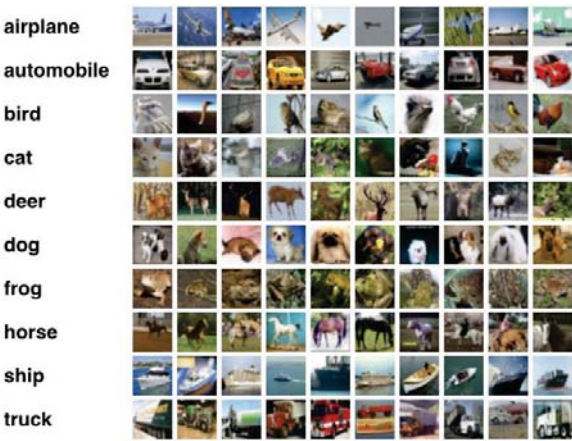


(a) 'Person'

(b) 'Not-person'

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

## Anomaly Detection

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

## Tiny Image Classification

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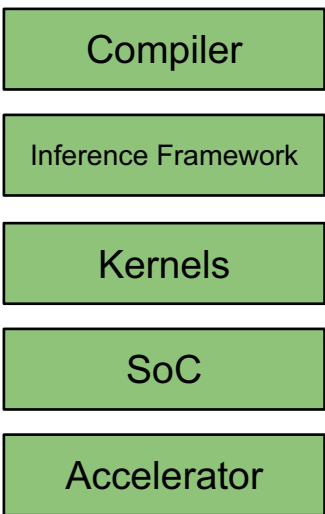
# BENCHMARK DESIGN CHALLENGES

- Number and diversity of use cases for tinyML
  - Which use cases should we focus on?

- **TinyML innovation is in HW, in SW, in tooling, in algorithms, etc.**
  - **How do we support fair comparison and innovation?**

