



Machine Learning Hardware

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.

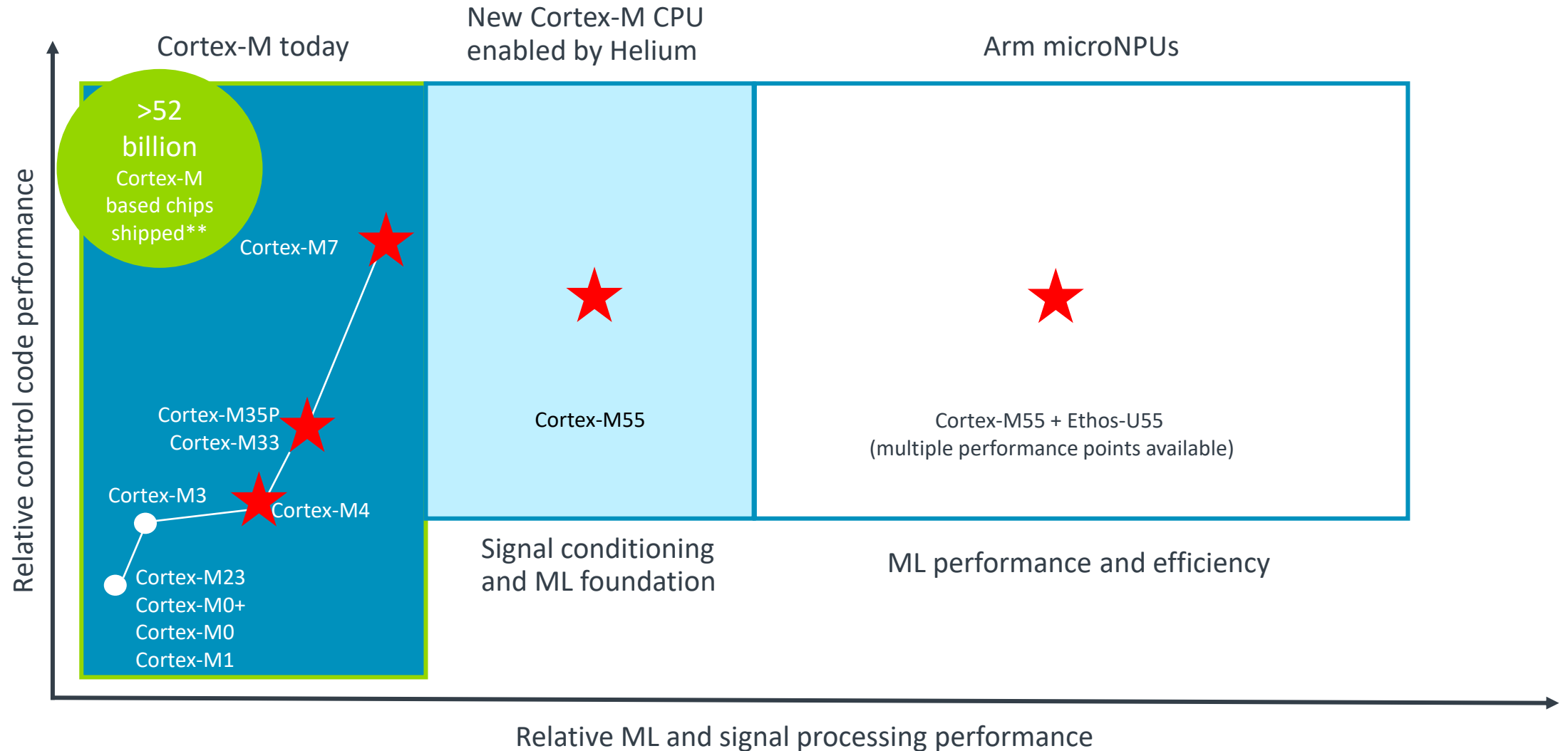


Vision

Images and video

Object detection, face unlock, object classification etc.

