# Transformer Background

- What is a transformer? Ref. [1] Vaswani et al. Attention is all you need, NIPS 2017

- A highly scalable network architecture based on self-attention

# Why Transformers?

- Potentially unified architecture for text, audio, and image

- Models based on transformers perform outstandingly in natural language processing (NLP) and computer vision (CV)

- Support wide use cases, not only image classification but also applications such as super resolution, segmentation, object detection, and much more

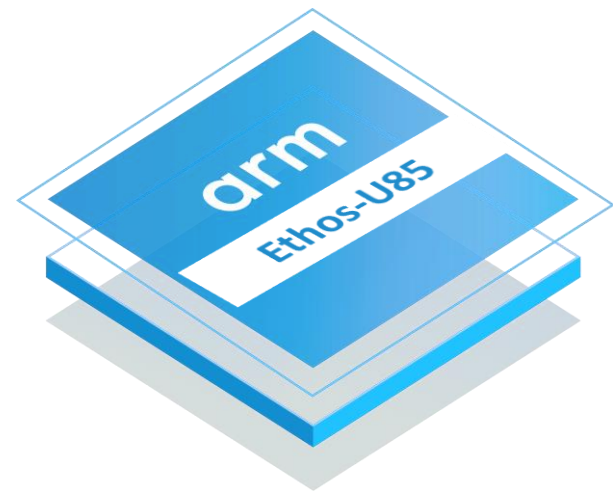# Transformer in Vision Applications

- While CNNs have inductive biases, e.g., locality and translation equivariance,

- The transformer uses self-attention to capture the dependencies within the input sequences

- Hence, models based on transformers are more extendable; i.e., work well in video understanding, image completion, multi-camera, and multi-modal domains

# Challenges in Deploying Transformer Models at the Edge

- Hardware is fragmented, ranging from CPU only, (CPU + GPU), or (CPU + accelerator), and others

  - What is the most suitable hardware solution for transformers?

- Efficiency is another challenge

  - How do you run transformer models with high power efficiency and low latency?

- Model size and memory usage

  - We need a toolset (with tutorials) to compress model size to a reasonable size so that it can be deployed at the edge.

# Arm Machine Learning Solution Supporting Vision Transformers

# Introducing Next Generation Arm NPU—
# What Makes it Attractive?

**Higher power efficiency**

- Targeting **20%** over current generation

**Increased performance**

- Configurations from **128** MACs/cycle to **2048** MACs/cycle

**Extended operator support**

- **Hardware accelerated transformer network support**

**Double MAC throughput**

- For **2/4** sparse layers

arm

Ethos-U85

# New Hardware Operators Accelerate Transformer Networks

- In addition to the operators currently supported by the original Ethos product family, the latest Arm Ethos-U85 includes native hardware support for transformer networks and DeeplabV3 semantic segmentation network, such as:
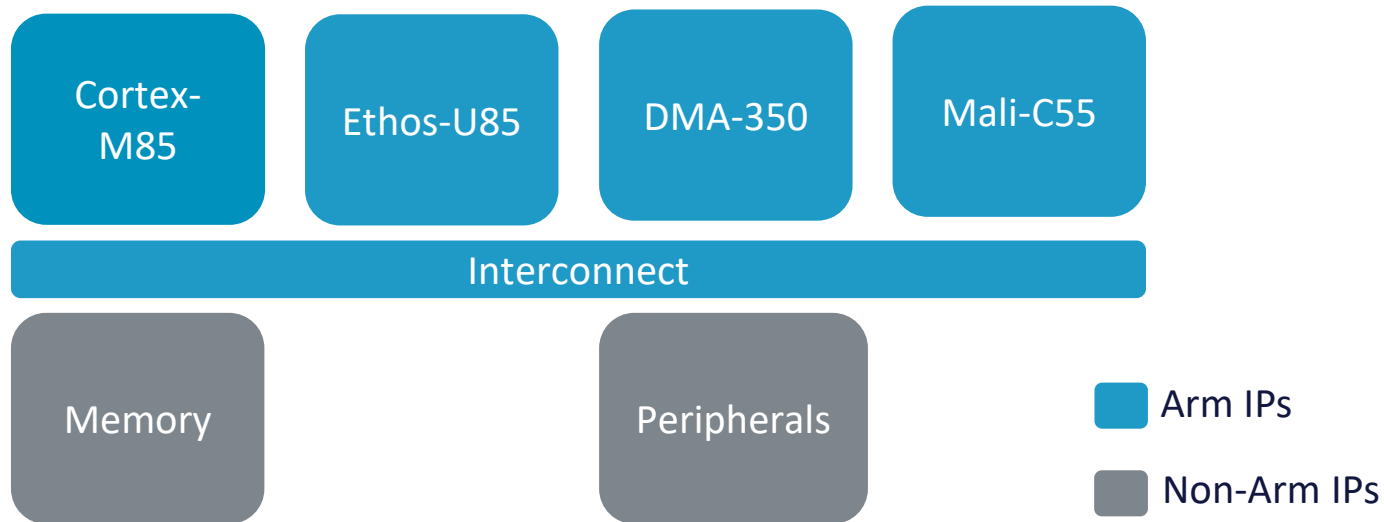
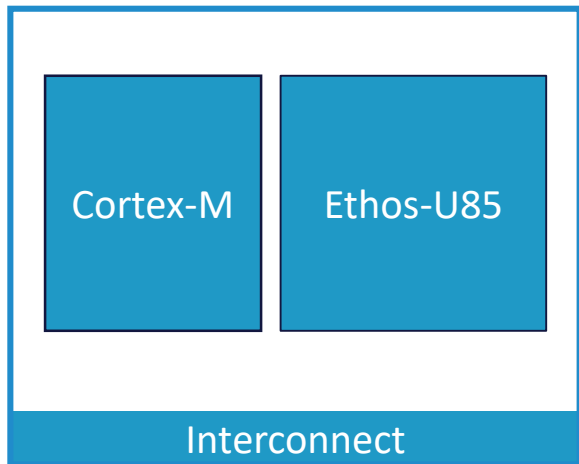| TRANSPOSE | GATHER | MATMUL | RESIZE BILINEAR | ARGMAX |