- Embedded system design is always a compromise
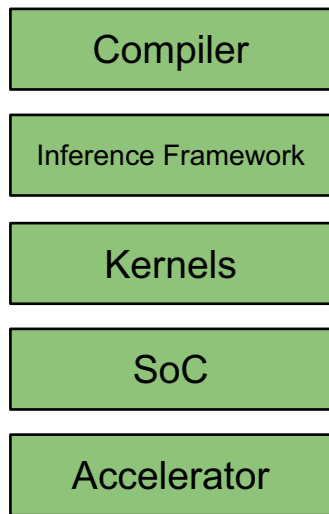  - What metrics to choose to fairly evaluate systems?
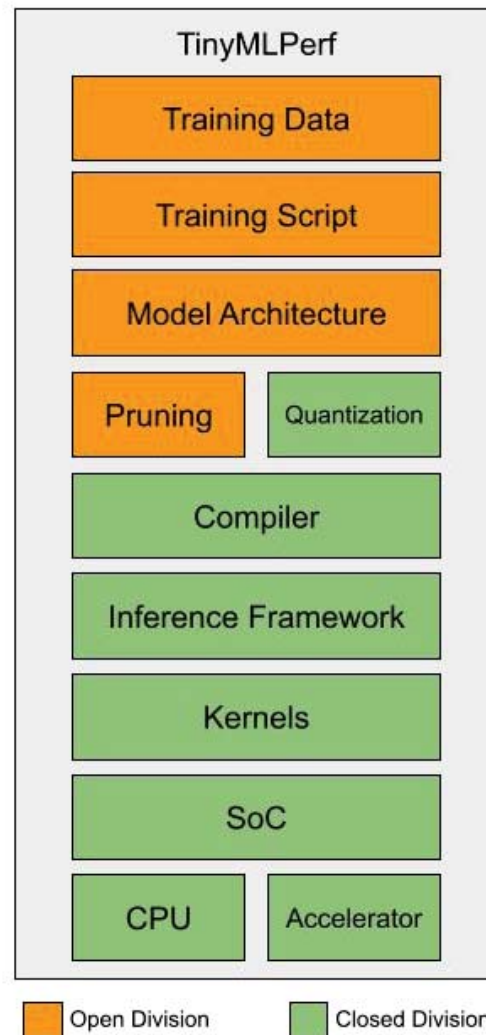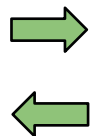
# Modular Design
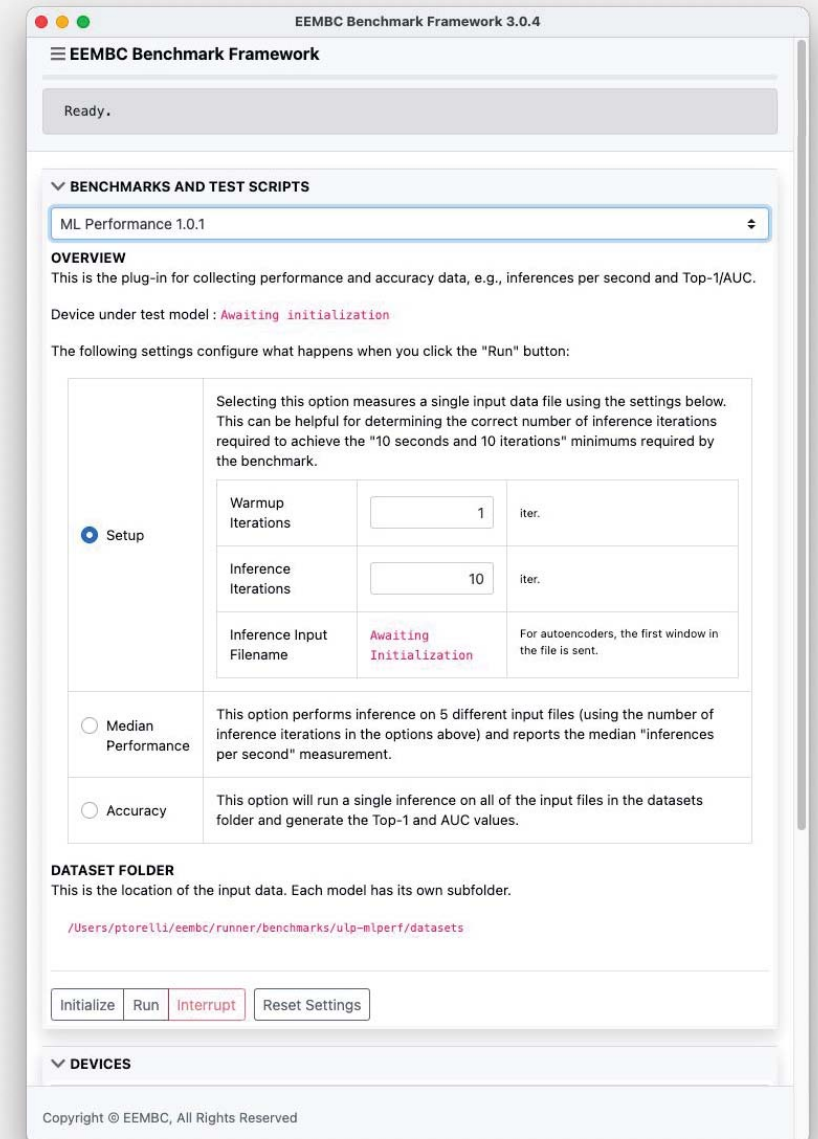
**Direct Comparison**
**CLOSED DIVISION**

**Demonstrate Improvement**
**OPEN DIVISION**

Closed
Submission A

Closed
Submission B

Open
Submission

TinyMLPerf

| Training Data |
| Training Script |
| Model Architecture |
| Pruning | Quantization |

| Compiler | | Compiler |
| Inference Framework | | Inference Framework |
| Kernels | VS. | Kernels |
| SoC | | SoC |
| Accelerator | | Accelerator |

| Compiler |
| Inference Framework |
| Kernels |
| SoC |
| CPU | Accelerator |

| Training Script |
| Model Architecture |

2X
Faster vs.
Reference

🟧 Open Division   🟩 Closed Division
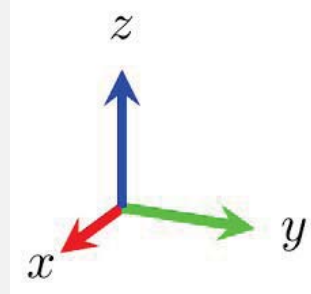
# BENCHMARK DESIGN CHALLENGES

- Number and diversity of use cases for tinyML
  - Which use cases should we focus on?

- TinyML innovation is in HW, in SW, in tooling, in algorithms, etc.
  - How do we support fair comparison and innovation?

- **Embedded system design is always a compromise**
  - **What metrics to choose to fairly evaluate systems?**

# METRICS OF INTEREST

- What metrics to choose, how to define
    - latency: ms/inference, or ms/task?
    - accuracy: which metric?
    - energy: what should be included?

