



arm

Comprehensive Arm Solutions for Innovative Machine Learning (ML) and Computer Vision (CV) Applications

Machine Learning is a Subset of Artificial Intelligence

AI means many things to many people

Artificial Intelligence

Machine Learning

Perception & Vision

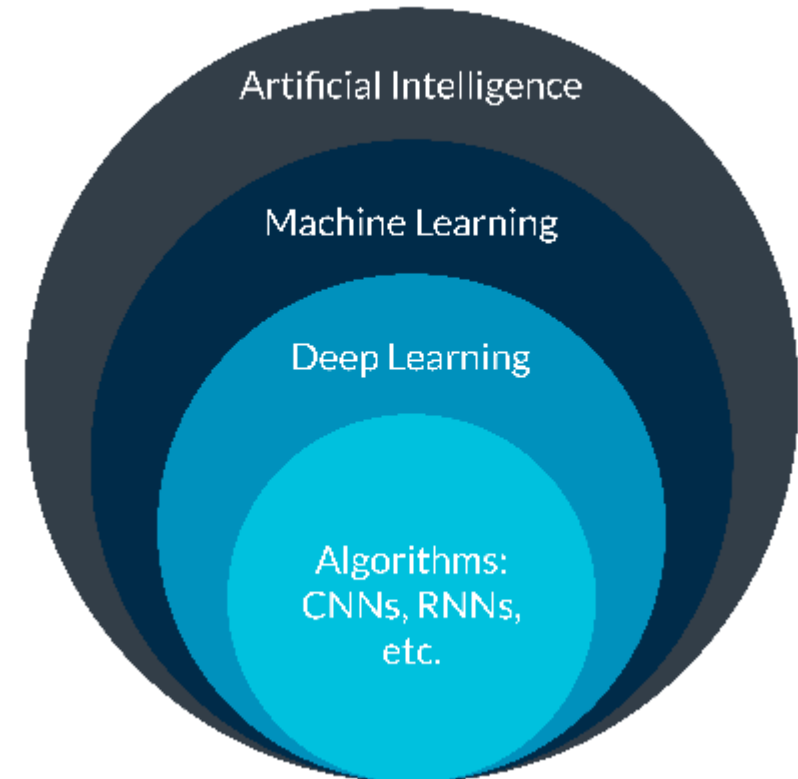
Natural Language Processing

Knowledge Representation

Planning & Navigation

Generalized Intelligence

ML itself has a lot of depth

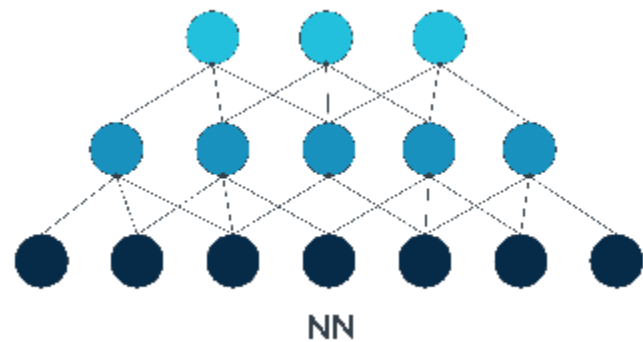


Why Artificial Intelligence(AI) Exploding Now

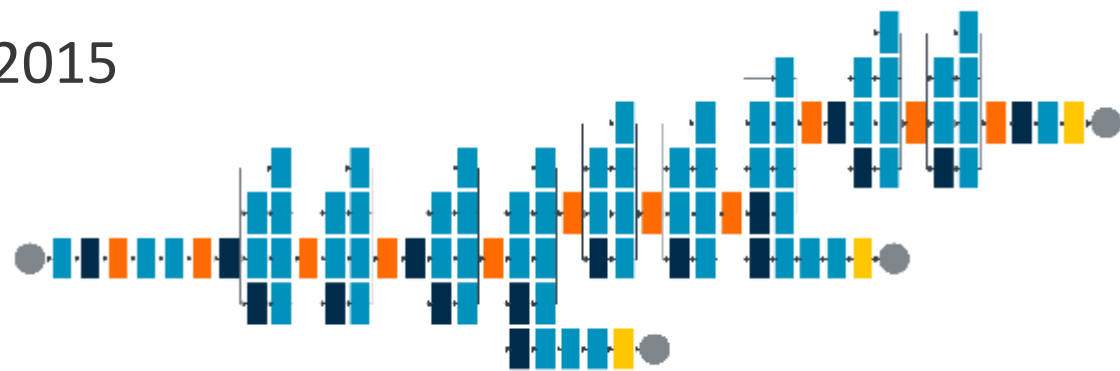
Availability of increased data sourced at the edge with ubiquitous powerful compute!

Compute

2010



2015



Data

2016 – 1 zettabyte



IP Traffic

2020 – 2.3 zettabyte



AI Presents Significant Opportunity for Innovation

VR/MR



Robotics



Drones



Shipping & logistics



IoT



Home, surveillance & analytics



Automotive



Mobile



Distributing Intelligence

Cloud-based training

High-performance processing



On-device learning

Security and privacy for your data



AI in your hand & cloud

Real-time inference for autonomous systems



Why is On-device ML Driving AI to the Edge?



