# What is tinyML and what is it used for?

## On-device machine learning applications in the single mW and below



### Vibration and motion

**Any 'signal'**

Predictive maintenance, sensor fusion, accelerometer, pressure, lidar/radar, speed, shock, vibration, pollution, density, viscosity, etc.



### Voice and sound

**Recognition and creation**

Keyword spotting, speech recognition, natural language processing, speech synthesis, sound recognition, etc.
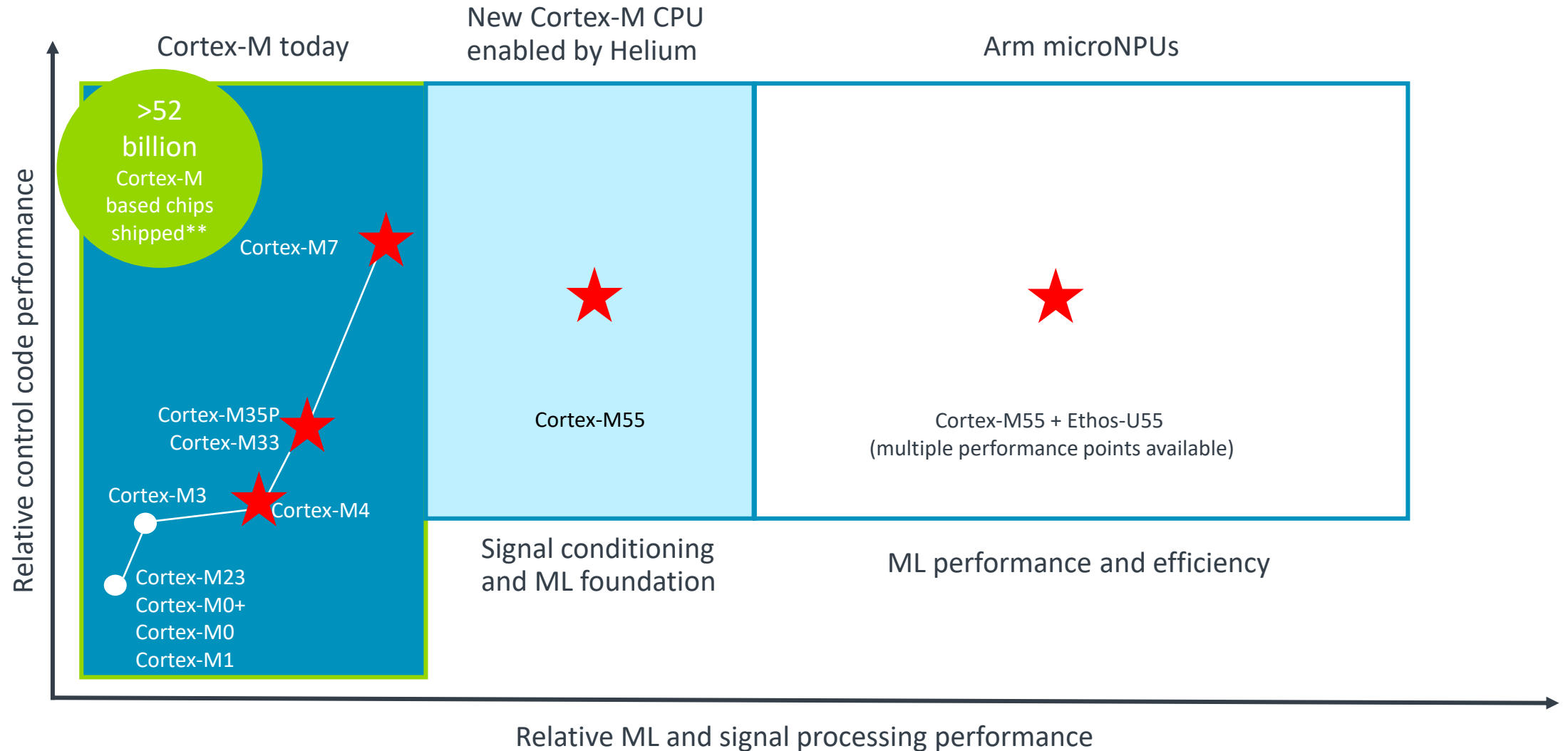


### Vision

**Images and video**

Object detection, face unlock, object classification etc.