- (Some) measurement difficulties
    - pre/post processing
    - single inference vs. more complex processing
    - accurate power measurement
    - accurate timestamping

- EEMBC Benchmark Runner

# MEASURING THOSE METRICS

## Latency

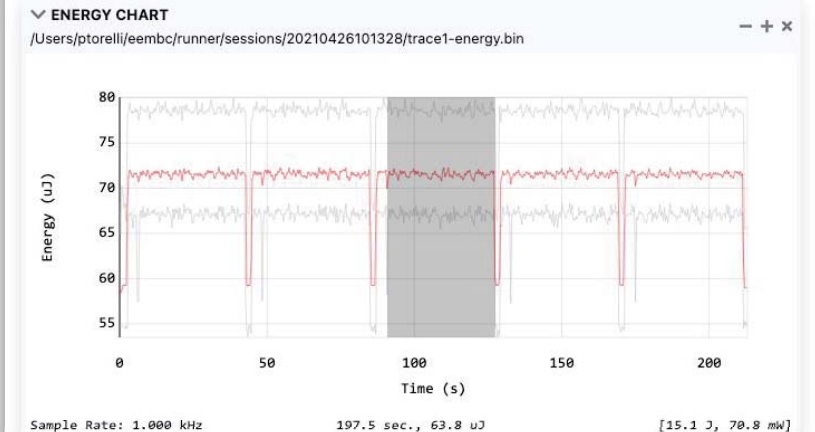Accurate timestamping
Evaluated on series of inferences
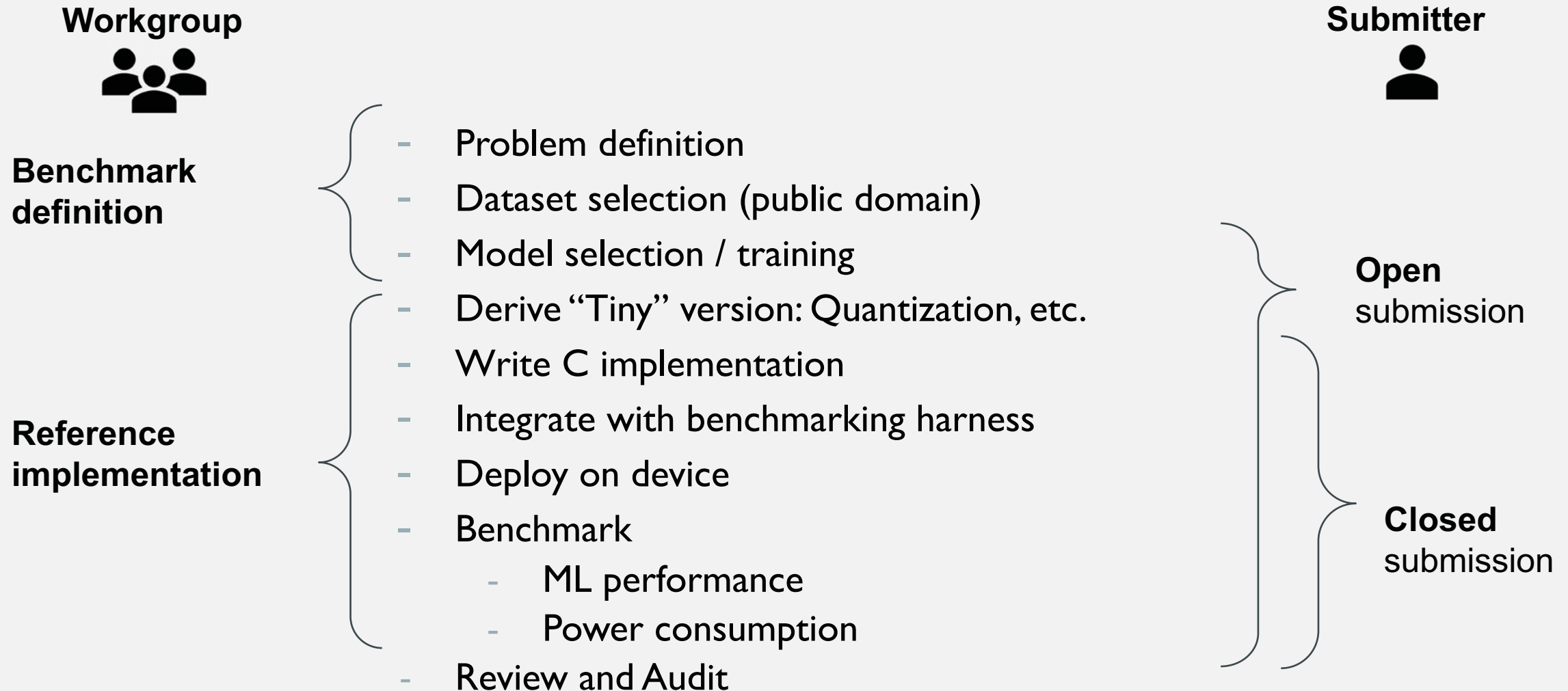Data-dependant execution?

