# tinyML EMEA Technical Forum 2021 Proceedings

June 7 – 10, 2021

Virtual Event

www.tinyML.org

# tinyML EMEA Technical Forum 2021
## June 7-10, 2021

**Runtime DNN Performance Scaling though Resource Management on Heterogeneous Embedded Platforms**

Lei Xun, PhD Candidate, University of Southampton, UK

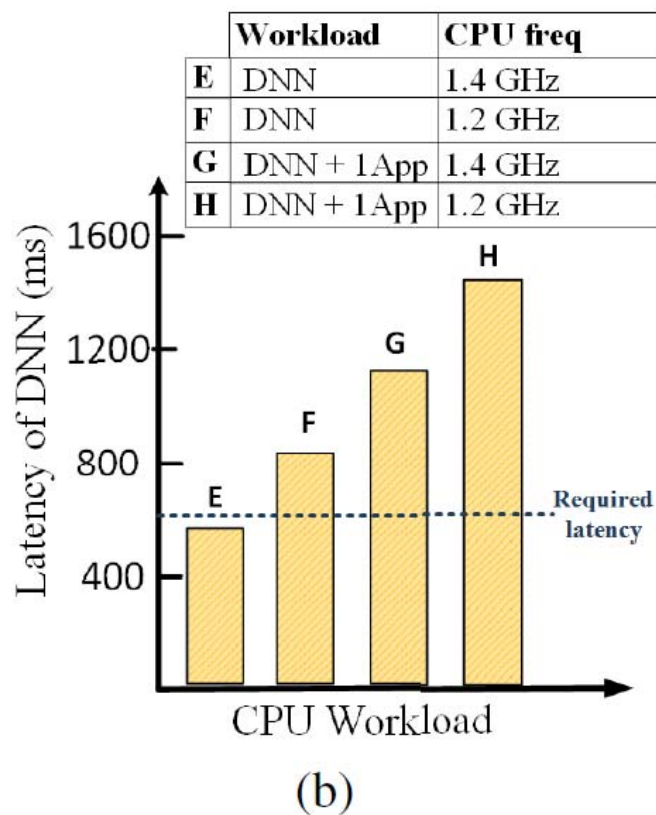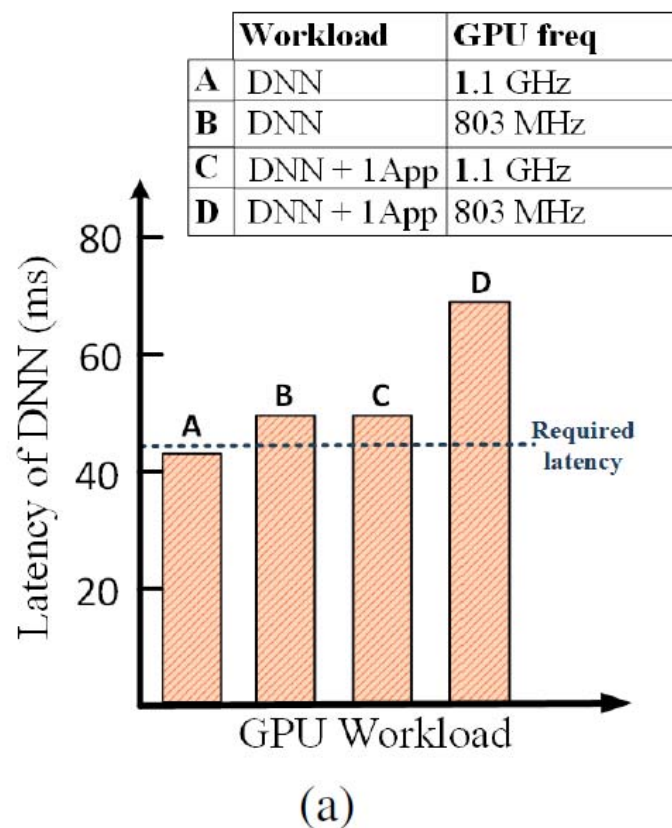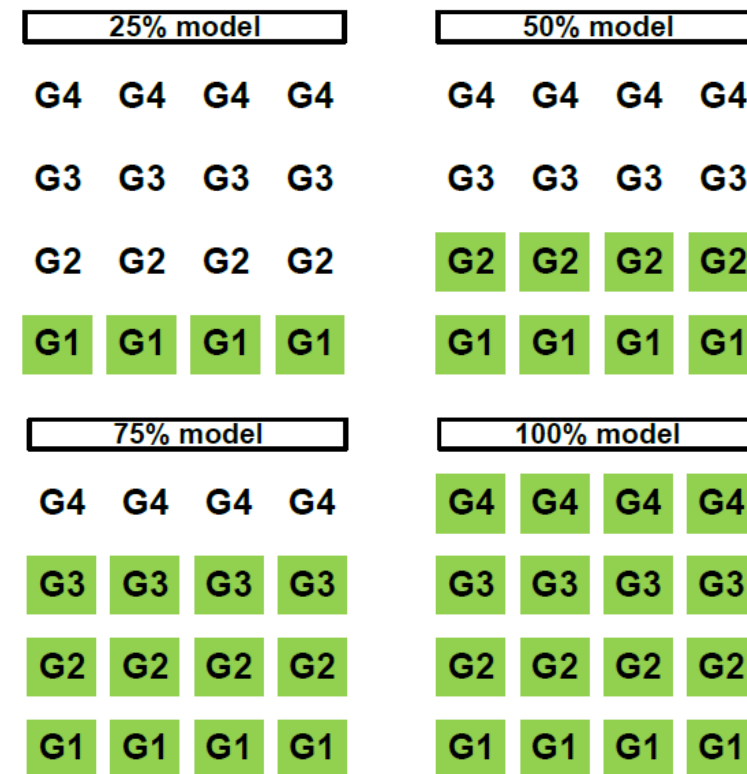Advisor: Geoff V. Merrett, Jonathon Hare, Bashir M Al-Hashimi