arm

# Pushing the Boundaries for Real-time On-device Processing



Cortex-M today

New Cortex-M CPU
enabled by Helium

Arm microNPUs

>52 billion Cortex-M based chips shipped**

Cortex-M7

Cortex-M35P
Cortex-M33

Cortex-M3

Cortex-M4

Cortex-M23
Cortex-M0+
Cortex-M0
Cortex-M1

Cortex-M55

Cortex-M55 + Ethos-U55
(multiple performance points available)

Relative control code performance

Signal conditioning
and ML foundation

ML performance and efficiency

Relative ML and signal processing performance

**Based on Arm data

★ Well suited for ML & DSP applications

Graph not to scale
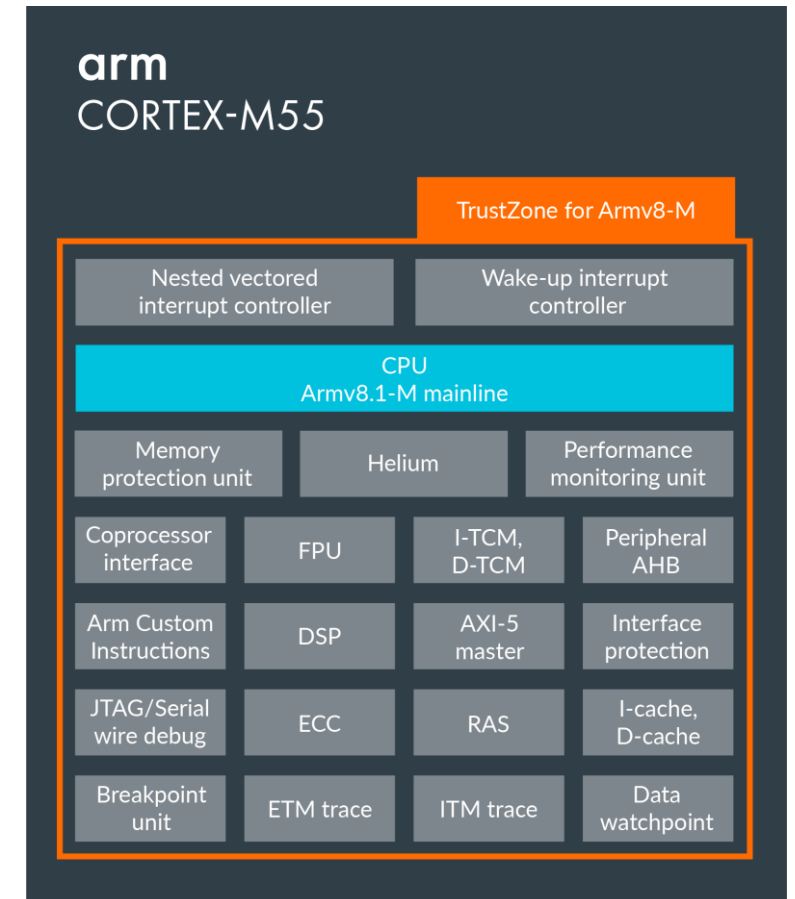
arm

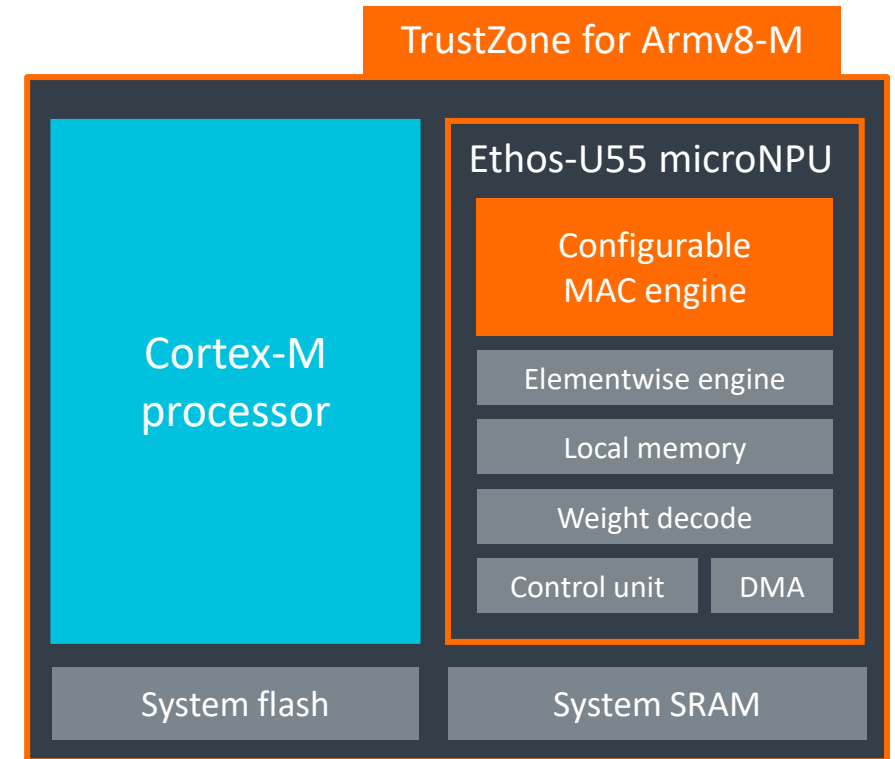# Cortex-M55: The Most AI-capable Cortex-M Processor