## Accuracy

Top-1 accuracy & AUC
CLOSED: meet threshold
OPEN: part of the metrics

## Energy

Power Monitor integration
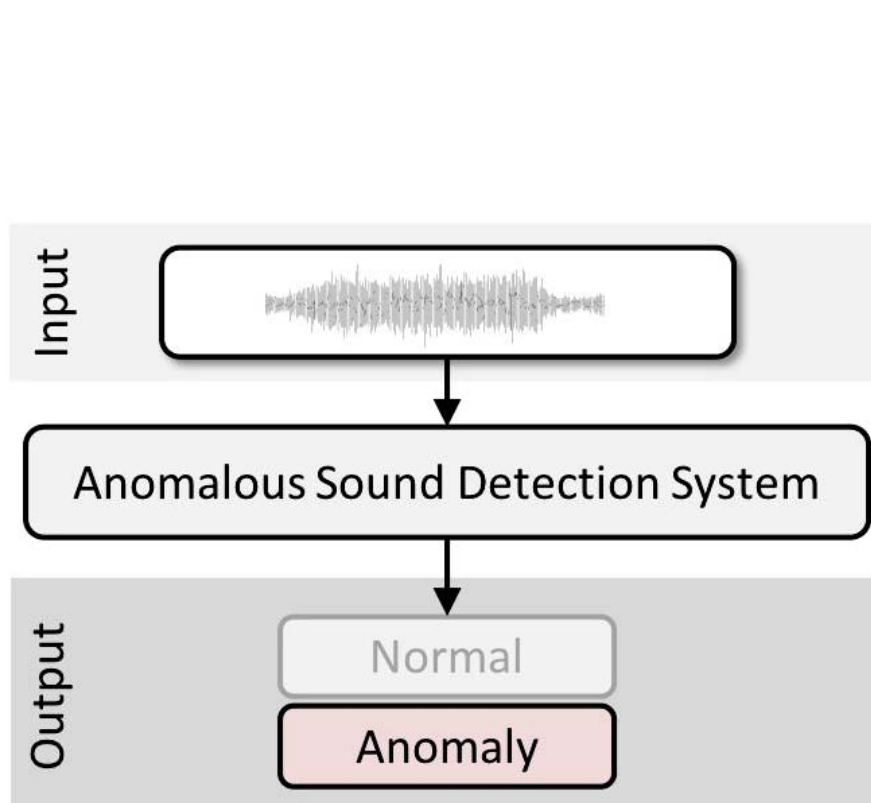No "cherry-picking"
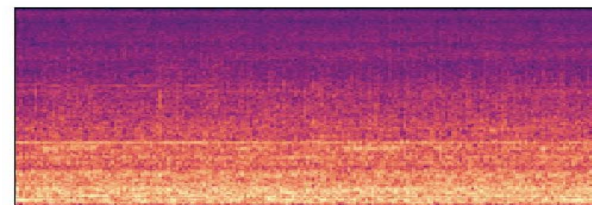Median result

# ALL TOGETHER: THE FULL BENCHMARK FLOW

**Workgroup**

**Submitter**

**Benchmark definition**
- Problem definition
- Dataset selection (public domain)
- Model selection / training

**Open** submission

**Reference implementation**
- Derive "Tiny" version: Quantization, etc.
- Write C implementation
- Integrate with benchmarking harness
- Deploy on device
- Benchmark
    - ML performance
    - Power consumption
- Review and Audit