# Arm Example Subsystem

- Pre-integrated and verified machine learning solution



© 2024 Arm Inc.

# How to Use Ethos-U85 in a System

- End Point AI: Cortex-M based system

- ML Island: Cortex-A based system

- Discrete NPU: Cortex-A only

**End Point AI: Cortex-M based system**
- Cortex-M
- Ethos-U85
- Interconnect
- System SRAM
- System Flash

**ML Island: Cortex-A based system**
- Cortex-A
- Cortex-M
- Ethos-U85
- Interconnect
- System SRAM
- System Flash
- DRAM

**Discrete NPU: Cortex-A only**
- Cortex-A
- Ethos-U85
- Interconnect
- System SRAM
- System Flash
- DRAM

Arm IPs       Non-Arm IPs

arm

# Software Flow on Arm Machine Learning Solution

- Cortex-M CPU with Ethos-U85



**HOST (OFFLINE)**

TF Frame-work → TF Quantization Tooling TFLite Converter → TFL flat file → NN Optimizer

**TARGET / DEVICE**

TFLu Runtime | Ethos-U85 Driver | **Ethos-U85 NPU**
CMSIS-NN Optimized Kernels | **Cortex-M CPU**
Ref. Kernels

Arm IPs

# Software Flow on Arm Machine Learning Solution

- Cortex-M + Cortex-A system

**HOST (OFFLINE)**

TF Frame-work → TF Quantization Tooling TFLite Converter → TFL flat file

TFL flat file ↔ NN Optimizer

**TARGET / DEVICE**

Wrapper App

.tflite flatbuffer

TFLiteμ runtime

Ethos-U85 Driver

Cortex-M + cache, MPU

Ethos-U85 NPU

Application

Inference API

Linux OS

Subsystem driver

Cortex-A + cache, MMU

Arm IPs

AXI bus

Address filter

AXI bus

SRAM

NPU carveout (boot time)

Linux/OS managed area

DRAM

© 2024 Arm Inc.

arm

12

# Software Flow on Arm Machine Learning Solution

- Cortex-A based system



HOST (OFFLINE)

TF Frame-work → TF Quantization Tooling TFLite Converter → TFL flat file ↔ NN Optimizer

TARGET / DEVICE

Application
TFLite delegate
Linux OS
NPU driver
Ethos-U85 NPU
Cortex-A + cache, MMU
AXI bus
Address filter
AXI bus
SRAM
NPU carveout (boot time)
Linux/OS managed area
DRAM

: Arm IPs

© 2024 Arm Inc.

13

# Arm Toolset Enables the Efficient Implementation of Transformers on Ethos

Arm transformer Tutorials, the Jupyter notebooks (.ipynb) showing how to quantize and compress transformer encoder and encoder-decoder models.