✓ First CPU based on Arm Helium technology

- Energy-efficient and configurable with vector processing capabilities

- Delivers up to 5x DSP performance and up to 15x ML performance*

- Versatile capability for both classical ML and NN inference

✓ Advanced memory interfaces for fast access to ML data and weights

✓ TrustZone support

✓ Extensive configurability



arm
CORTEX-M55

| TrustZone for Armv8-M | |
|---|---|
| Nested vectored interrupt controller | Wake-up interrupt controller |

CPU
Armv8.1-M mainline

| Memory protection unit | Helium | Performance monitoring unit |
|---|---|---|

| Coprocessor interface | FPU | I-TCM, D-TCM | Peripheral AHB |
|---|---|---|---|
| Arm Custom Instructions | DSP | AXI-5 master | Interface protection |
| JTAG/Serial wire debug | ECC | RAS | I-cache, D-cache |
| Breakpoint unit | ETM trace | ITM trace | Data watchpoint |

*Compared to previous Cortex-M generations

arm

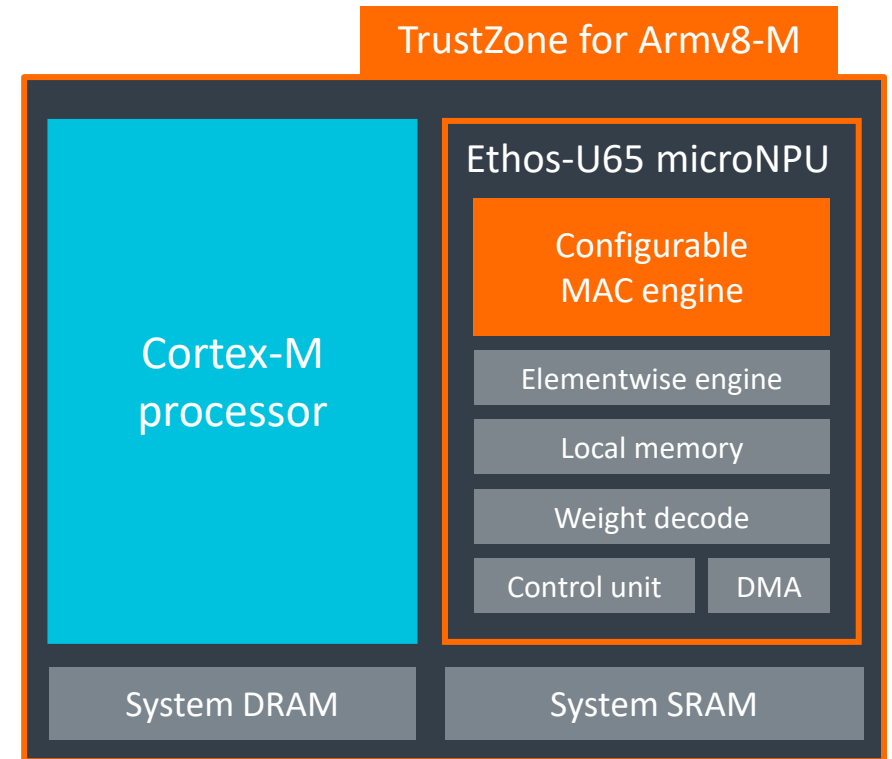# Ethos-U55: The first microNPU for Cortex-M

✓ Highest efficiency and small memory footprint

✓ 32, 64, 128, or 256 unit multiply-accumulate (MAC) engine

✓ Weight decoder and DMA for on-the-fly weight decompression

✓ Tooling available for offline optimization

✓ Works with a range of Cortex-M processors:

