# tinyML Talks Strategic Partners
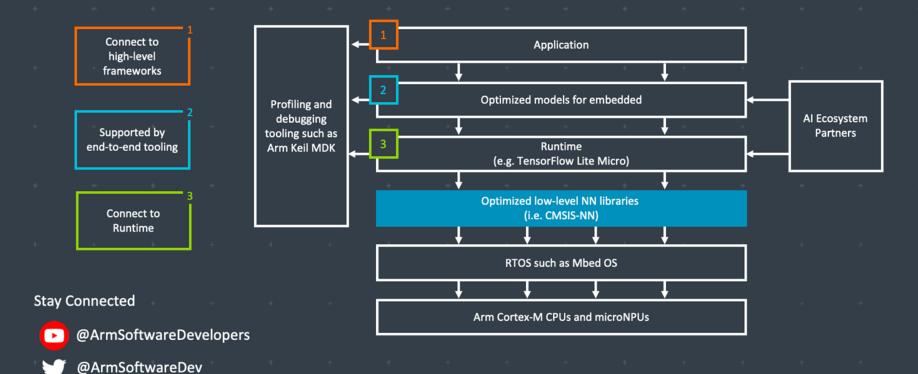
Additional Sponsorships available – contact Olga@tinyML.org for info

# 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**

# WE USE AI TO MAKE OTHER AI FASTER, SMALLER AND MORE POWER EFFICIENT

**Automatically compress** SOTA models like MobileNet to <200KB with **little to no drop in accuracy** for inference on resource-limited MCUs

**Reduce** model optimization trial & error from weeks to days using Deeplite's **design space exploration**
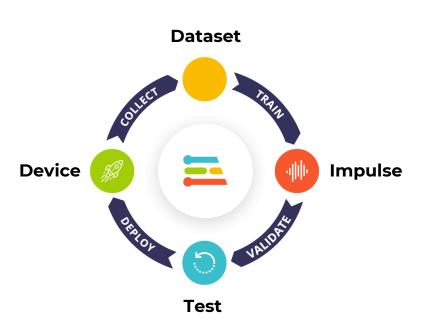
**Deploy more** models to your device without sacrificing performance or battery life with our **easy-to-use software**
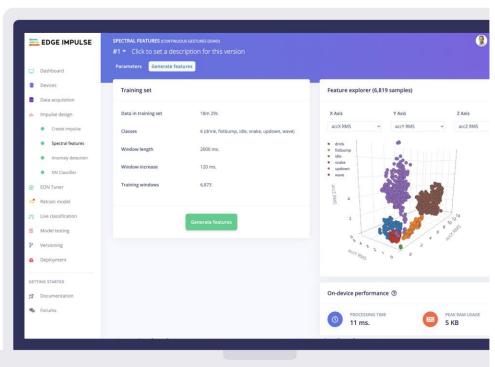
BECOME BETA USER **bit.ly/testdeeplite**

mobility**Xlab**  **arm**  AI 100 CBINSIGHTS

emza
visual sense

# The Eye in IoT
## Edge AI Visual Sensors

info@emza-vs.com

**CMOS Imaging Sensor**

- Ultra Low power CMOS imager
- Ai + IR capable

**Computer Vision Algorithms**

**IoT System on Chip**

WiseEye

- Machine Learning edge computing silicon
- <1mW always-on power consumption
- Computer Vision hardware accelerators

- Machine Learning algorithm
- <1MB memory footprint
- Microcontrollers computing power
- Trained algorithm
- Processing of low-res images
- Human detection and other classifiers

# Enabling the next generation of Sensor and Hearable products to process rich data with energy efficiency



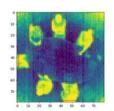Visible Image

Sound

IR Image

Radar

Bio-sensor

Gyro/Accel

Wearables / Hearables

Battery-powered consumer electronics

IoT Sensors

GREENWAVES TECHNOLOGIES

# ⚡Grovety Inc.

## SOFTWARE DEVELOPMENT SERVICES FOR TINYML SOLUTIONS

**1**

**Development tools**
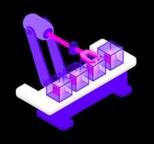SDK, IDE, compilers, leveraging on TVM, uTVM & LLVM
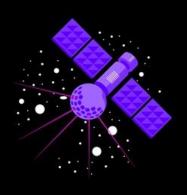
**2**