Pushing the Boundaries for Real-time On-device Processing

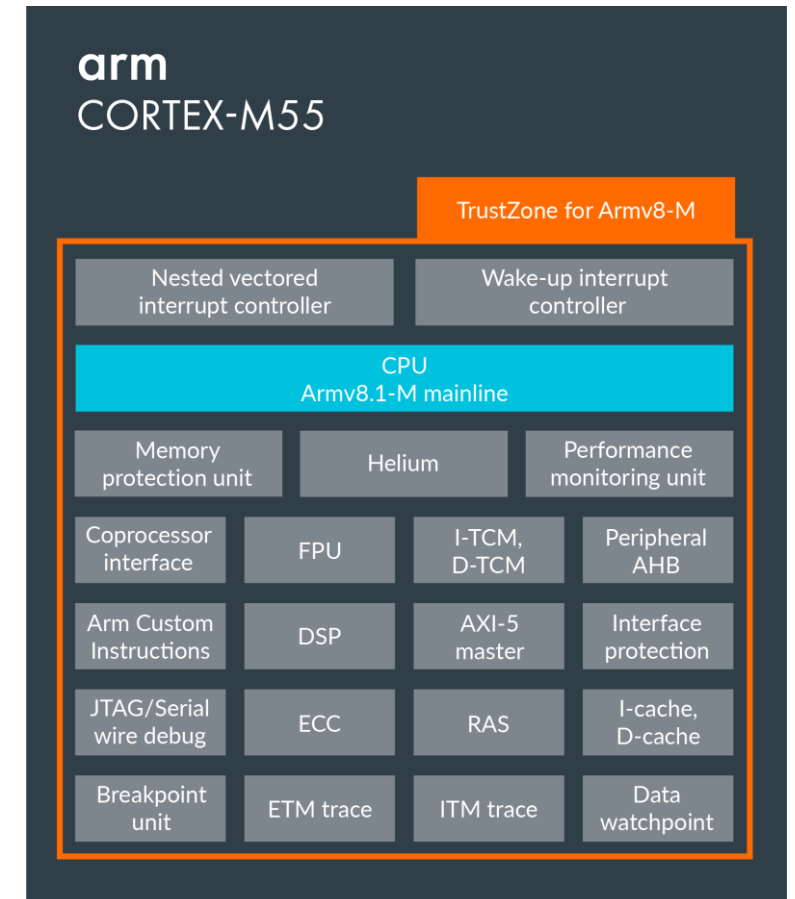


**Based on Arm data



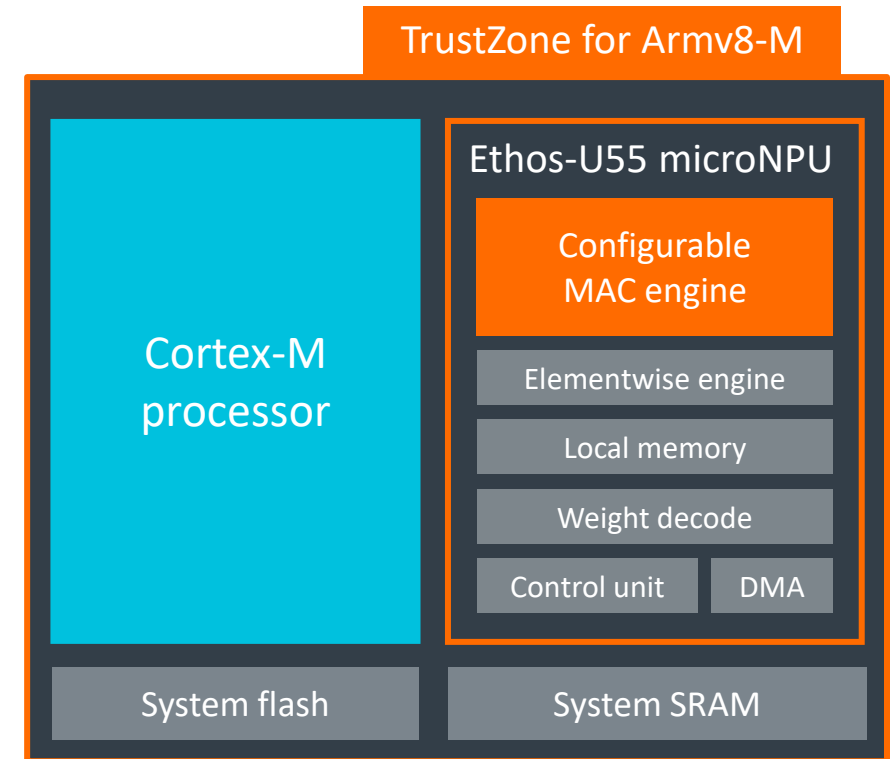
Cortex-M55: The Most AI-capable Cortex-M Processor