Bandwidth



Power



Cost



Latency



Privacy

Arm ML Platform Enables



Efficiency

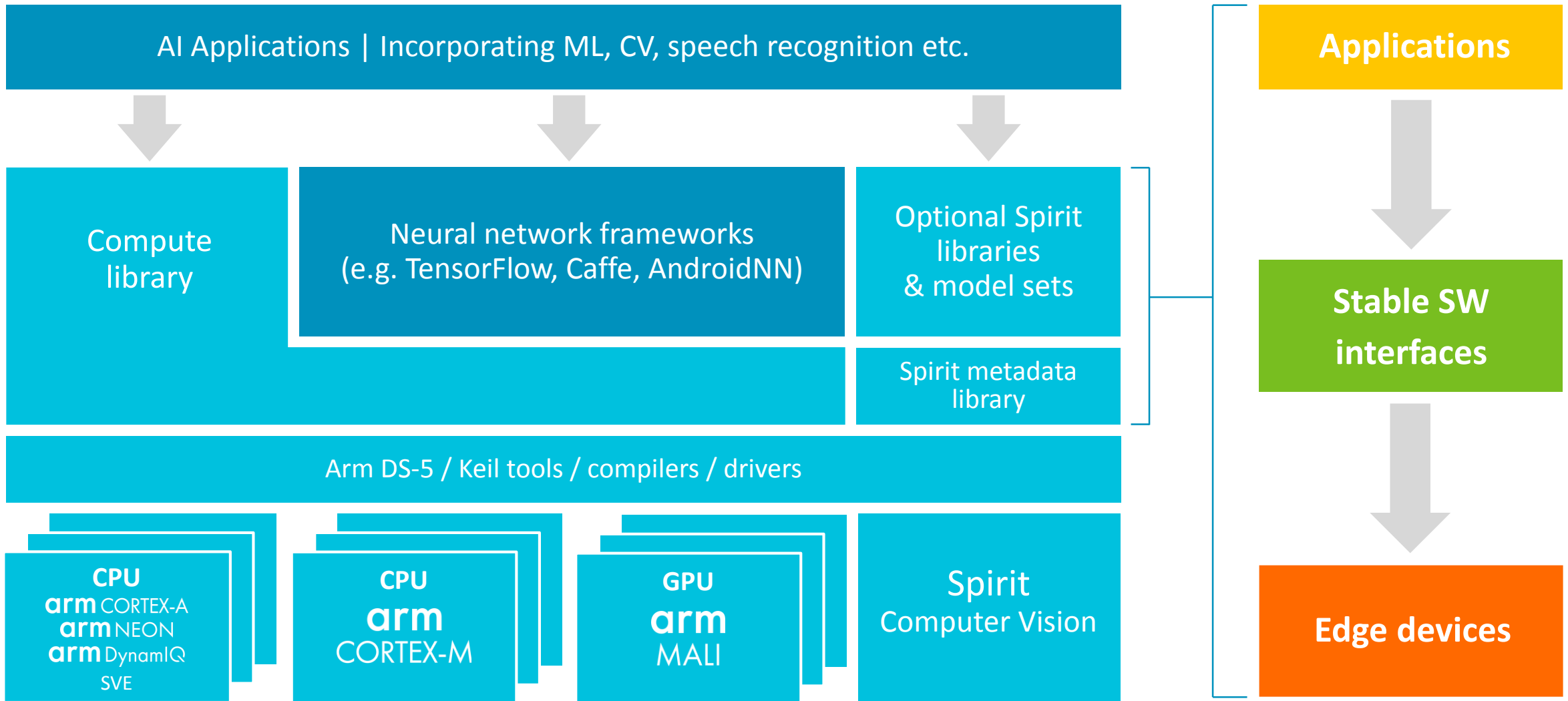


Flexibility



Freedom

Arm's ML Computing Platform



Components of Arm ML Platform

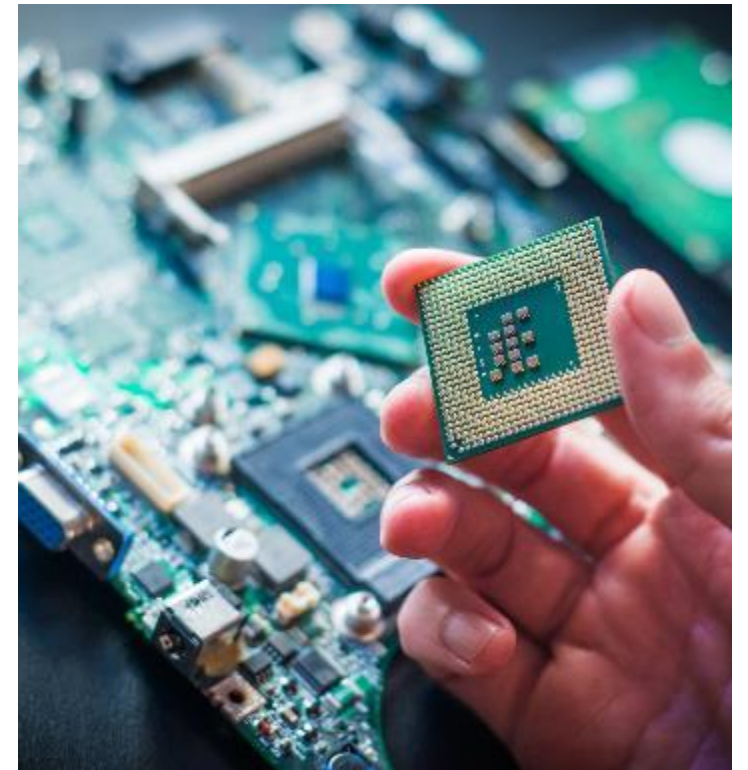
Software



Specialized Acceleration



Hardware



Software Development

```
elif _operation == "MIRROR_Y":
    mirror_mod.use_x = False
    mirror_mod.use_y = True
    mirror_mod.use_z = False
elif _operation == "MIRROR_Z":
    mirror_mod.use_x = False
    mirror_mod.use_y = False
    mirror_mod.use_z = True

#selection at the end -add back the deselected mirror modifier object
mirror_ob.select= 1
modifier_ob.select=1
bpy.context.scene.objects.active = modifier_ob
print("Selected" + str(modifier_ob)) # modifier ob is the active ob
#mirror_ob.select = 0
done = bpy.context.selected_objects[0]
bpy.data.objects[mirror_ob.name].select = 1
```

Arm Compute Library

Faster, advanced processing

What is the Arm Compute Library?