June 10, 2021

# Motivation for dynamic DNNs

- DNNs are typically compressed before deployed on embedded platform

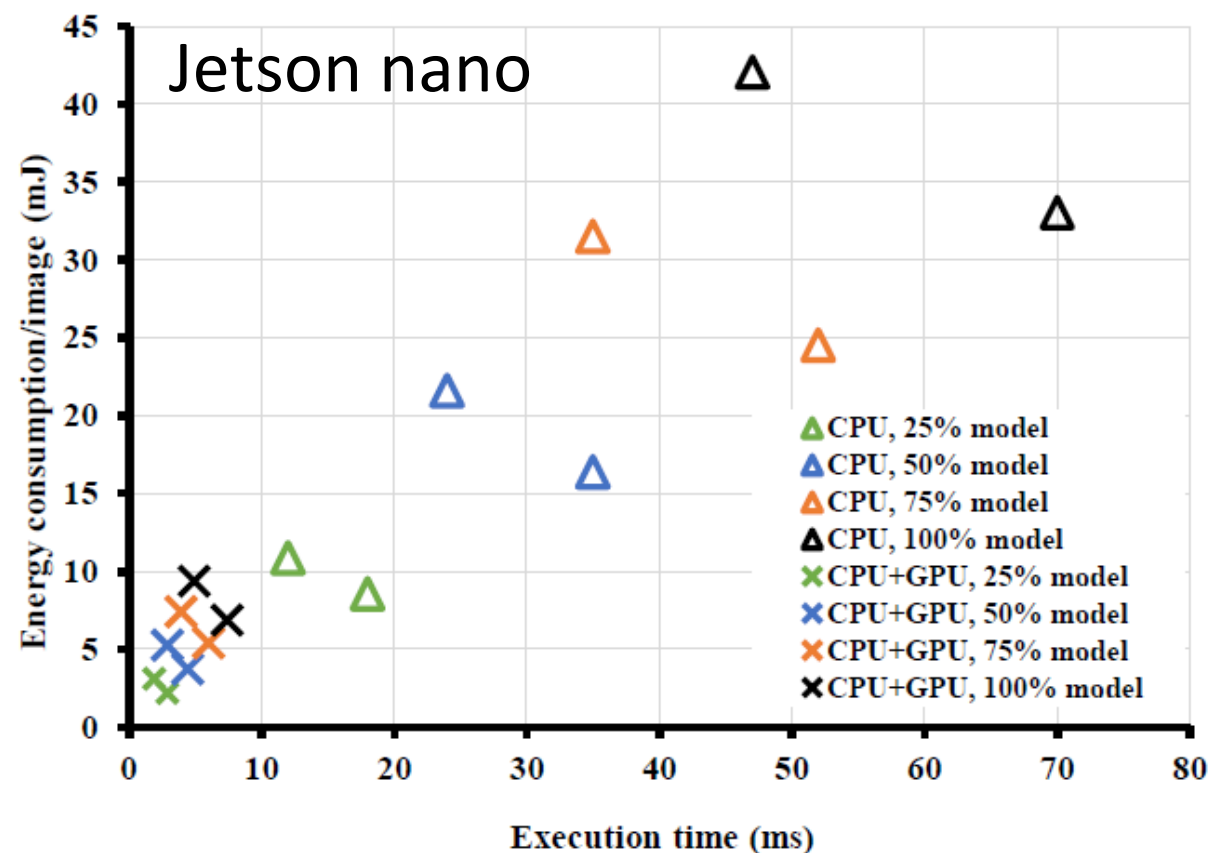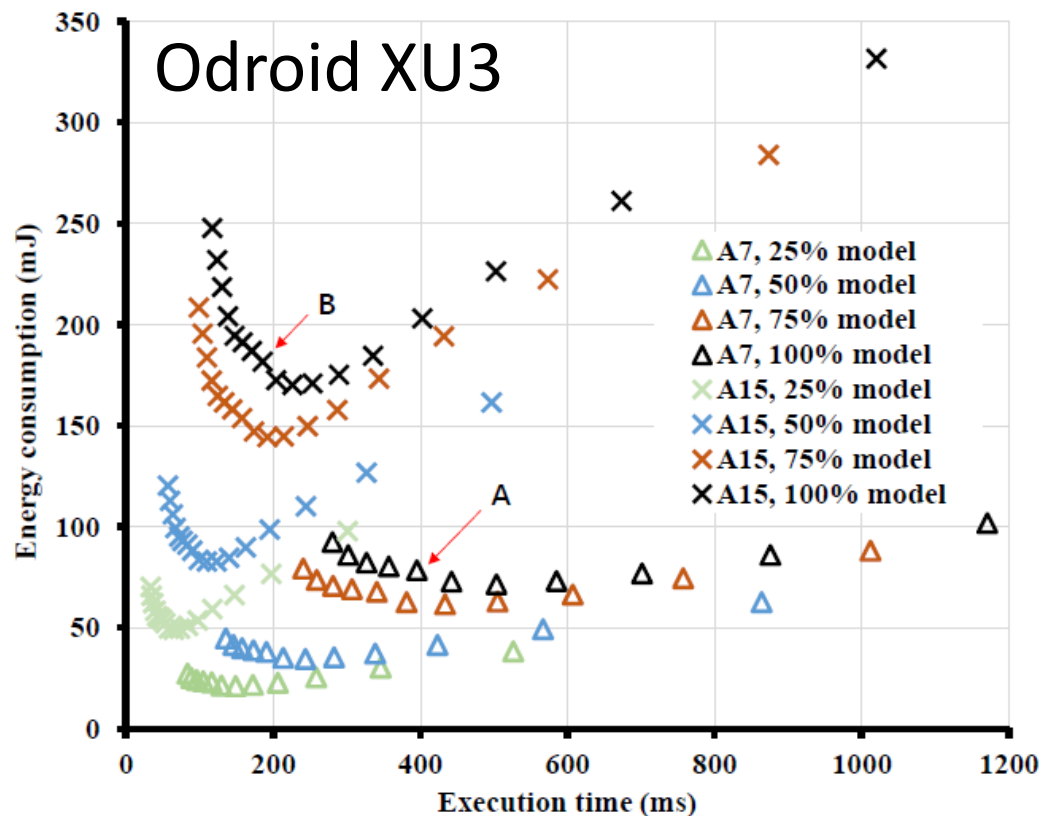- However, the assumed hardware resources may not be available at runtime

| | Workload | GPU freq |
|---|---|---|
| A | DNN | 1.1 GHz |
| B | DNN | 803 MHz |
| C | DNN + 1App | 1.1 GHz |
| D | DNN + 1App | 803 MHz |

| | Workload | CPU freq |
|---|---|---|
| E | DNN | 1.4 GHz |
| F | DNN | 1.2 GHz |
| G | DNN + 1App | 1.4 GHz |
| H | DNN + 1App | 1.2 GHz |