**Firmware**
Drivers, BSP, protocols, etc.

**arm** AI PARTNER

# Distributed infrastructure for TinyML apps

**HOTG**
Decoupling intelligence

**Develop at warp speed**

**Automate deployments**

**Device orchestration**

**HOTG is building the distributed infrastructure to pave the way for AI enabled edge applications**

# LatentAI

## Adaptive AI for the Intelligent Edge

Latentai.com

# maxim integrated™

# Maxim Integrated: Enabling Edge Intelligence
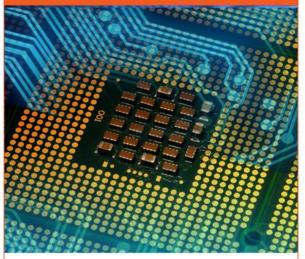
## Advanced AI Acceleration IC

The new MAX78000 implements AI inferences at low energy levels, enabling complex audio and video inferencing to run on small batteries. Now the edge can see and hear like never before.

www.maximintegrated.com/MAX78000

## Low Power Cortex M4 Micros

Large (3MB flash + 1MB SRAM) and small (256KB flash + 96KB SRAM, 1.6mm x 1.6mm) Cortex M4 microcontrollers enable algorithms and neural networks to run at wearable power levels.

www.maximintegrated.com/microcontrollers

## Sensors and Signal Conditioning

Health sensors measure PPG and ECG signals critical to understanding vital signs. Signal chain products enable measuring even the most sensitive signals.

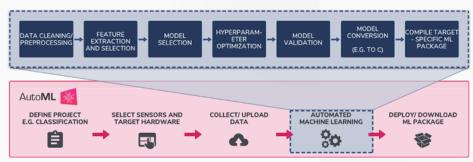www.maximintegrated.com/sensors

maxim integrated™

# Qeexo AutoML

Automated Machine Learning Platform that builds tinyML solutions for the Edge using sensor data

## Key Features

- Supports 17 ML methods:
  - Multi-class algorithms:  GBM, XGBoost, Random Forest, Logistic Regression, Gaussian Naive Bayes, Decision Tree, Polynomial SVM, RBF SVM, SVM, CNN, RNN, CRNN, ANN
  - Single-class algorithms: Local Outlier Factor, One Class SVM, One Class Random Forest, Isolation Forest
- Labels, records, validates, and visualizes time-series sensor data
- On-device inference optimized for low latency, low power consumption, and small memory footprint applications
- Supports Arm® Cortex™- M0 to M4 class MCUs

## End-to-End Machine Learning Platform



For more information, visit: www.qeexo.com

## Target Markets/Applications

- Industrial Predictive Maintenance
- Smart Home
- Wearables
- Automotive
- Mobile
- IoT