Functions for CV and deep-learning algorithms

Optimized for Arm CPU and GPU

OS and platform agnostic

No fee, MIT license

Delivers faster processing

4.6x faster than stock OpenCV on NEON

Offers OpenCV and Open VX compatibility

Use as a plug-in backend for your own runtime implementation

Available now: <https://developer.arm.com/technologies/compute-library>

Arm Compute Library Partners

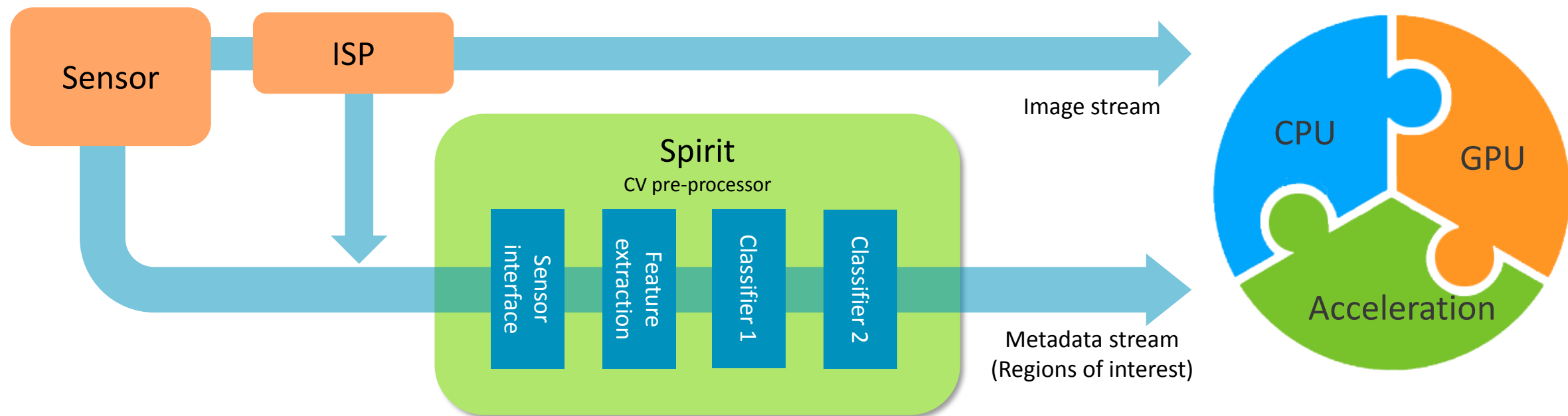




Specialized Acceleration

Computer Vision (CV)

Spirit for Object Detection and Localization



Comparison with Neural Network Framework Solutions



Spirit: Key Features

Object localization (>200k locations in full HD)

Size variation (~20x for full HD)

Scalable to 4K without performance compromise

Real time, 60 fps, no dropped frames

Invariance to optical distortions

Invariance to illumination conditions

High occlusion tolerance

Suitable for stationary and moving cameras



Comparison with a DSP

Spirit uses a form of HOG*/ SVM* baked into an efficient hardware design

Using a DSP to achieve the same performance

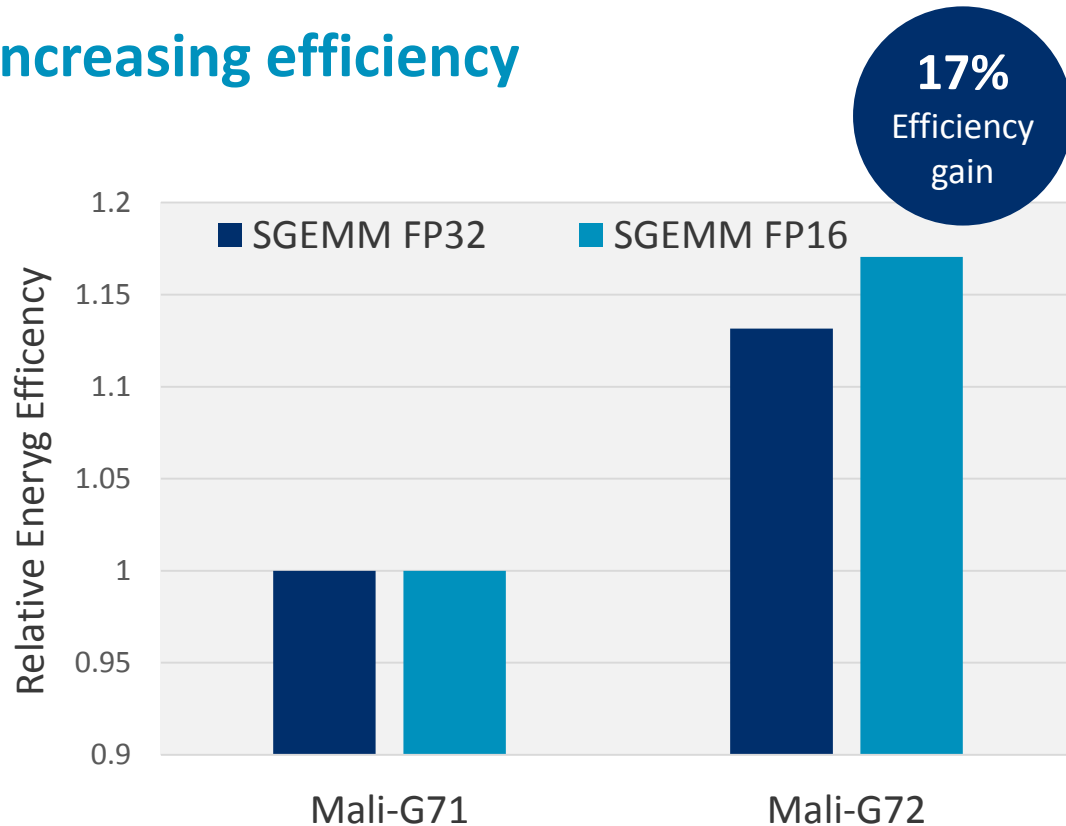
E.g. Pedestrian Detection on a DSP

- Processed at VGA resolution
- Achieving 5fps
- Operating at 40MHz/50MHz
- Scaling this to Spirit performance levels of 1080p60 would require the DSP to run at 3.24GHz