- Cortex-M55
- Cortex-M7
- Cortex-M33
- Cortex-M4

**TrustZone for Armv8-M**

Cortex-M processor

**Ethos-U55 microNPU**

Configurable MAC engine

Elementwise engine

Local memory

Weight decode

Control unit | DMA

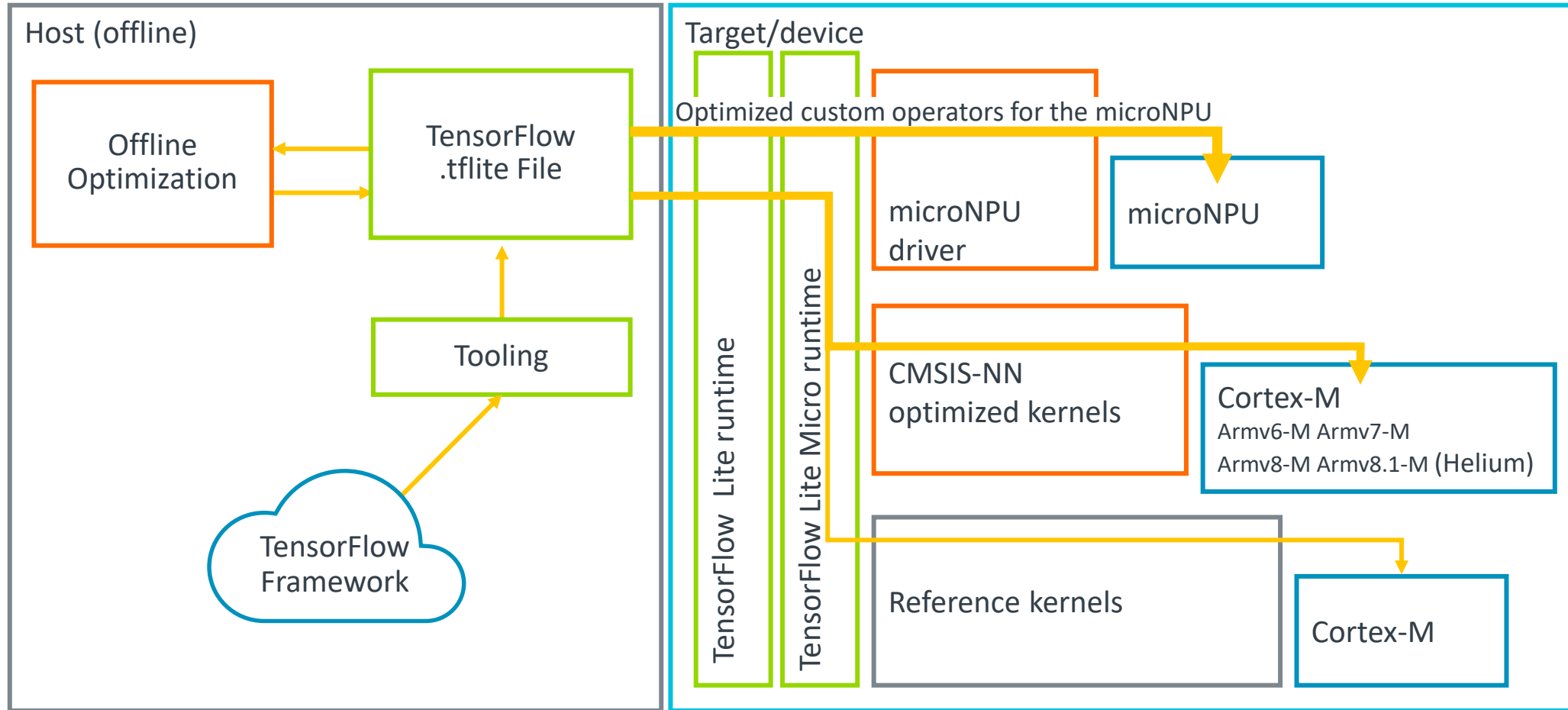System flash