**Qualcomm**
AI research

# 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

# 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

# BROAD AND SCALABLE EDGE COMPUTING PORTFOLIO

## Microcontrollers & Microprocessors

**Arm® Core**

**RENESAS RA**
Arm® Cortex®-M 32-bit MCUs
Arm ecosystem, advanced security, Intelligent IoT

**RENESAS RZ**
Arm®-based High-end 32 & 64-bit MPUs
High-resolution HMI, Industrial network & real-time control

**RENESAS RE**
Arm® Cortex®-M0+ Ultra-low Power 32-bit MCUs
Innovative process tech (SOTB), Energy harvesting

**Renesas Synergy™**
Arm®-based 32-bit MCUs for Qualified Platform
Qualified software and tools

**Renesas Core**

**RENESAS RL78**
Ultra-low Energy 8 & 16-bit MCUs
Bluetooth® Low Energy, SubGHz, LoRa®-based Solutions

**RENESAS RX**
High Power Efficiently 32-bit MCUs
Motor control, Capacitive touch, Functional safety, GUI

**RENESAS RH850**
40nm/28nm process Automotive 32-bit MCUs
Rich functional safety and embedded security features

## Core technologies

**AI**
A broad set of high-power and energy-efficient embedded processors

**INNOVATION**

**Security & Safety**
Comprehensive technology and support that meet the industry's stringent standards

**Digital & Analog & Power Solution**
Winning Combinations that combine our complementary product portfolios

**Cloud Native**
Cross-platforms working with partners in different verticals and organizations

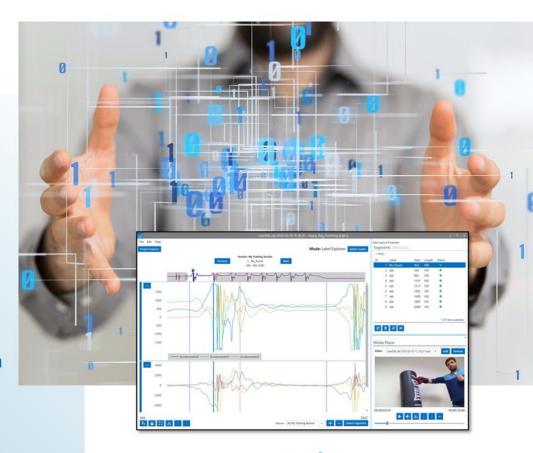BIG IDEAS FOR EVERY SPACE **RENESAS**

seeed studio

The IoT Hardware Enabler

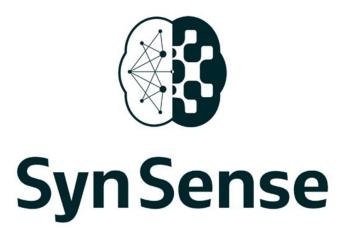# Build Smart IoT Sensor Devices From Data

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

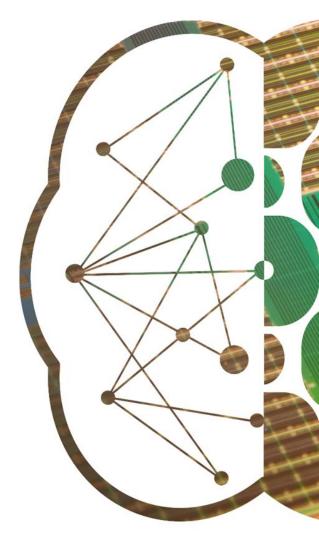We enable the creation of production-grade smart sensor devices.

sensiml.com

# SynSense

**SynSense** builds **sensing and inference** hardware for **ultra-low-power** (sub-mW) **embedded, mobile and edge** devices. We design systems for **real-time always-on smart sensing**, for audio, vision, IMUs, bio-signals and more.

https://SynSense.ai

# Next tinyML Talks

| Date | Presenter | Topic / Title |
|------|-----------|---------------|
| Tuesday, January 11 | Tim Callahan<br>Staff Software Engineer, Google | CFU Playground: Customize Your ML Processor for Your Specific TinyML Model |

Webcast start time is 8:00 am Pacific time

Please contact talks@tinyml.org if you are interested in presenting

# tinyML Summit 2022
## Miniature dreams can come true...
March 28-30, 2022
Hyatt Regency San Francisco Airport
https://www.tinyml.org/event/summit-2022/

Registration will be open on **December 15**, 2021.

Deadline for poster submission is **December 17**.

*The Best Product of the Year and the Best Innovation of the Year awards are open for nominations between **November 15  and February 28**.*

# tinyML Research Symposium 2022

March 28, 2022

https://www.tinyml.org/event/research-symposium-2022

Call for papers – Submission deadline is **December 17**, 2021.