- ✓ 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



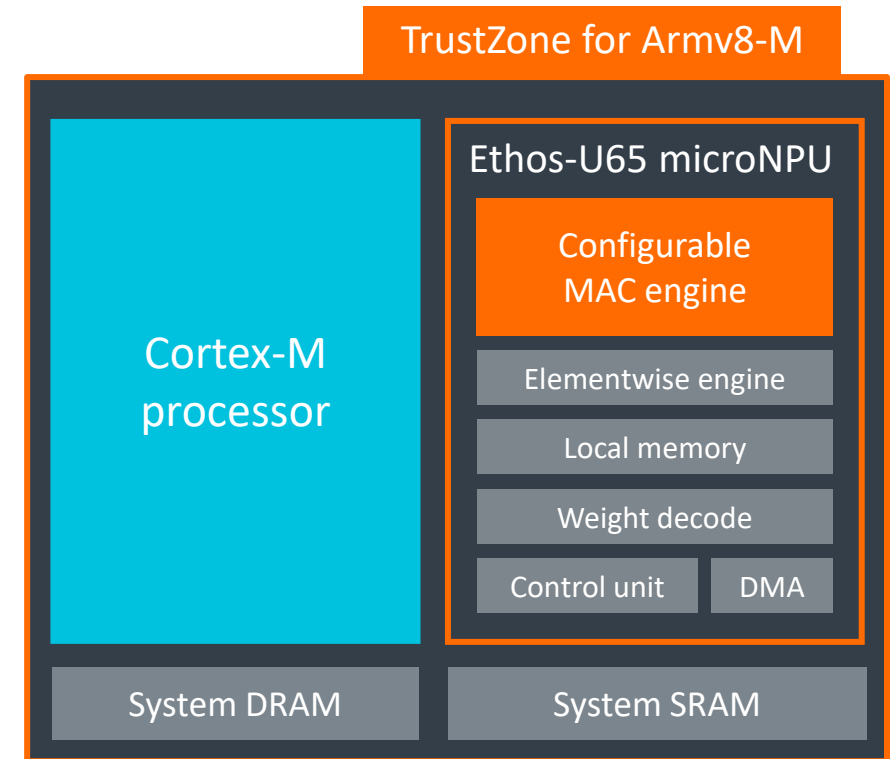
Ethos-U55: The first microNPU for Cortex-M