System SRAM

arm

# Ethos-U65: The second generation microNPU

- ✓ 256 or 512 unit multiply-accumulate (MAC) engine

- ✓ DMA update for DRAM as well as flash support

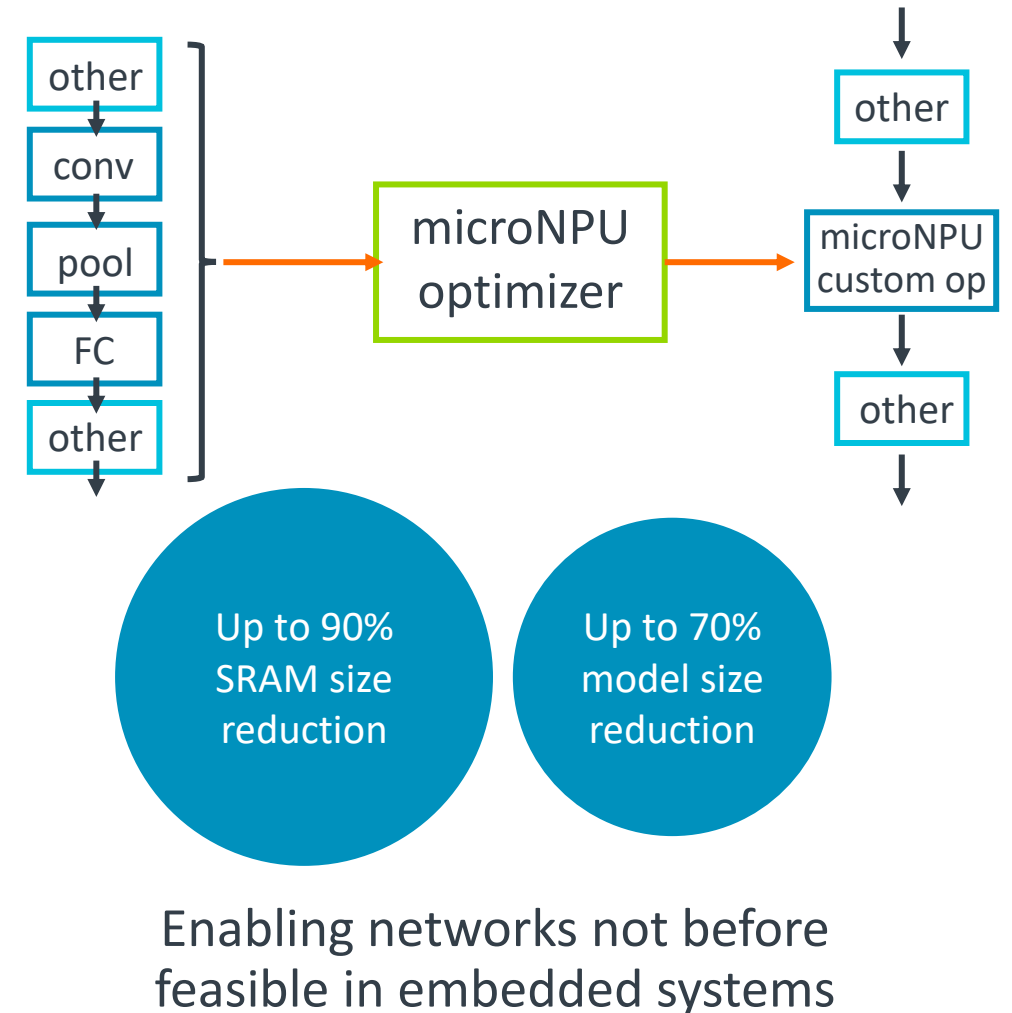- ✓ Can be an M-class subsystem inside an A-class system