(a)

(b)

- Dynamic DNNs can be executed partially to trade-off accuracy for latency/power/energy reduction

Lei Xun, Long Tran-Thanh, Bashir M Al-Hashimi, and Geoff V Merrett (2019) *Incremental training and group convolution pruning for runtime DNN performance scaling on heterogeneous embedded platforms.* in **Workshop on Machine Learning for CAD (MLCAD).**

# Accuracy/latency/energy trade-offs

- Subnetworks are shown in different colors, computing elements are shown in different symbols) and frequency scaling are shown in points
- Operating points example: On Odroid XU3, A has the best trade-off under 100mJ and 400ms requirements, B is the best for 200mJ and 200ms
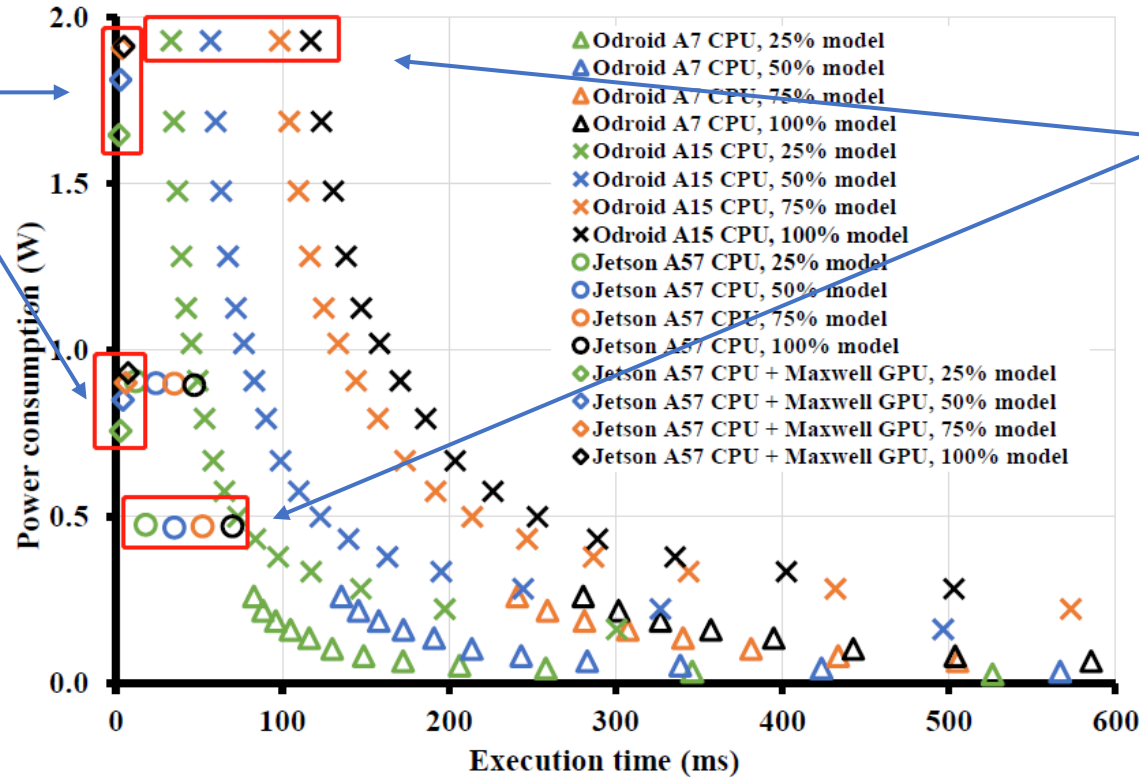


Lei Xun, Long Tran-Thanh, Bashir M Al-Hashimi, and Geoff V Merrett (2020) *Optimising resource management for embedded machine learning*. In **Design, Automation and Test in Europe Conference (DATE).**

# Accuracy/latency/power trade-offs

Now we know that dynamic DNNs can trade-off accuracy for latency and energy, what about power?