- ✓ 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

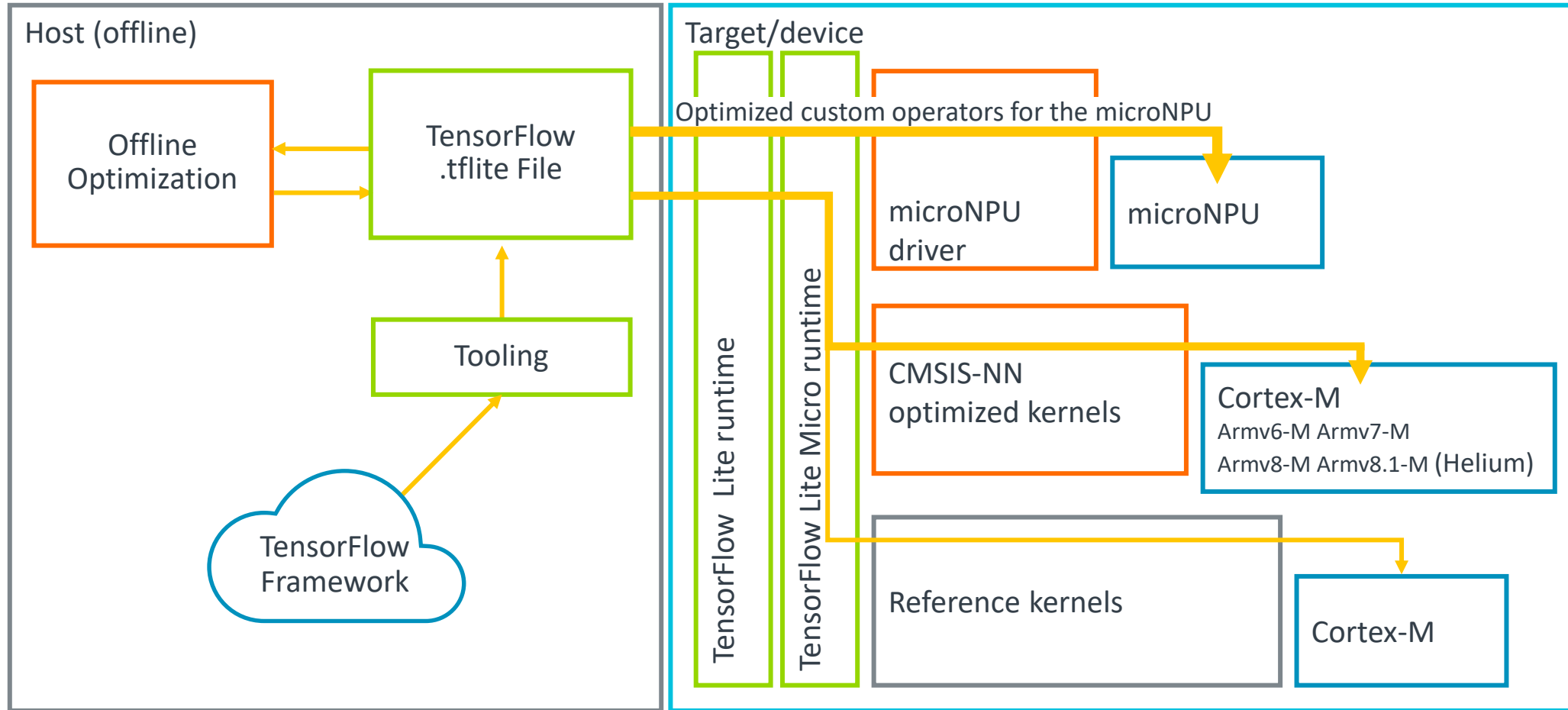


Ethos-U65: The second generation microNPU

- ✓ 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

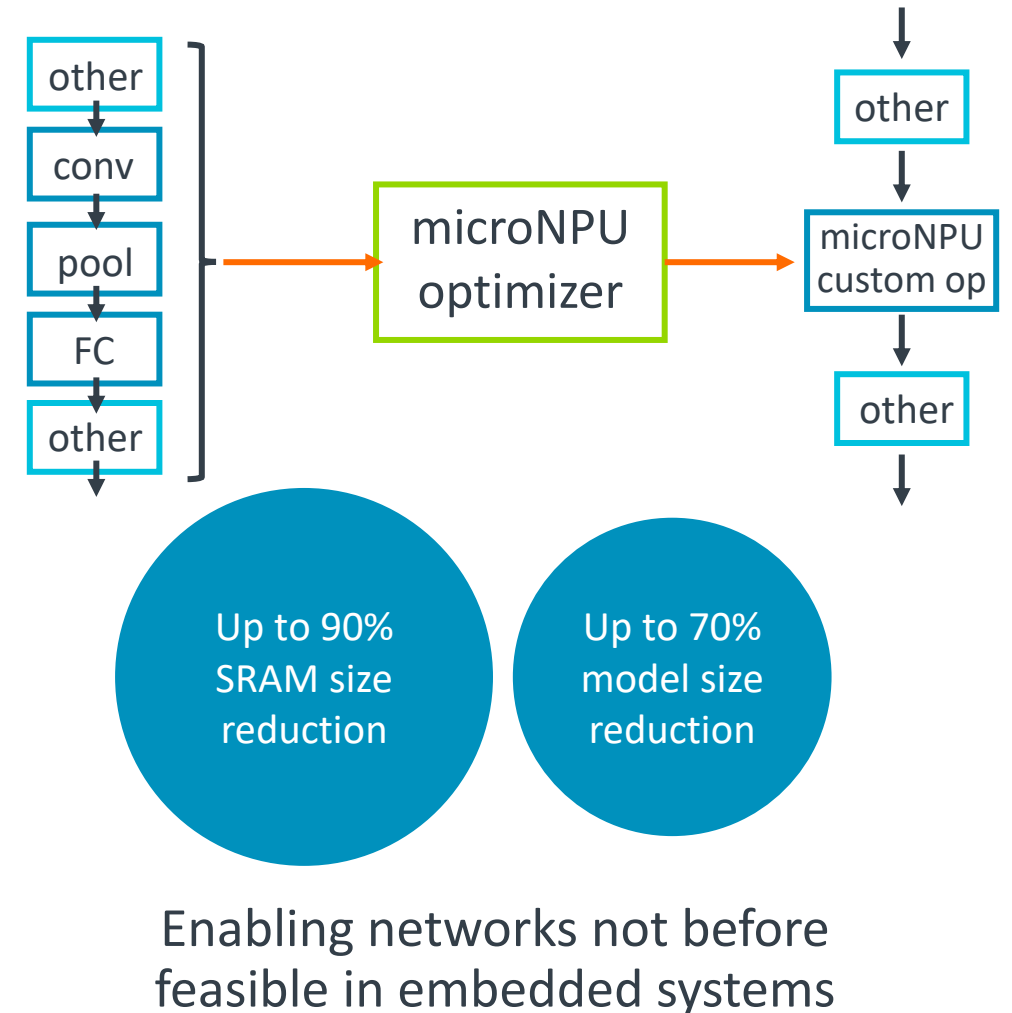


Mapping of NNs to Ethos-U using TensorFlow Lite



The Vela Optimizer