**Closed** submission

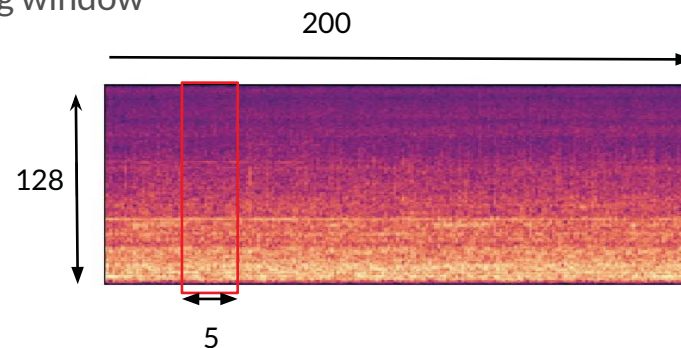# ALL TOGETHER: ANOMALY DETECTION EXAMPLE

# EXAMPLE: ANOMALY DETECTION
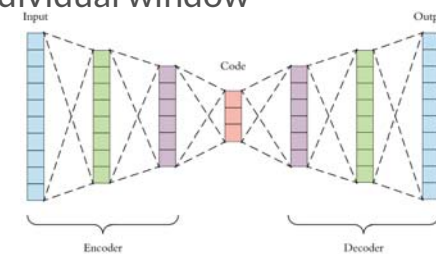
- Pre-processing: Spectrogram

- Sliding window

- Anomaly score on individual window AutoEncoder

- Post-processing: average score

# REFERENCE IMPLEMENTATION
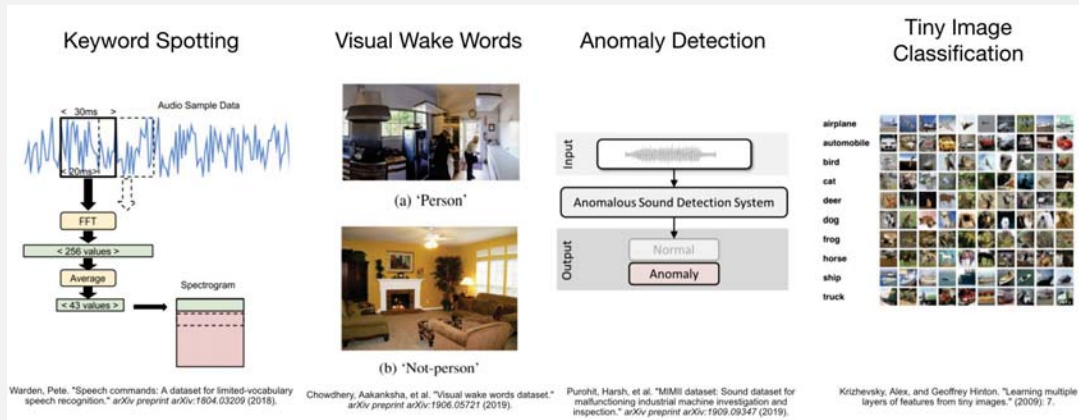
## ST NUCLEO-L4R5ZI



 +  + 

- Open platform
- Widely known, available, and affordable
- ARM Cortex-M4 with FPU
- 2 MB Flash
- 640KB SRAM

- Open source
- Portability of reference implementations

- Open source
- Full toolchain with everything we need
- Improving quantization support