ML on Mali GPUs

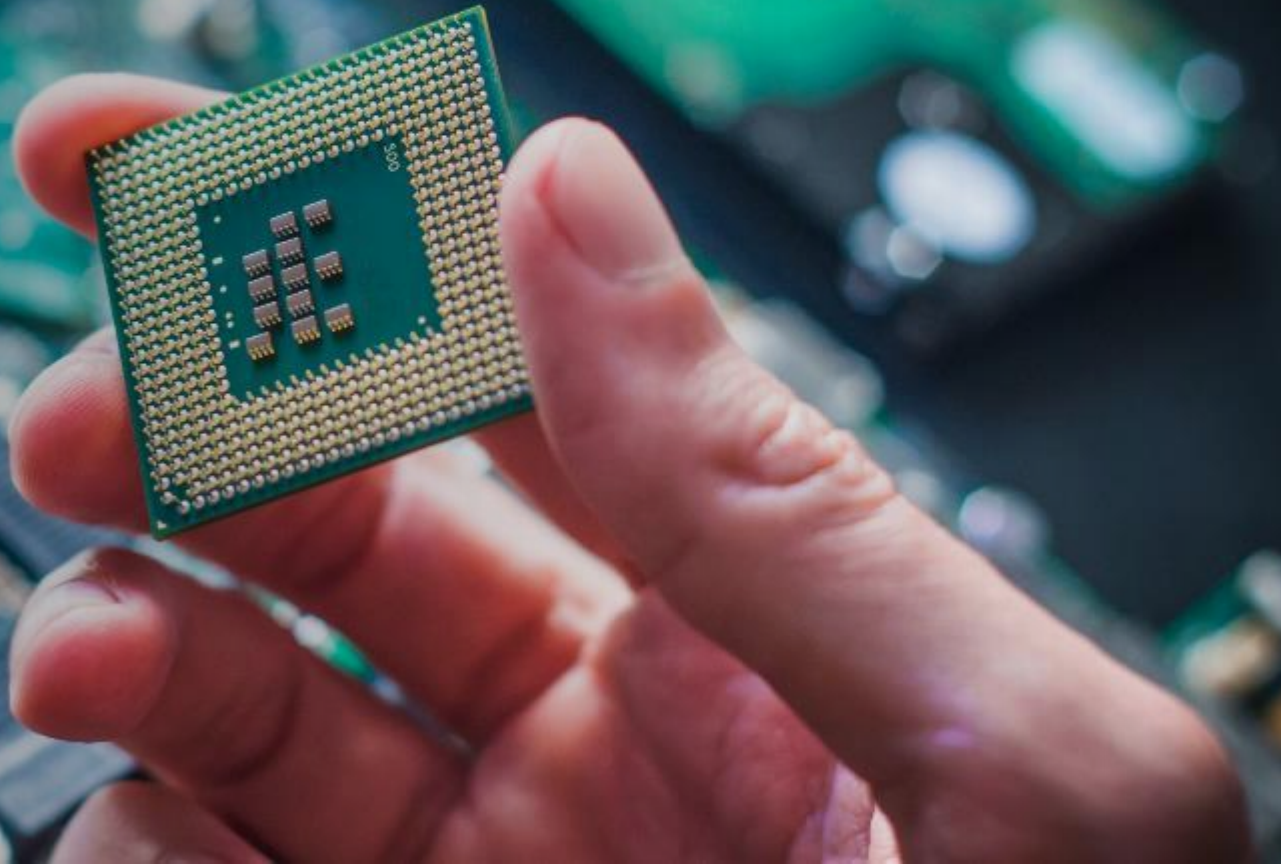
Mali GPUs: Increasing ML Throughput and Efficiency

Increasing efficiency



- GEMM depicts core functionality of ML algorithms
- Mali-G72 has several optimizations to improve ML inference
 - Less power-hungry FMA unit
 - Bigger L1 cache in the execution engine
- Mali-G72 is the most efficient Mali GPU for machine learning

Hardware



AI Applications at the Edge on Arm



Detect plant diseases



Sort cucumbers



Detect Caltrain delays

Instruction Sets for AI

Cortex-A

- Additional dot product instructions (Cortex-A55 and Cortex-A75)
- New Scalable Vector Extensions (SVE) instructions

Cortex-M

- Optimized CMSIS-DSP libraries for matrix multiplication

Closely-coupled acceleration

- Improved performance and efficiency (for broader use cases)
- Flexibility in multi-core computing with Arm DynamIQ technology

DynamiQ: New Cluster Design for New Cores

Arm DynamiQ big.LITTLE systems:

- Greater product differentiation and scalability
- Improved energy efficiency and performance
- SW compatibility with Energy Aware Scheduling (EAS)

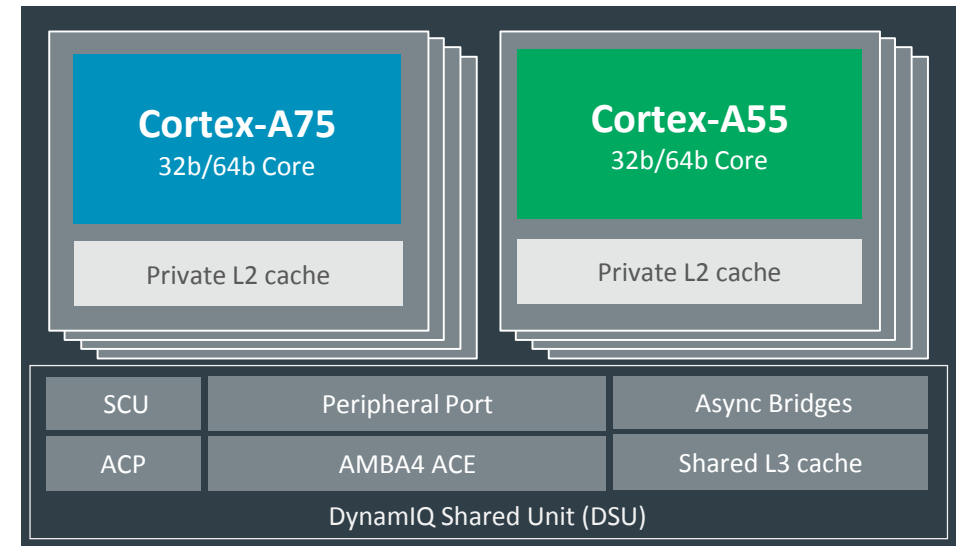
Private L2 and shared L3 caches

- Local cache close to processors
- L3 cache shared between all cores

DynamiQ Shared Unit (DSU)

- Contains L3, Snoop Control Unit (SCU) and all cluster interfaces

Additional instructions for ML



1b+7L



2b+6L



4b+4L



1b+2L



1b+3L



1b+4L

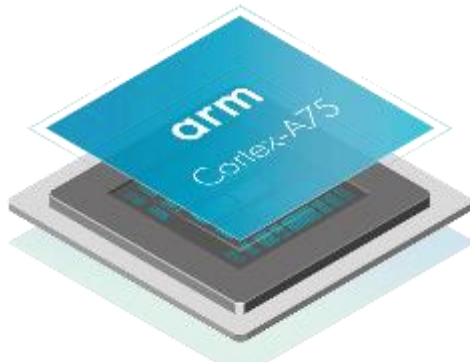
Example: DynamiQ big.LITTLE configurations