The power of GPU

- Scales with frequency scaling

- Scales with dynamic DNN, this provide us new opportunities to meet power target

The power of a single-core CPU:

- Scales with frequency scaling

- Does not scale with dynamic DNN, since the computation intensity does not change, only the latency changes



△ Odroid A7 CPU, 25% model
△ Odroid A7 CPU, 50% model
△ Odroid A7 CPU, 75% model
△ Odroid A7 CPU, 100% model
✕ Odroid A15 CPU, 25% model
✕ Odroid A15 CPU, 50% model
✕ Odroid A15 CPU, 75% model
✕ Odroid A15 CPU, 100% model
○ Jetson A57 CPU, 25% model
○ Jetson A57 CPU, 50% model
○ Jetson A57 CPU, 75% model
○ Jetson A57 CPU, 100% model
◇ Jetson A57 CPU + Maxwell GPU, 25% model
◇ Jetson A57 CPU + Maxwell GPU, 50% model
◇ Jetson A57 CPU + Maxwell GPU, 75% model
◇ Jetson A57 CPU + Maxwell GPU, 100% model

To know more about dynamic DNNs, please check out our SOTA dynamic DNN paper:
Wei Lou*, Lei Xun*, Mohammadamin Sabetsarvestani, Jia Bi, Jonathon Hare, Geoff V Merrett (2021) *Dynamic-OFA: Runtime DNN architecture switching for performance scaling on heterogeneous embedded platforms*. In **Conference on Computer Vision and Pattern Recognition Workshops (CVPR'W).** https://arxiv.org/abs/2105.03596

# Premier Sponsor

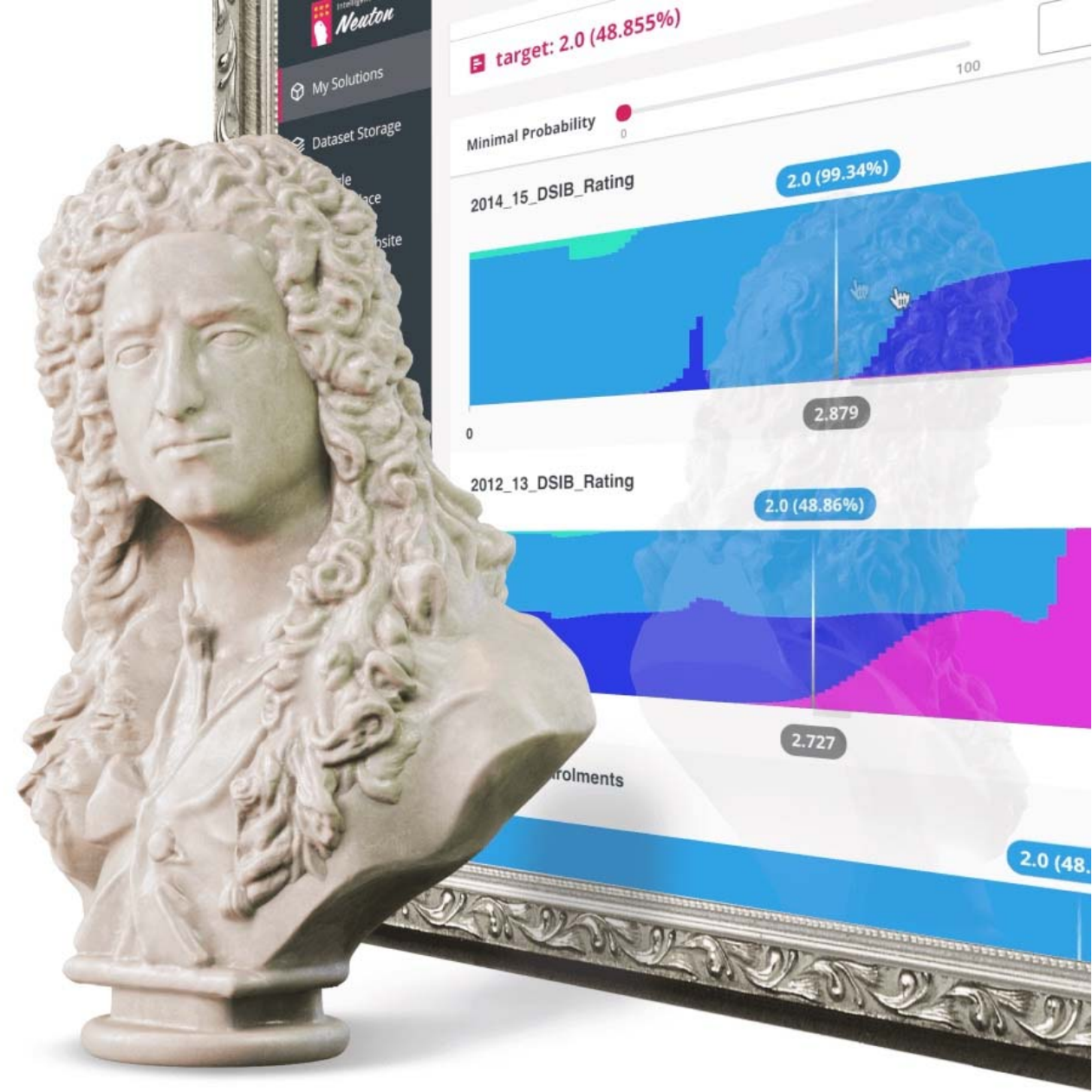![Intelligent Agent Neuton logo]