# TINYMLPERF:WHERE NEXT?



V0.1.0 is out
bootstrapping the process

Just the beginning …
Grow the community
Refine and define use cases
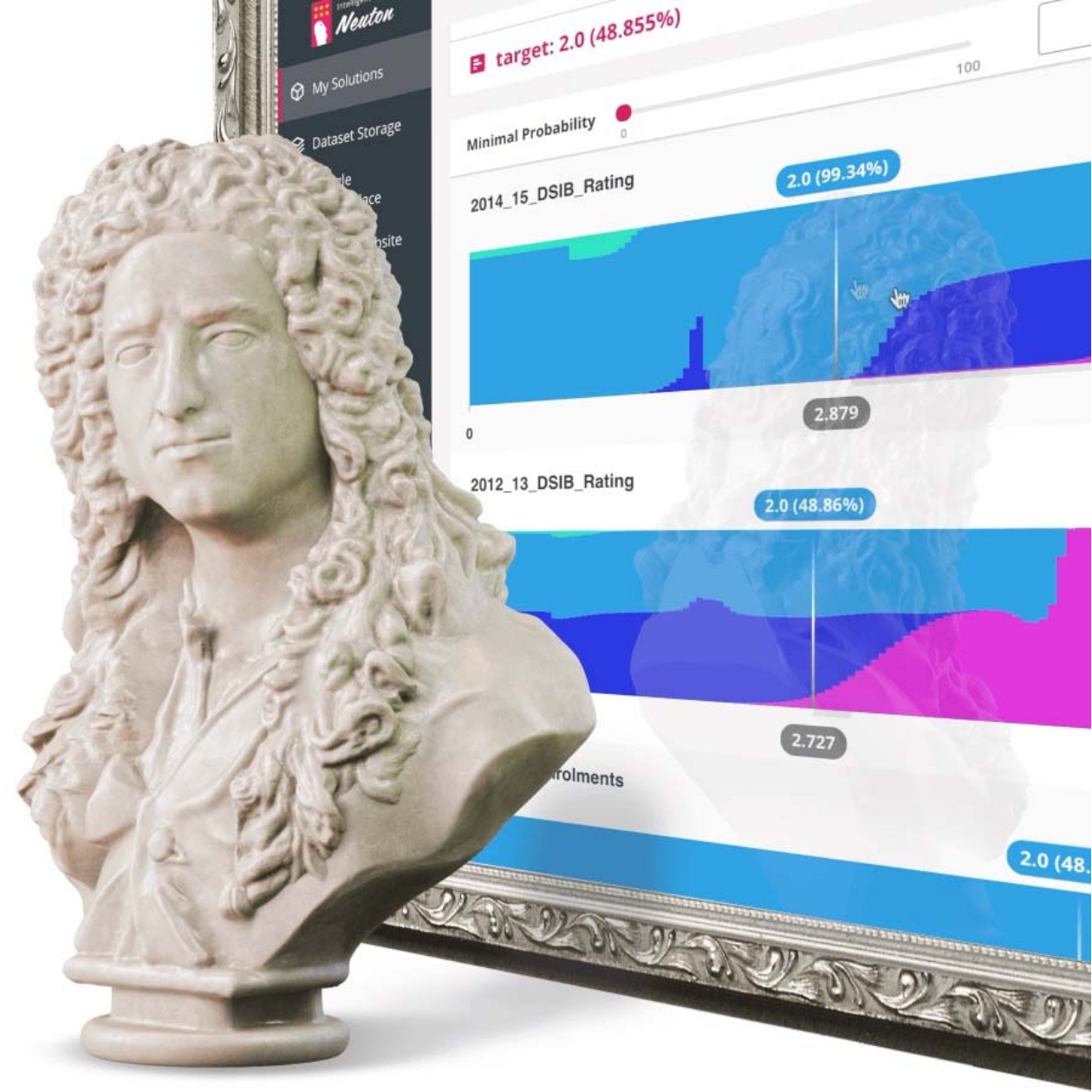towards V0.2.0