More sponsorships are available: sponsorships@tinyML.org

# Reminders

Slides & Videos will be posted tomorrow

Please use the Q&A window for your questions

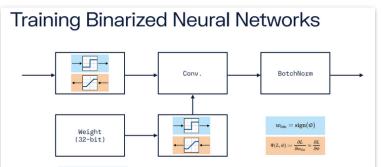tinyml.org/forums    youtube.com/tinyml

# Cedric Nugteren



Cedric Nugteren is a software engineer focussed on writing efficient code for deep learning applications. After he received his MSc and PhD from Eindhoven University of Technology he optimized GPU and CPU code for various companies using C++, OpenCL and CUDA. Then, he worked for 4 years on deep learning for autonomous driving at TomTom, after which he joined Plumerai where he is now writing fast code for the smallest microcontrollers.

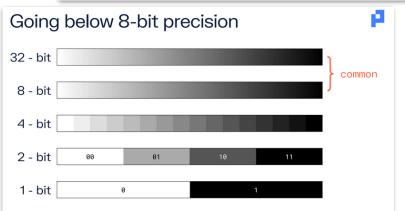# Demoing the world's fastest inference engine for Arm Cortex-M

Plumerai

Cedric Nugteren - cedric@plumerai.com

# You might know us from: BNNs?

## Training Binarized Neural Networks
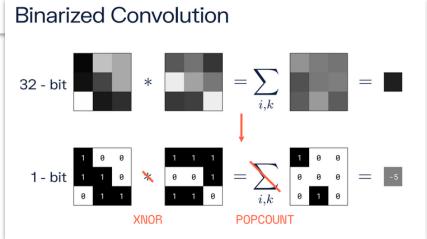


- On the forward pass, weights are binarized
- On the backward pass, a Straight-Through Estimator approximates the gradient

$w_{\text{bin}} = \text{sign}(\tilde{w})$

$\Phi(L, \tilde{w}) := \frac{\partial L}{\partial w_{\text{bin}}} \approx \frac{\partial L}{\partial \tilde{w}}$

Helwegen et al., 2019, NeurIPS – Latent Weights Do Not Exist: Rethinking Binarized Neural Network Optimization

## Binarized Neural Networks for efficient deep learning

Larq is an ecosystem of open-source Python packages for building, training and deploying Binarized Neural Networks to enable efficient inference on mobile and edge devices.

Get started with Larq

## Going below 8-bit precision



32 - bit

8 - bit

} common

4 - bit

2 - bit  00  01  10  11

1 - bit  0  1

## Binarized Convolution



32 - bit  *  = $\sum_{i,k}$  =

1 - bit  *  = $\sum_{i,k}$  = -5

XNOR    POPCOUNT

# You might know us from: Person detection?

Person Presence Detection

Example:
Smart doorbell

Example:
Smart offices

On STM32...

Binary size: 895 KiB     Binary size: 895 KiB

# You might know us from: Our own IP core?



Tuesday, November 9, 2021

Person presence detection on the Lattice CrossLink-NX FPGA



The Lattice CrossLink-NX Voice & Vision Machine Learning Board with the CrossLink-NX LIFCL-40 FPGA

Camera

Ikva IP core

# Or from: the world's fastest Cortex-M inference?

Monday, October 4, 2021

## The world's fastest deep learning inference software for Arm Cortex-M

New: Try out our inference engine with your own model!

At Plumerai we enable our customers to [...] on tiny embedded hardware. We're proud to announce that our infere[...] [...]ontrollers is the fastest and most memory-efficient in the world, for both B[...] deep learning models. Our inference software is an essential component of our solution, since it directs resource management akin to an operating system. It has **40% lower latency** and requires **49% less RAM** than TensorFlow Lite for Microcontrollers with Arm's CMSIS-NN kernels while retaining the same accuracy. It also outperforms any other deep learning inference software for Arm Cortex-M:

| | Inference time | RAM usage |
|---|---|---|
| TensorFlow Lite for Microcontrollers 2.5 (with CMSIS-NN) | 129 ms | 155 KiB |
| Edge Impulse's EON | 120 ms | 153 KiB |
| MIT's TinyEngine [1] | 124 ms | 98 KiB |
| STMicroelectronics' X-CUBE-[...] | 103 ms | 109 KiB |
| → **Plumerai's inference software** | **77 ms** | **80 KiB** |

**1. What is an inference engine?**

**3. Live demo of public benchmarking service**

**4. What did we do to become so efficient?**

**2. Are we really that efficient?**

28

0. How did we get here?

# How did we get here?

| Early days of Plumerai | Quite soon after | Last year | Today |

Our goal: run complex computer vision tasks on tiny devices efficiently

Need to cover the entire stack for high efficiency

BNNs do not only have binarized layers, but also INT8…

The world's fastest INT8 inference engine



EfficientDet-level accuracy on an MCU.