# Automated TinyML

Zero-code SaaS solution

**Create tiny models, ready for embedding, in just a few clicks!**

Compare the benchmarks of our compact models to those of TensorFlow and other leading neural network frameworks.
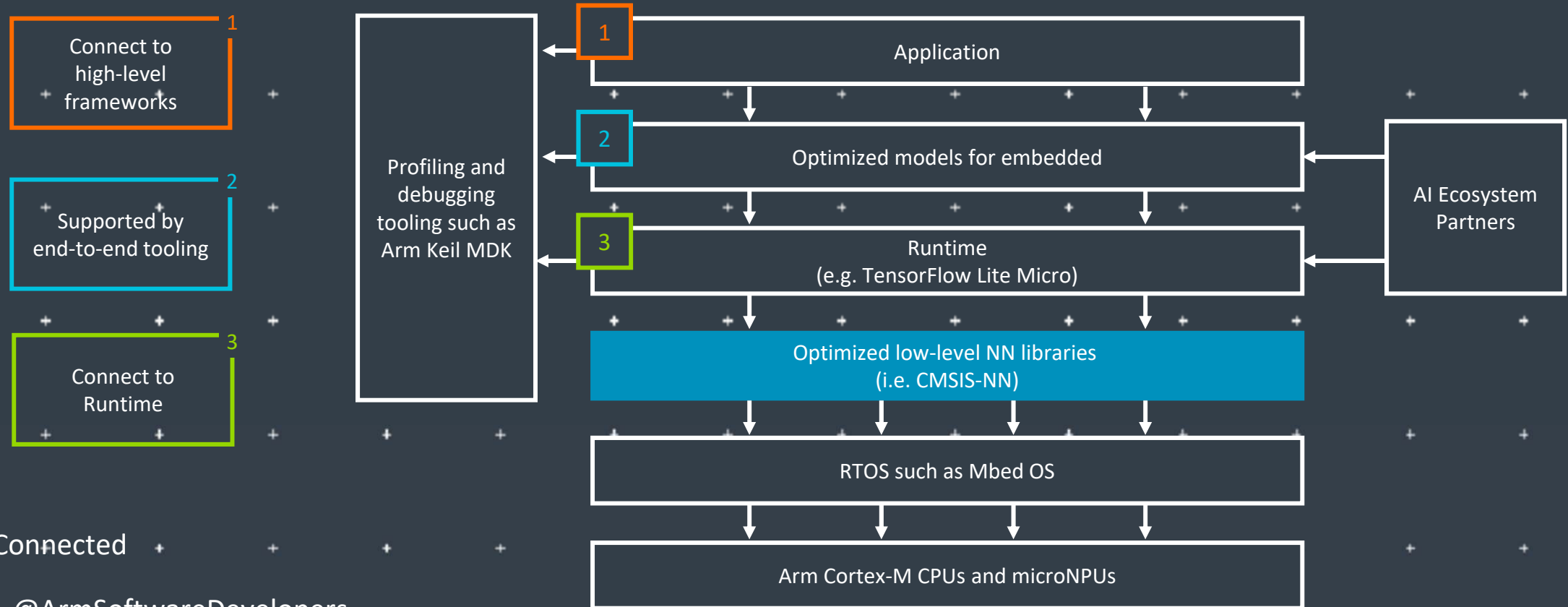
*Build Fast. Build Once. Never Compromise.*