New DynamIQ-based CPUs for New Possibilities

Cortex-A75 processor

>50%

more performance
compared to current devices

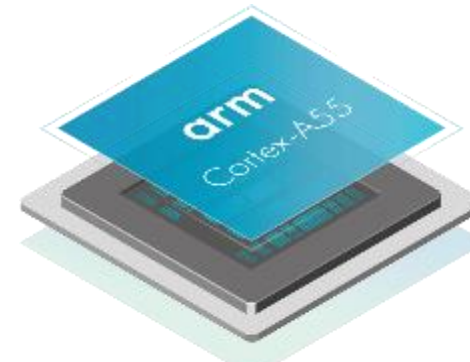


Estimated device performance using SPECINT2006, final device results may vary
Comparison using Cortex-A73 at 2.4GHz vs Cortex-A75 at 3GHz

Cortex-A55 processor

2.5x

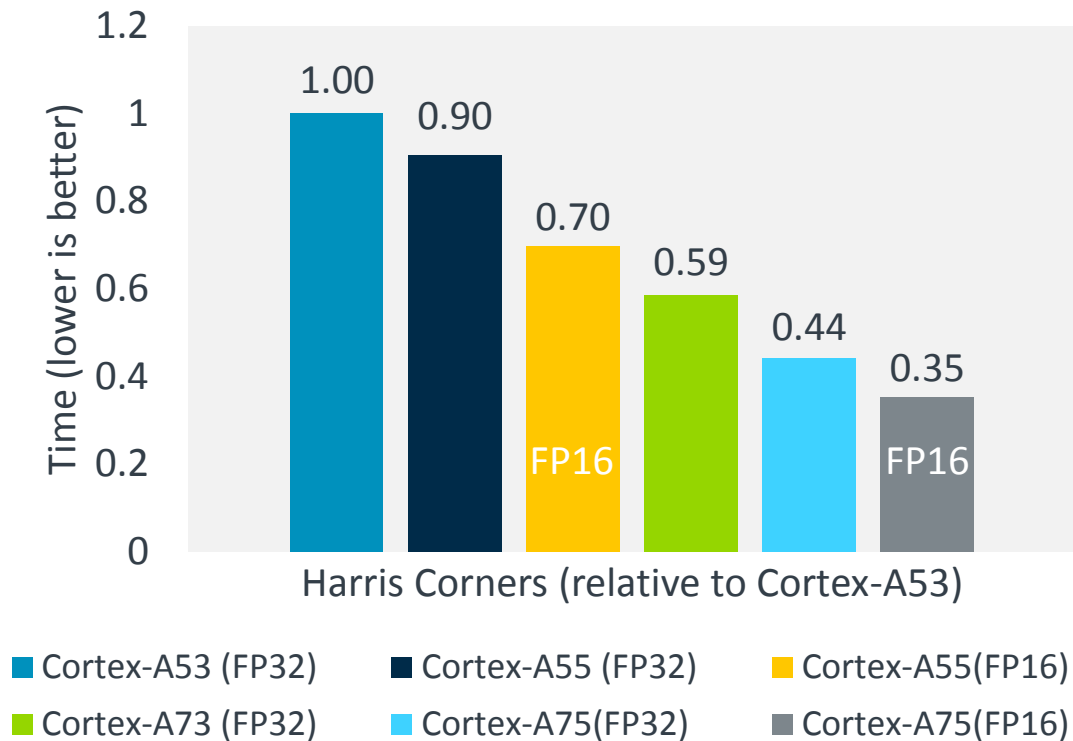
higher power efficiency
compared to current devices



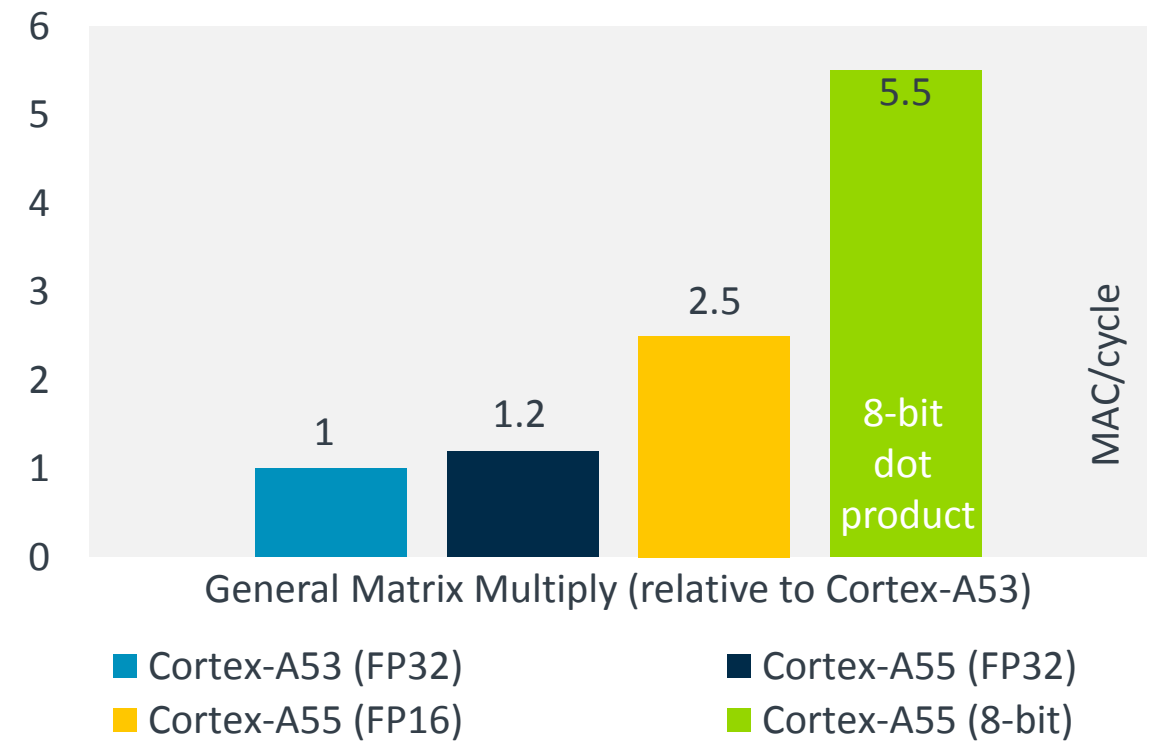
Comparison using Cortex-A53 in 28nm devices vs Cortex-A55 in 16nm devices