# Premier Sponsor

**Executive Sponsors**

# Arm: The Software and Hardware Foundation for tinyML

| | 1 | | | Application |
|---|---|---|---|---|

**Connect to high-level frameworks** — 1

**Supported by end-to-end tooling** — 2

**Connect to Runtime** — 3

Profiling and debugging tooling such as Arm Keil MDK

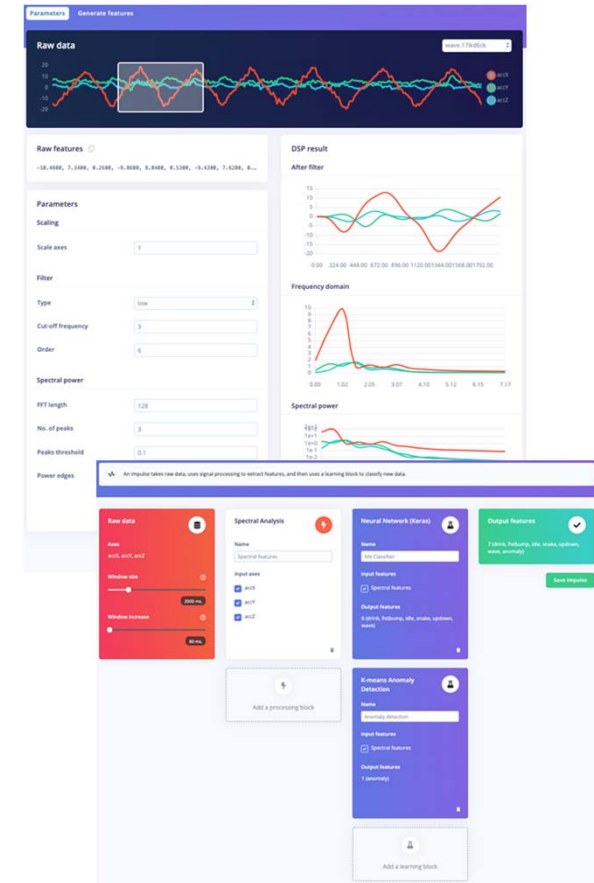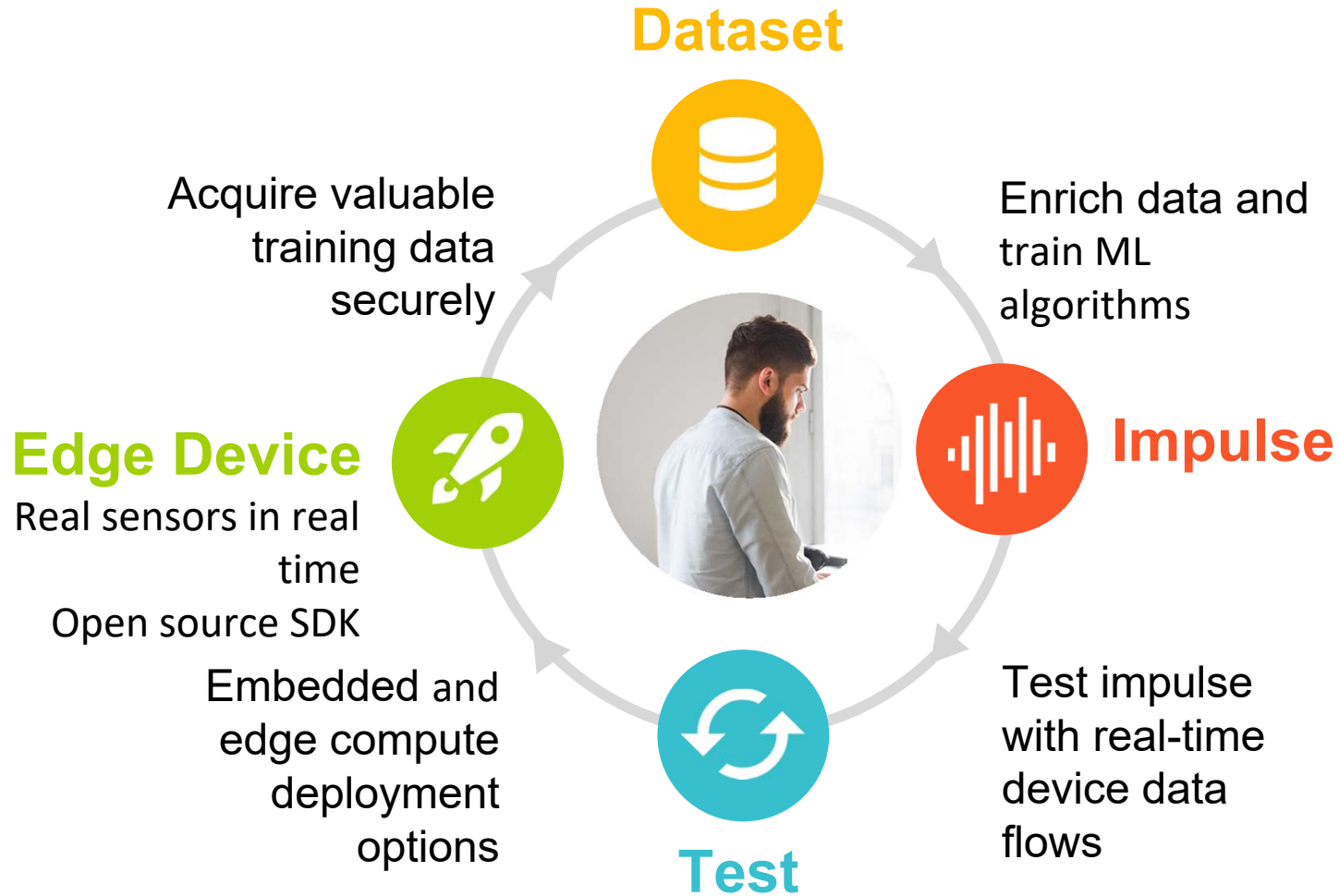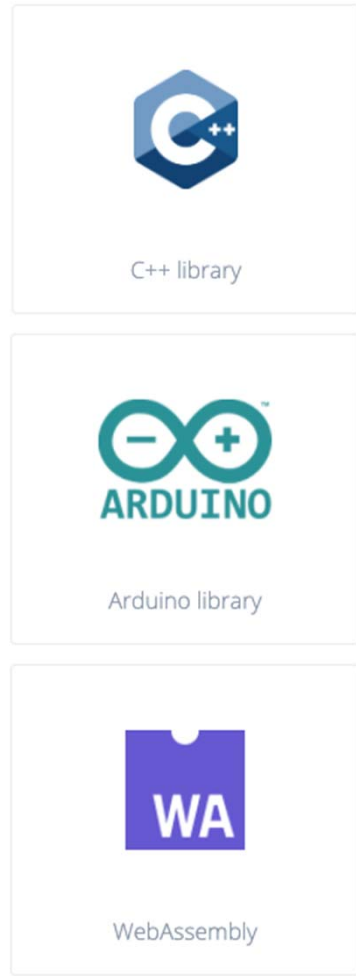| 1 | Application |
|---|---|
| 2 | Optimized models for embedded |
| 3 | Runtime (e.g. TensorFlow Lite Micro) |
| | Optimized low-level NN libraries (i.e. CMSIS-NN) |
| | RTOS such as Mbed OS |
| | Arm Cortex-M CPUs and microNPUs |