# Executive Sponsors

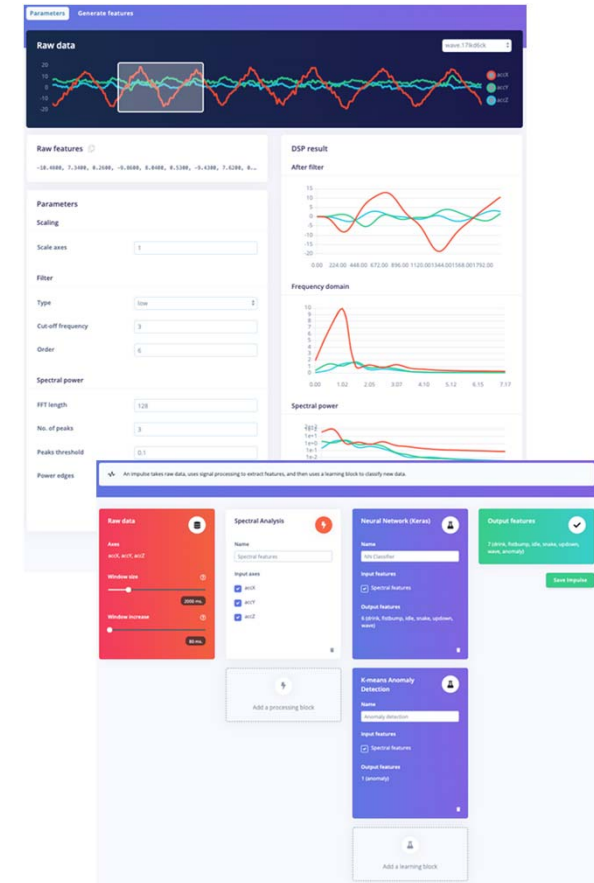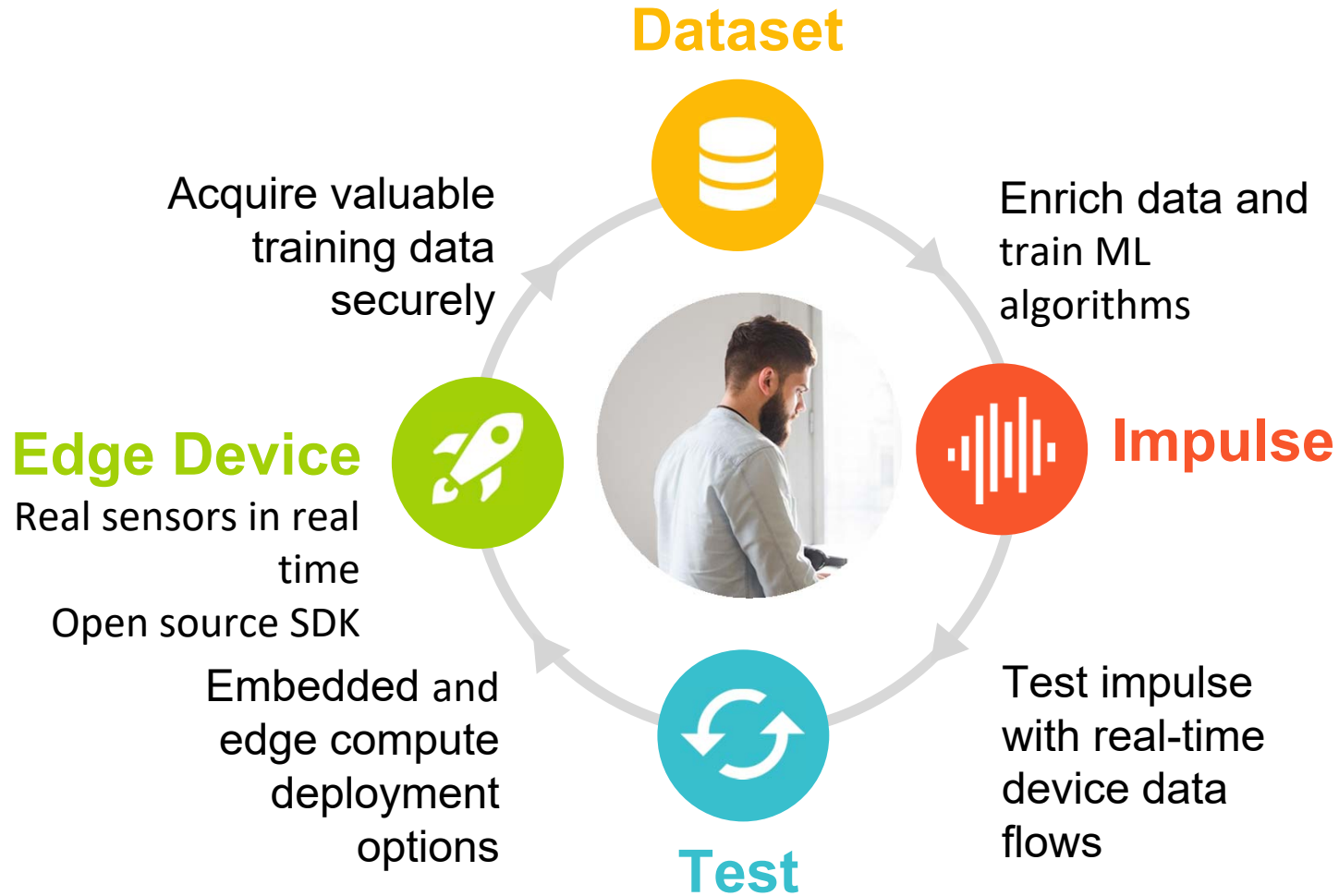# Arm: The Software and Hardware Foundation for tinyML