# Vision Transformer Example Implementation
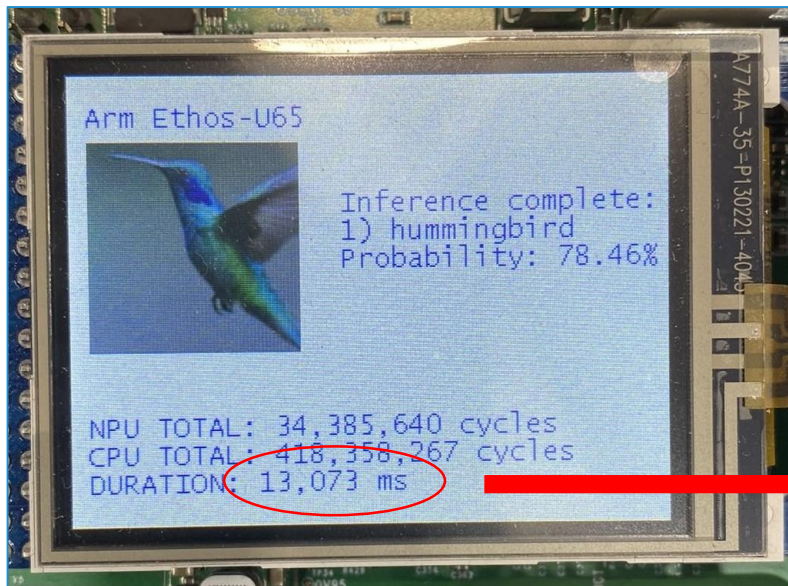
arm

# DEiT Tiny Runs on Ethos-U85

Demo is to compare how much faster the latest Ethos-U85 runs a transformer network compared to the previous Ethos, since there is no fall back for those operators with Ethos-U85
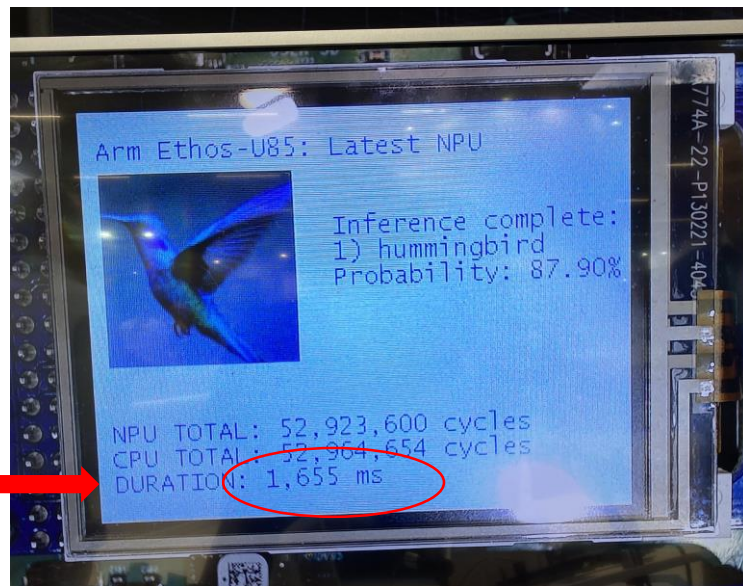
images → Arm MPS3 board with the latest Ethos-U85

images → Arm MPS3 board with previous Ethos-U

Output

Ethos-U85 - hummingbird - execution speed

Previous Ethos - hummingbird - execution speed

# Up to 8X Acceleration in Inference time

- Previous Ethos



- The Latest Ethos-U85



For more details, please visit Arm booth at **#409**.

# Summary

- Machine learning (ML) is everywhere, and its landscape is evolving from CNNs to transformer-based models

- Arm just launched the latest NPU in the Arm Ethos product family to extend the support of accelerating transformers at the edge

- Finally, "Edge AI runs on Arm."

# Resources

Arm Ethos-U product page
https://www.arm.com/products/silicon-ip-cpu?families=ethos%20npus

Arm transformer tutorials
https://github.com/ARM-software/ML-zoo/tree/master/tutorials/transformer_tutorials

Arm keyword-transformer
https://github.com/ARM-software/keyword-transformer

**Please visit Arm booth #409 at the 2024 Embedded Vision Summit for more demos:**

"The Newly Launched Arm Ethos-U85 NPU"

"Renesas RZ/V2H- Qual-core Cortex-A55 Vision AI MPU"

"Arm-Himax, the High-efficiency Embedded Computer Vision"

# Reference

- Reference [1]: A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in Proceedings of the 31st International Conference on Neural Information Processing Systems, 2017, pp. 6000–6010

# Thank You

arm