AI Ecosystem Partners

Stay Connected

▶ @ArmSoftwareDevelopers

🐦 @ArmSoftwareDev

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

arm

# TinyML for all developers

C++ library

Arduino library

WebAssembly

**Dataset**

Acquire valuable training data securely

Enrich data and train ML algorithms

**Edge Device**
Real sensors in real time
Open source SDK
Embedded and edge compute deployment options

**Impulse**

Test impulse with real-time device data flows

**Test**

www.edgeimpulse.com

**Advancing AI research to make efficient AI ubiquitous**

**Power efficiency**

Model design, compression, quantization, algorithms, efficient hardware, software tool

**Personalization**

Continuous learning, contextual, always-on, privacy-preserved, distributed learning

**Efficient learning**

Robust learning through minimal data, unsupervised learning, on-device learning

**A platform to scale AI across the industry**

**Perception**
Object detection, speech recognition, contextual fusion

**Reasoning**
Scene understanding, language understanding, behavior prediction

**Action**
Reinforcement learning for decision making

Edge cloud

Cloud

IoT/IIoT

Automotive

Mobile

# SYNTIANT

Syntiant Corp. is moving artificial intelligence and machine learning from the cloud to edge devices. Syntiant's chip solutions merge deep learning with semiconductor design to produce ultra-low-power, high performance, deep neural network processors. These network processors enable always-on applications in battery-powered devices, such as smartphones, smart speakers, earbuds, hearing aids, and laptops. Syntiant's Neural Decision Processors$^{TM}$ offer wake word, command word, and event detection in a chip for always-on voice and sensor applications.

Founded in 2017 and headquartered in Irvine, California, the company is backed by Amazon, Applied Materials, Atlantic Bridge Capital, Bosch, Intel Capital, Microsoft, Motorola, and others. Syntiant was recently named a CES® 2021 Best of Innovation Awards Honoree, shipped over 10M units worldwide, and unveiled the NDP120 part of the NDP10x family of inference engines for low-power applications.

www.syntiant.com          @Syntiantcorp

**Platinum Sponsors**

Part of your life. Part of tomorrow.

www.infineon.com

**Gold Sponsors**

# Build Smart IoT Sensor Devices From Data

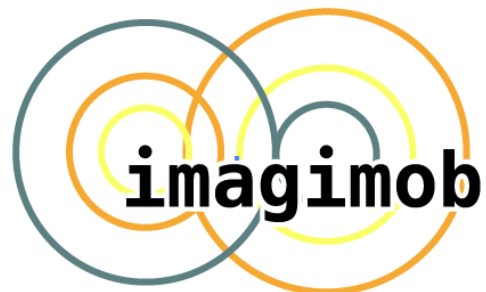SensiML pioneered TinyML software tools that auto generate AI code for the intelligent edge.

- End-to-end AI workflow
- Multi-user auto-labeling of time-series data
- Code transparency and customization at each step in the pipeline

We enable the creation of production-grade smart sensor devices.

**sensiml.com**

# Silver Sponsors

# Copyright Notice

## www.tinyML.org