**1** Connect to high-level frameworks
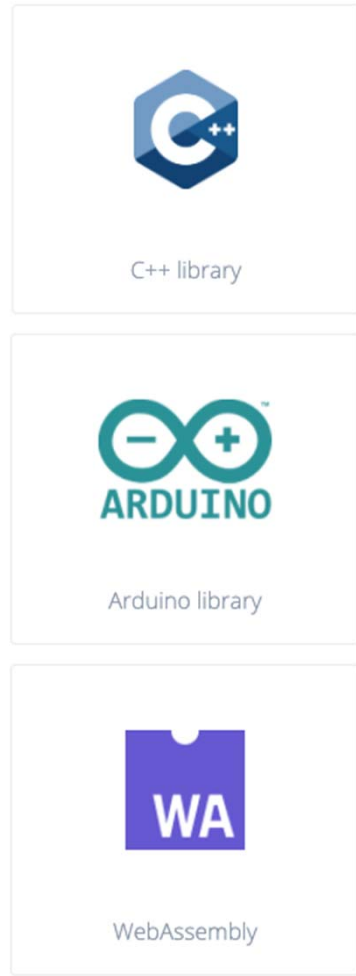
**2** Supported by end-to-end tooling

**3** Connect to Runtime

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

**Impulse**

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

Test impulse with real-time device data flows

**Test**

www.edgeimpulse.com

# Advancing AI research to make efficient AI ubiquitous

**Qualcomm AI research**

### 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

IoT/IIoT

Automotive

Cloud

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™ 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

# RealityAI®

## Add Advanced Sensing to your Product with Edge AI / TinyML

https://reality.ai    info@reality.ai    @SensorAI    Reality AI

## Pre-built Edge AI sensing modules, plus tools to build your own

### Reality AI solutions

> Prebuilt sound recognition models for indoor and outdoor use cases

> Solution for industrial anomaly detection

> Pre-built automotive solution that lets cars "see with sound"

### Reality AI Tools® software

> Build prototypes, then turn them into real products

> Explain ML models and relate the function to the physics

> Optimize the hardware, including sensor selection and placement

# Gold Sponsors

# LatentAI

## Adaptive AI for the Intelligent Edge

Latentai.com

# Build Smart IoT Sensor Devices From Data

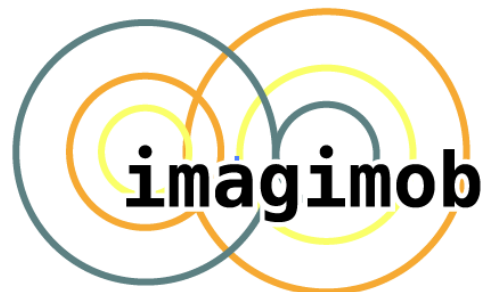SensiML pioneered TinyML software tools that auto generate AI code for the intelligent edge.

- 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