Enhanced Architecture for Emerging Use Cases

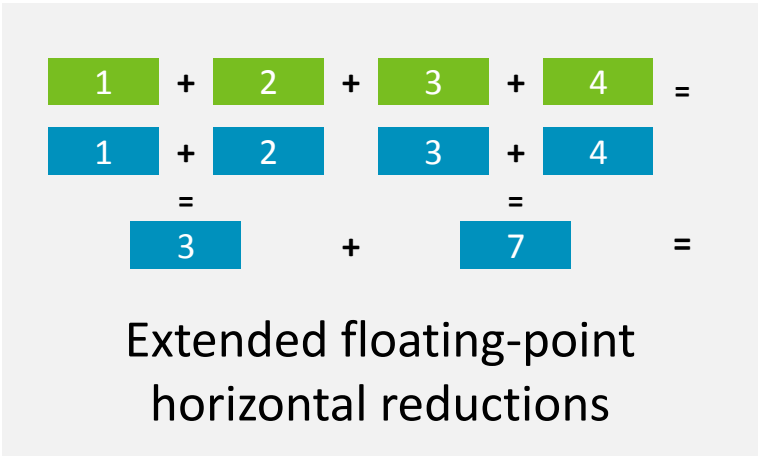
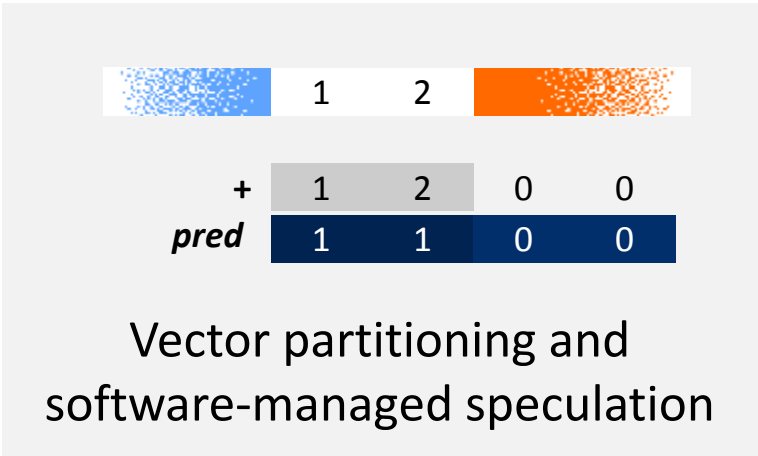
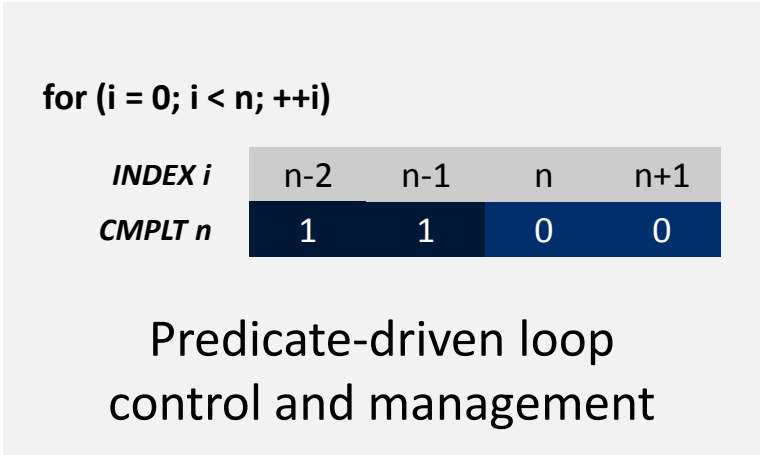
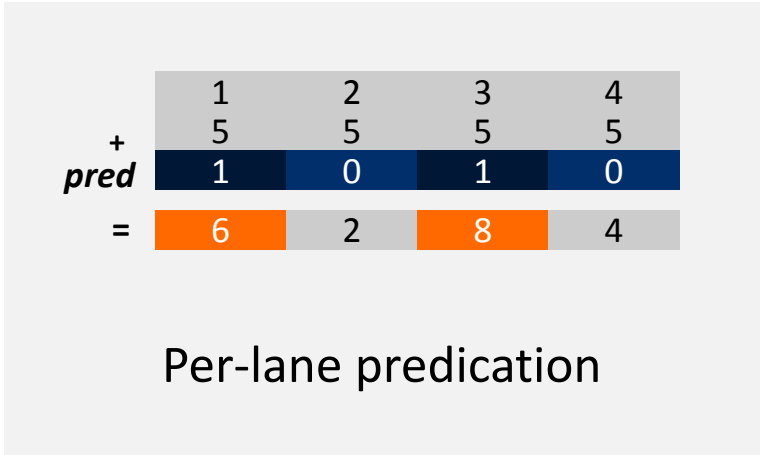
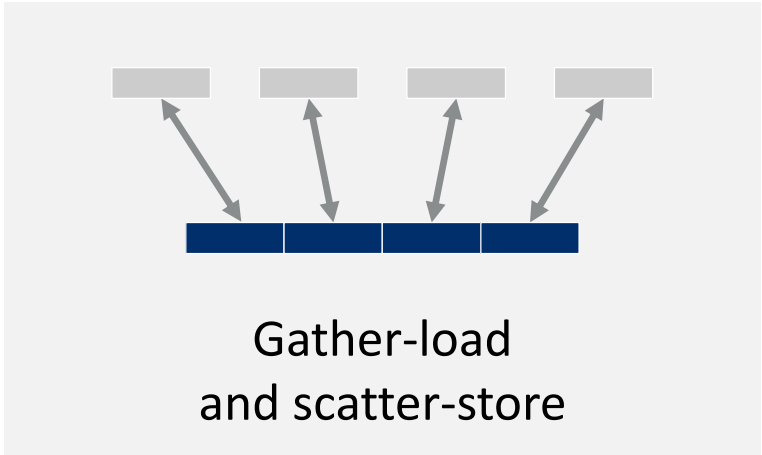
Computer Vision