$3 MCU

$10,000 GPU

Going below 8-bit precision

| 32 - bit | | common |
| 8 - bit | |
| 4 - bit | |
| 2 - bit | 00 01 10 11 |
| 1 - bit | 0 1 |

| Plumerai Data Pipeline |
| Plumerai BNN Models |
| Plumerai Inference Stack |
| Plumerai Hardware |

Faster models

Less RAM

- Better accuracy
- Lower energy usage
- Cheaper and smaller
- Room for other apps

# 1. What is an inference engine?

# The machine learning flow

### Pick a model
Pick a new model or retrain an existing one.

### Convert
Convert a TensorFlow model into a compressed flat buffer with the TensorFlow Lite Converter.

### Deploy
Take the compressed .tflite file and load it into a mobile or embedded device.

**Deploy INT8 quantized model on device**

**Run optimized code**

# The tasks of an inference engine



1. Execute the layers of the model in the correct order

2. Plan the activations and weights in memory efficiently

3. Provide optimized INT8 code for each layer type (e.g. convolution, fully connected)

# An inference engine example: TFLM

### TensorFlow Lite for Microcontrollers

On this page
Why microcontrollers are important
Supported platforms
Explore the examples
Workflow
Limitations
Next steps

TensorFlow Lite for Microcontrollers is designed to run machine learning models on microcontrollers and other devices with only few kilobytes of memory. The core runtime just fits in 16 KB on an Arm Cortex M3 and can run many basic models. It doesn't require operating system support, any standard C or C++ libraries, or dynamic memory allocation.

1. Interpreter: TensorFlow Lite for Microcontrollers

2. Memory planner: TensorFlow Lite for Microcontrollers

3. Optimized INT8 code: ARM CMSIS-NN

Neural Network Application Code

| Convolution | Pooling |
| Fully connected | Activations |

NNFunctions

| Data type conversion |
| Activation tables |

NNSupportFunctions

# 2. Are we really that efficient?

# A closer look at the results

What most people are using

Our inference engine

| | Inference time | RAM usage |
|---|---|---|
| TensorFlow Lite for Microcontrollers 2.5 (with CMSIS-NN) | 129 ms | 155 KiB |
| Edge Impulse's EON | 120 ms | 153 KiB |
| MIT's TinyEngine [1] | 124 ms | 98 KiB |
| STMicroelectronics' X-CUBE-AI | 103 ms | 109 KiB |
| **Plumerai's inference software** | **77 ms** | **80 KiB** |

**Model**: MobileNetV2 [2] [3] (alpha=0.30, resolution=80x80, classes=1000)

**Board**: STM32F746G-Discovery at 216 MHz with 320 KiB RAM and 1 MiB flash

Also tested: microTVM, but ran out of memory

**No tricks:** no binarization or pruning, accuracy remains the same in this table

# Just good on MobileNetV2?



TFLite for M... ...CMSIS-...

Edge Impulse's EON

MIT's TinyEngine [1]

STMicroelectronics' X...