**TrustZone for Armv8-M**

Ethos-U65 microNPU

Cortex-M processor

Configurable MAC engine

Elementwise engine

Local memory

Weight decode

Control unit | DMA

System DRAM | System SRAM

arm

# Mapping of NNs to Ethos-U using TensorFlow Lite
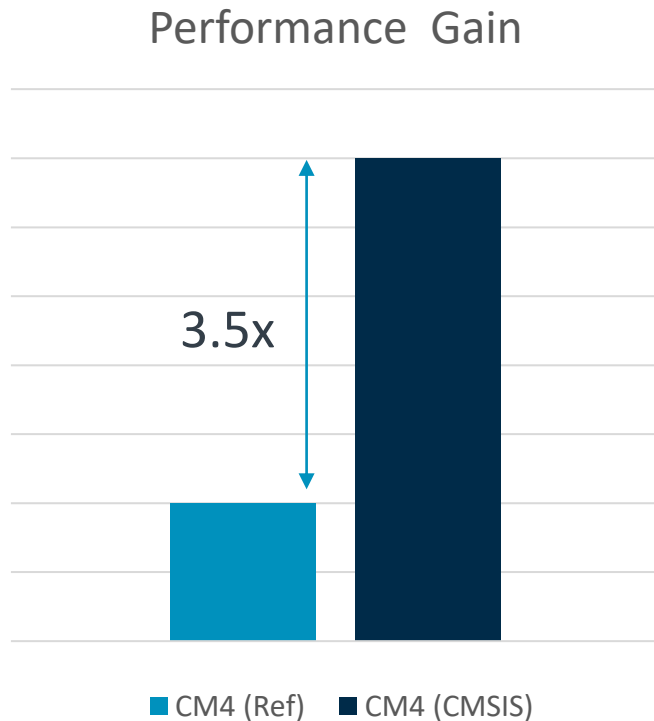
arm

# The Vela Optimizer

- Open source compiler
- Reads a tflite file and identifies subgraphs
- Optimizes scheduling of subgraphs
- Loss-less compression of weights
- Generates commands for microNPU
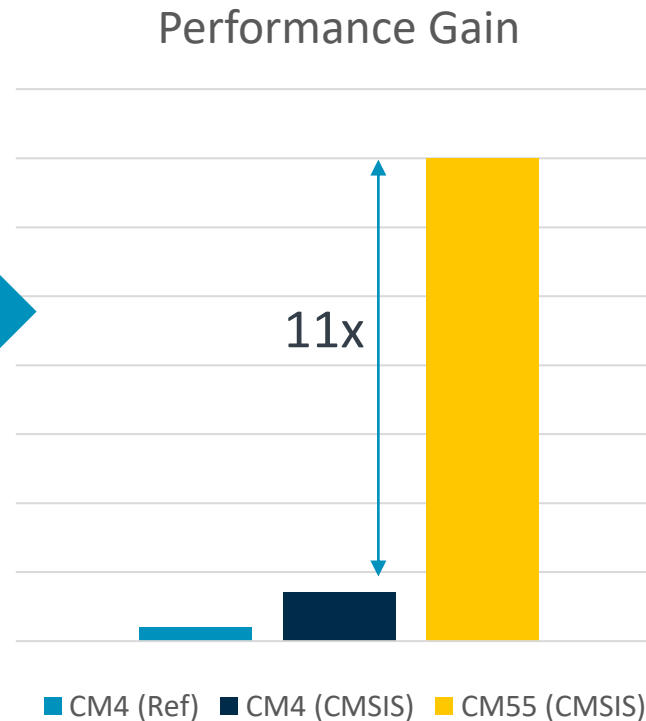- Writes out a modified tflite file

other

conv

pool

FC

other

microNPU optimizer

other

microNPU custom op

other

Up to 90% SRAM size reduction

Up to 70% model size reduction

Enabling networks not before feasible in embedded systems

arm

# Neural Network performance Across ARM IPs

Wav2Letter

## Efficient Software (CMSIS-NN)

Performance Gain

3.5x

■ CM4 (Ref)   ■ CM4 (CMSIS)

## AI Capable Cortex-M55

Performance Gain

11x

■ CM4 (Ref)   ■ CM4 (CMSIS)   ■ CM55 (CMSIS)

## AI Dedicated U55 256 MAC/cycle

Performance Gain

Total gain of almost 1000x

25x

■ CM4 (Ref)   ■ CM4 (CMSIS)
■ CM55 (CMSIS)   ■ CM55+U55

arm

# Full example: Typical ML Workload for a Voice Assistant

**Speed to inference**



50x

6x

Cortex-M7   Cortex-M55   Cortex-M55 + Ethos-U55

**Energy efficiency**

25x

7x

Cortex-M7   Cortex-M55   Cortex-M55 + Ethos-U55

✓ Faster responses

✓ Smaller form-factors

✓ Improved accuracy

Latency and energy spent for all tasks listed combined: voice activity detection, noise cancellation, two-mic beamforming, echo cancellation, equalizing, mixing, keyword spotting, OPUS decode, and automatic speech recognition.

arm

# Broadest Range of ML-optimized Processing Solutions



Cortex-A, Mali and
Ethos-N or Ethos-U65

Cortex-M and
Ethos-U55

Cortex-M55

Cortex-M today

| Vibration detection | Sensor fusion | Keyword detection | Anomaly detection | Object detection | Gesture detection | Biometric awareness | Speech recognition | Object classification | Real-time recognition |

**arm**