Machine Learning



Introducing the Scalable Vector Extension (SVE)



Instruction Sets for AI

Cortex-A

- Additional dot product instructions (Cortex-A55 and Cortex-A75)
- New SVE instructions

Cortex-M

- Optimized CMSIS-DSP libraries for matrix multiplication

Closely-coupled acceleration

- Improved performance and efficiency (for broader use cases)
- Flexibility in multi-core computing with DynamIQ technology

Software Optimizations: Cortex-M Example

Convolution

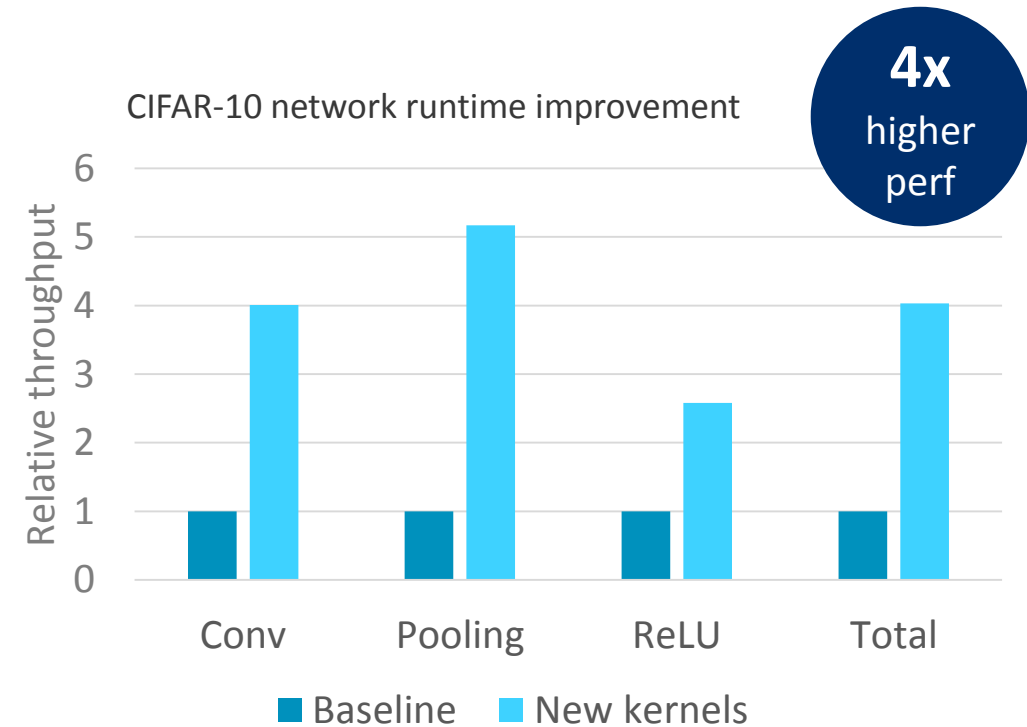
- Use of partial im2col to reduce the memory footprint
- Optimized data dimension layout (Height-Width-Channel) to save im2col overhead

Pooling

- Split into x-pooling and y-pooling instead of window-based
- 5.1X improvements compared to Caffe-like implementation

Activation

- ReLU: use SIMD within a register, 2.6X improvement compared to Caffe-like implementation
- Sigmoid and Tanh: use table look-up



*Baseline uses CMSIS 1D Conv and Caffe-like Pooling/ReLU

The new kernels will be integrated into future versions of CMSIS

Instruction sets for AI

Cortex-A

- Additional dot product instructions (Cortex-A55 and Cortex-A75)
- New SVE instructions

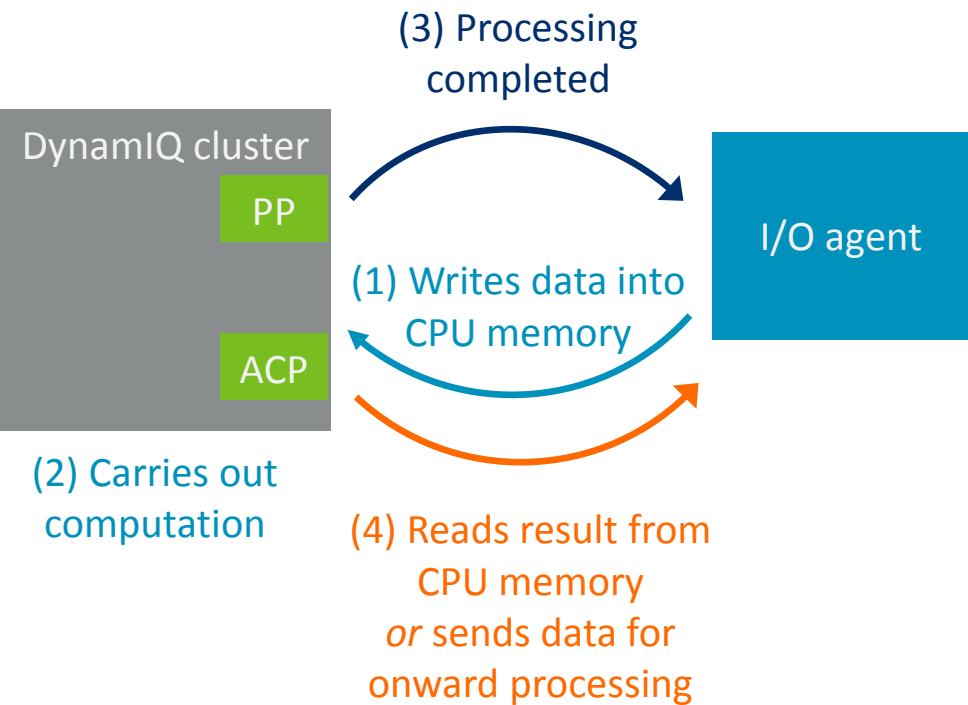
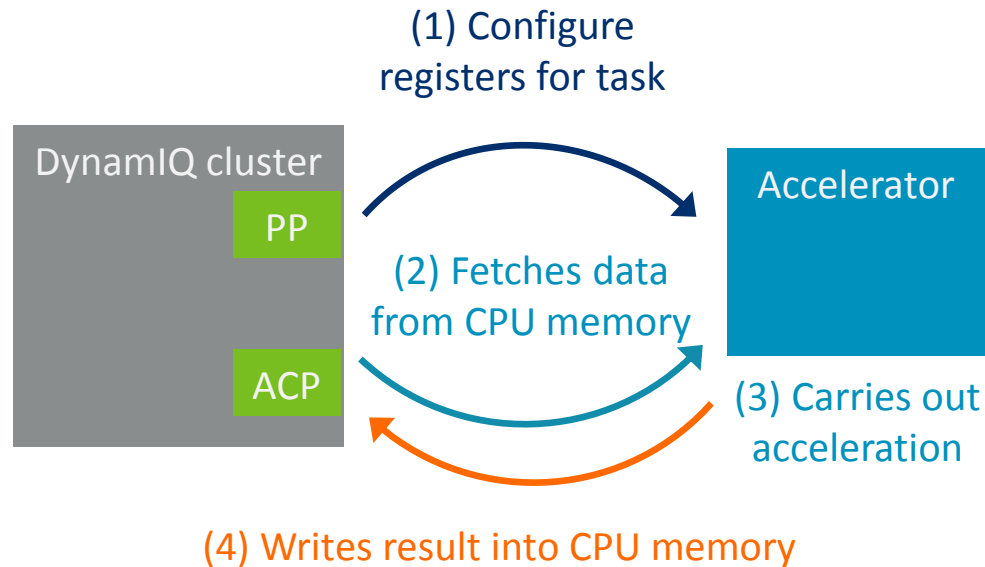
Cortex-M