→ **Plumerai's inference**

**Model**: MobileNetV2 [2] ...tion=80x80, classes=1000)

**Board**: STM32F746G-... ...th 320 KiB RAM and 1 MiB flash

```
arm_ml_zoo:
https://github.com/ARM-software/ML-zoo

tflm_examples:
https://github.com/tensorflow/tflite-micro/tree/main/tensorflow/lite/micro/examples/

mcunet:
https://github.com/mit-han-lab/tinyml/tree/master/mcunet

philippvk_stm_demos:
https://github.com/PhilippvK/stm32-tflm-demos/tree/main/examples

himax_yolo:
https://github.com/HimaxWiseEyePlus/Yolo-Fastest/tree/master/ModelZoo

tfhub:
https://tfhub.dev/s?deployment-format=lite&tf-version=tf2
https://www.tensorflow.org/lite/guide/hosted_models
https://www.tensorflow.org/lite/examples/object_detection/overview#performance_benchmarks

tf_tpu:
https://github.com/tensorflow/tpu/tree/master/models/official/

keras_applications (from code with 96x96x3 input and no 'top'):
https://keras.io/api/applications/

plumerai:
https://blog.plumerai.com/2021/10/cortex-m-inference-software/

mlperftiny:
https://github.com/mlcommons/tiny/tree/master/v0.5/training
```

# More off-the-shelf models

# More off-the-shelf models



Average speed-up factor: 1.53x

On Cortex-M7

Average RAM reduction factor: 1.45x

**Remember**: accuracy remains the same, only speed and memory requirements change

# A closer look at the MLPerf Tiny models

# 3. Live demo of
# public benchmarking service

# Public benchmarking service: try it yourself!



**Any model that runs with TFLM is supported!**

**Note:** M4 & M7 and ST boards are just examples: our inference engine runs across many platforms

Visit https://plumerai.com/benchmark to try it with your own model

# 4. What did we do to become so efficient?

# How to beat the competition?



1. Better memory planning

2. Optimized and model-specific INT8 code for Cortex-M

44

# Memory planning: a (rotated) game of Tetris



RAM size →

Time (layer execution) →

Objective of the game:
Use as little RAM as possible

Each Tetris block is a tensor

# Memory planning for an example model



Mem requirement: ~750KB

RAM size →

0 KB              640 KB

Time (layer execution) →

Conv2D

Conv2D

DWConv

Add

Add

Dense

46    *(simplified view: only activations are planned here, not weights)*

# A much better memory plan



(simplified view: only activations are planned here, not weights)

# Even better: lower granularity planning



*(simplified view: only activations are planned here, not weights)*

# Memory planning at Plumerai: summary

Lower RAM usage:

1. Smarter tensor placement

2. Lower granularity planning

RAM savings highly model dependent!

# Optimized INT8 code for speed



Optimized Conv2D code: im2col + GEMM

Optimized code for special cases, e.g. 1x1 Conv2D

Optimized code for other op (e.g. ADD):

Example code optimizations:

- Hand-written assembly (if needed)
- Specialization for Cortex-M4 or M7 capabilities
- Register-count aware optimisations
- Template-based loop unrolling
- Weight memory layout pre-processing

# Model-agnostic - vs - model-specific



Example: (micro) TVM

| ① | ② | ③ | ④ | ⑤ | ⑥ | ⑦ |
|---|---|---|---|---|---|---|
| TF/PyTorch/ONNX | Relay (High-level IR) | TE (Computation definition) | AutoTVM/AutoScheduler (Auto-tuning module) | TE + Schedule (Optimization specification) | TIR (Low-level IR) | Machine Code |

*Image taken from: https://tvm.apache.org/docs/tutorial/introduction.html*

# Model-specific code generation

```c
// Some function with unknown num-channels
void some_loop(int8_t* src, int8_t* dst,
               int num_channels) {
    for (int i = 0; i < num_channels; ++i) {
        dst[i] = src[i] * 16;
    }
}
```

```
1   some_loop(signed char*, signed char*, int):
2           cmp      r2, #0
3           ble      .L1
4           subs     r0, r0, #1
5           subs     r1, r1, #1
6           add      r2, r2, r0
7   .L3:
8           ldrb     r3, [r0, #1]!    @ zero_extendq
9           lsls     r3, r3, #4
10          cmp      r0, r2
11          strb     r3, [r1, #1]!
12          bne      .L3
13  .L1:
14          bx       lr
```