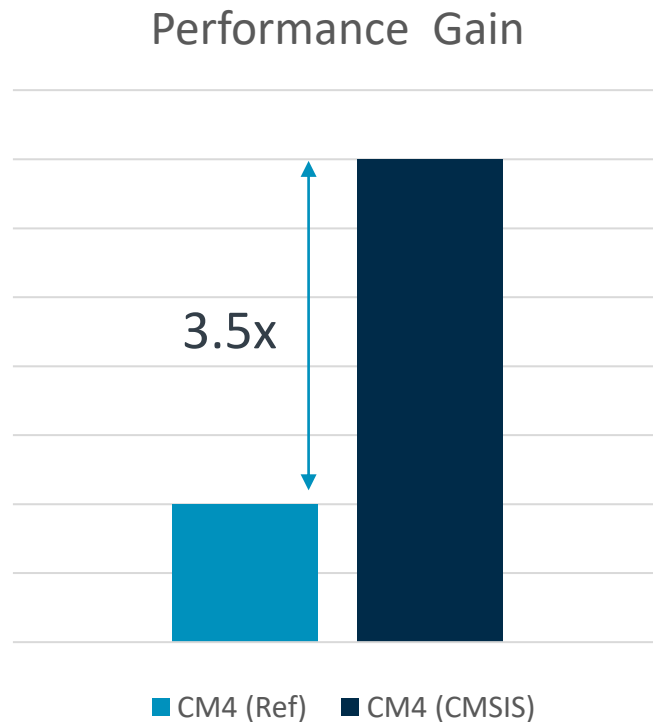
- Open source compiler
- Reads a tflite file and identifies subgraphs
- Optimizes scheduling of subgraphs
- Loss-less compression of weights
- Generates commands for microNPU
- Writes out a modified tflite file



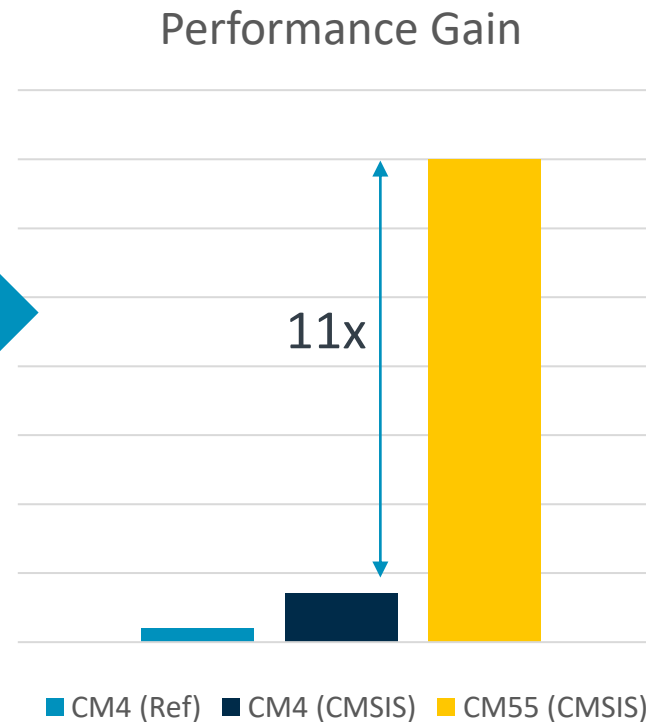
Neural Network performance Across ARM IPs