- Optimized CMSIS-DSP libraries for matrix multiplication

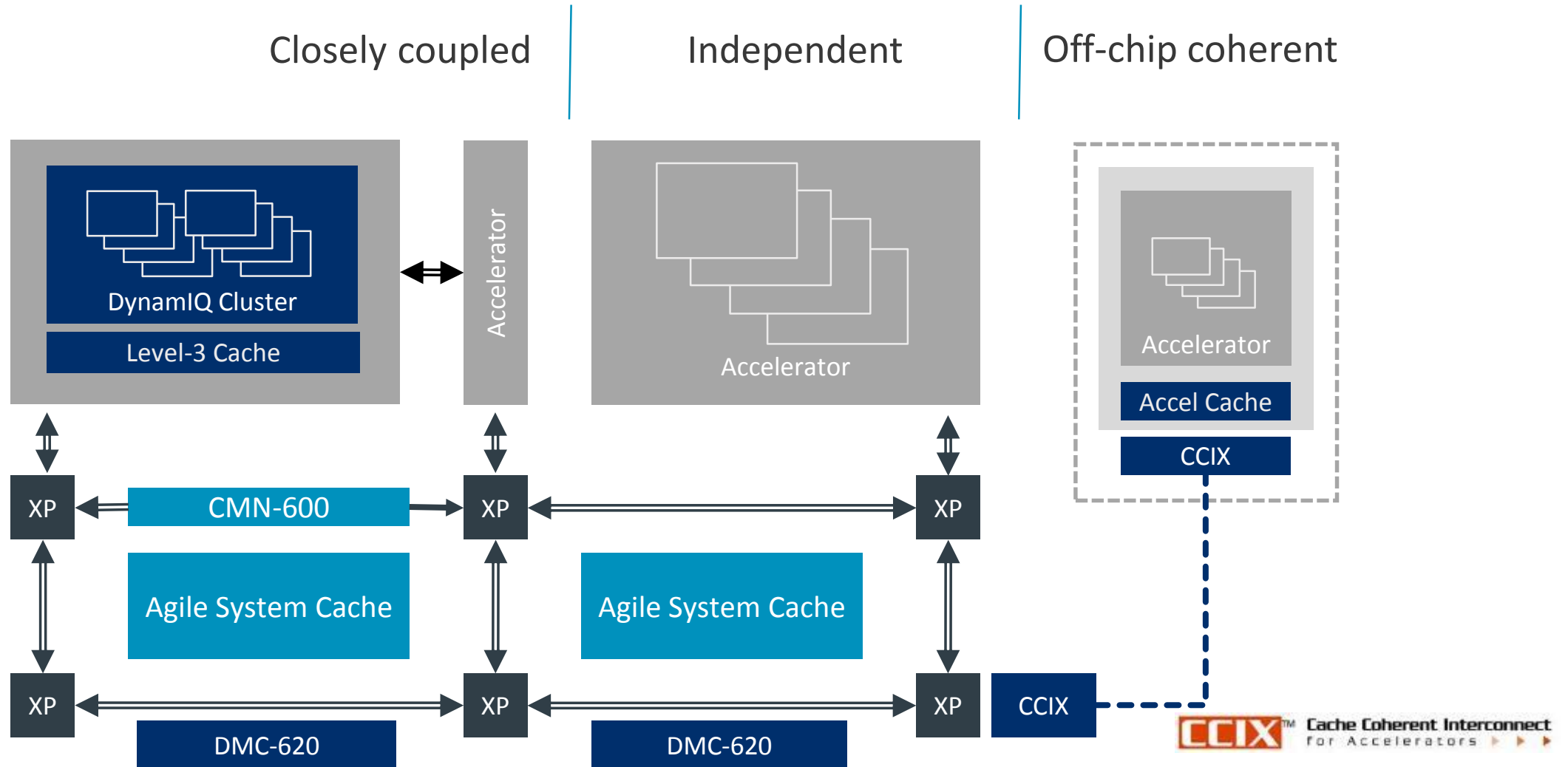
Closely-coupled acceleration

- Improved performance and efficiency (for broader use cases)
- Flexibility in multi-core computing with DynamIQ technology

Interface to Acceleration Logic



Flexible Acceleration Platform



Arm ML Platform Enables



Efficiency



Flexibility



Freedom

Thank You!

Danke!

Merci!

謝謝!

ありがとう!

Gracias!

Kiitos!

arm

arm

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