Generic code with compare and branch instructions

```c
// Same function but with 3 channels hard-coded
void some_loop(int8_t* src, int8_t* dst) {
    for (int i = 0; i < 3; ++i) {
        dst[i] = src[i] * 16;
    }
}
```

```
15  some_loop(signed char*, signed char*):
16          ldrb     r3, [r0]         @ zero_extendqisi2
17          lsls     r3, r3, #4
18          strb     r3, [r1]
19          ldrb     r3, [r0, #1]     @ zero_extendqisi2
20          lsls     r3, r3, #4
21          strb     r3, [r1, #1]
22          ldrb     r3, [r0, #2]     @ zero_extendq
23          lsls     r3, r3, #4
24          strb     r3, [r1, #2]
25          bx       lr
```
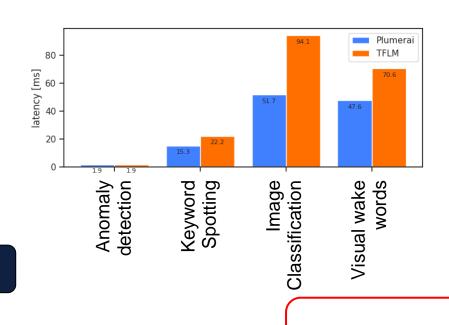
Unrolled code with only add, load and stores

# Better speed at Plumerai: summary

Lower latency, better speed:

1. Optimized code for Cortex-M

2. Model-specific code generation



Latency savings dependent on the layers and layer configurations

You've got mail!

# 5. Conclusion

# The world's fastest Cortex-M inference



Monday, October 4, 2021

## The world's fastest deep learning
## software for Arm Cortex-M

New: Try out our inference engine with your own model!

At Plumerai we enable our customers to perform increasingly complex AI tasks on tiny
We're proud to announce that our inference software for Arm Cortex-M microcontroller
memory-efficient in the world, for both Binarized Neural Networks and for 8-bit deep le
inference software is an essential component of our solution, since it directs resource
operating system. It h                                                      TensorF
Microcontrollers with                                                      cy. It als
deep learning inferen

| | rence ti | |
|---|---|---|
| TensorFlow L | 129 | |
| Edge Impulse | 120 ms | 153 KiB |
| MIT's TinyEn | 124 ms | 98 KiB |
| STMicroelect | 103 ms | 109 KiB |
| **Plumerai's i** | **77 ms** | **80 KiB** |

Model data → Optimized code → Compiled binary → Code running on device

# What can Plumerai mean for you?



Monday, October 4, 2021

The world's fastest deep learning inference software for Arm Cortex-M

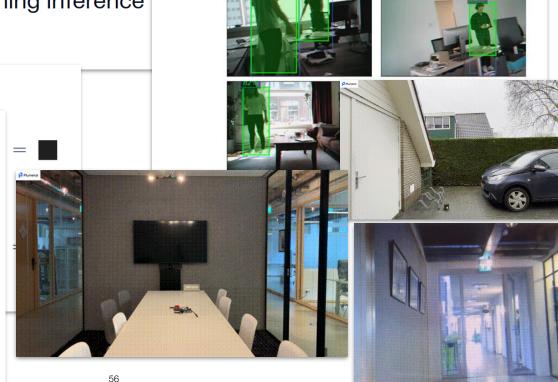New: Try out our inference engine with your own model!

Binarized Convolution

Person detection

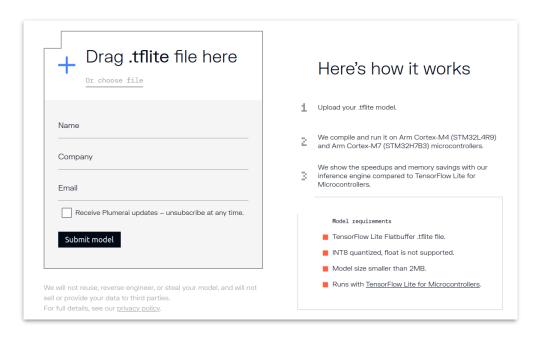- Plumerai Data Pipeline
- Plumerai BNN Models
- Plumerai Inference Stack
- Plumerai Hardware

# Public benchmarking service: try it yourself!



Visit plumerai.com/benchmark to try it with your own model

Contact hello@plumerai.com for help or other questions

# Copyright Notice

## www.tinyML.org