Wav2Letter

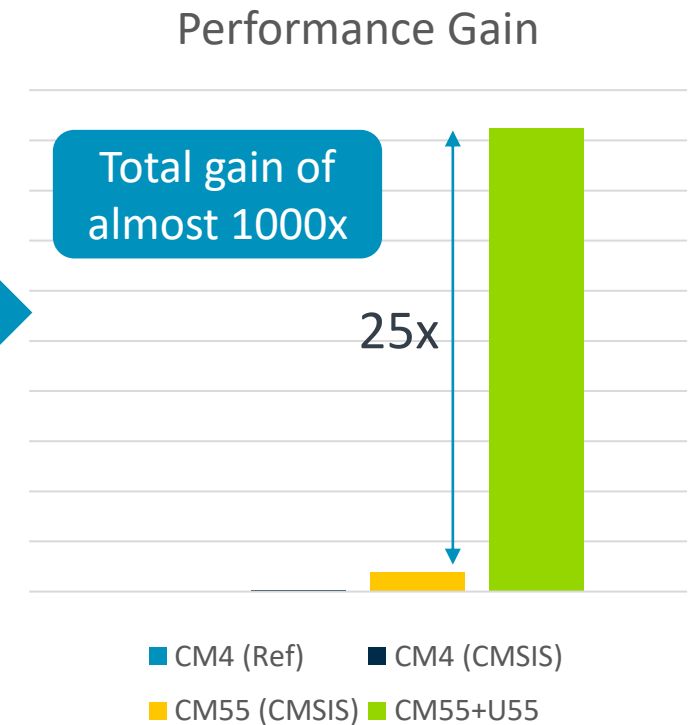
Efficient Software (CMSIS-NN)



AI Capable Cortex-M55

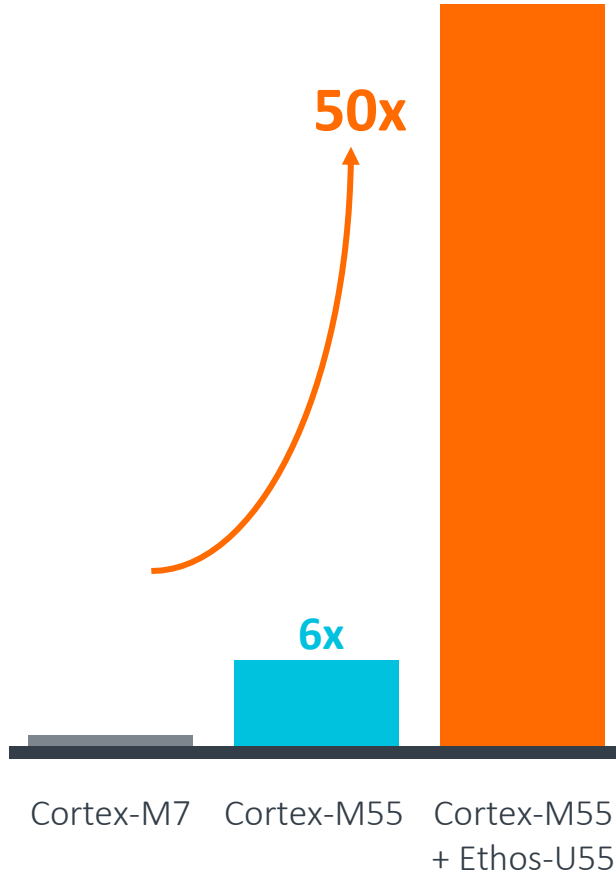


AI Dedicated U55 256 MAC/cycle

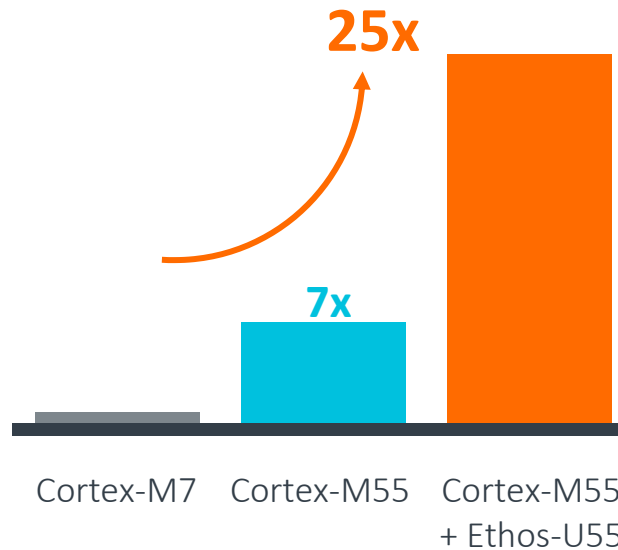


Full example: Typical ML Workload for a Voice Assistant

Speed to inference



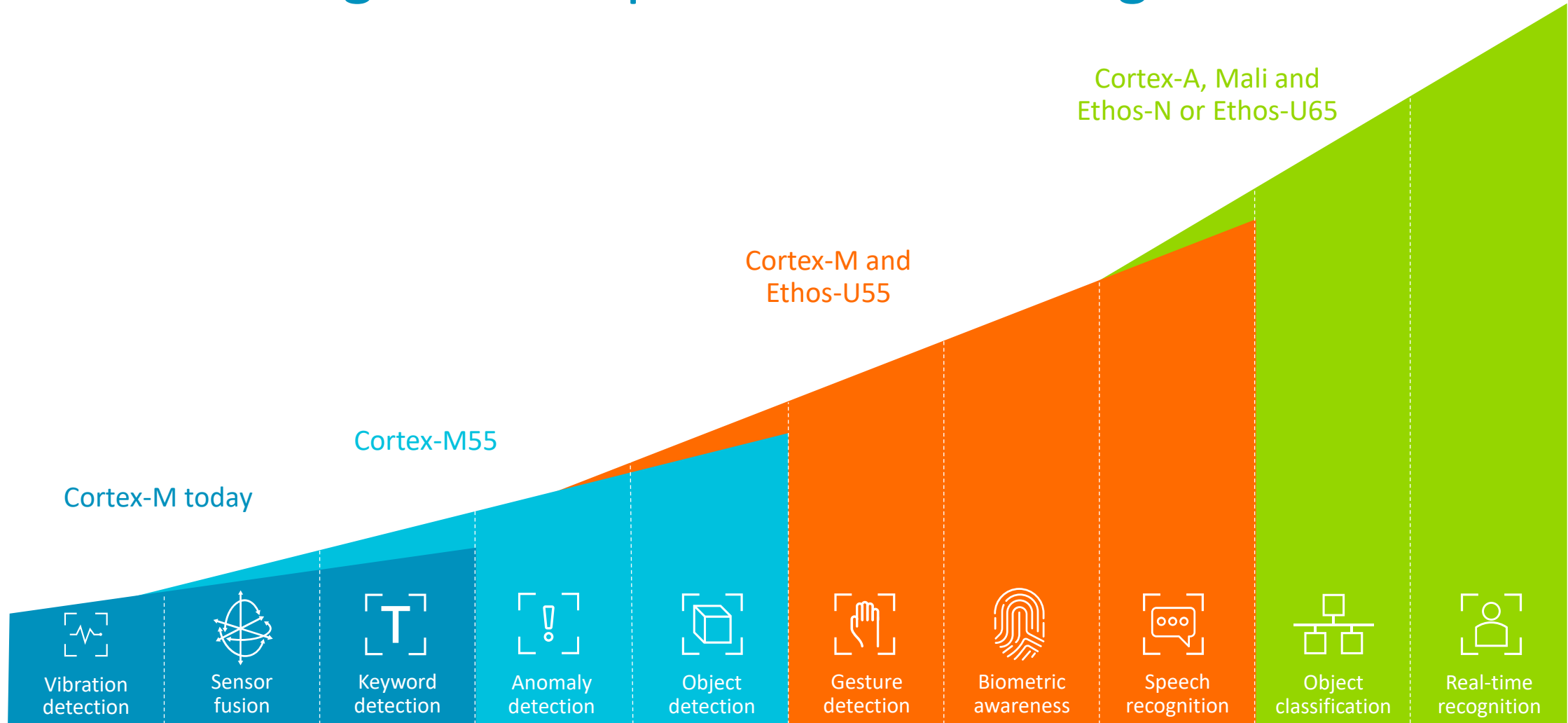
Energy efficiency



- ✓ 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.

Broadest Range of ML-optimized Processing Solutions





The Arm trademarks featured in this presentation are registered trademarks or trademarks of Arm Limited (or its subsidiaries) in the US and/or elsewhere. All rights reserved. All other marks featured may be trademarks of their respective owners.

www.arm.com/company/policies/trademarks