



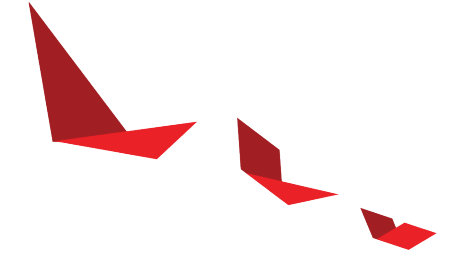
# Design Rationale of Two Generations of AI Engines

Kees Vissers  
Fellow

September 7, 2021

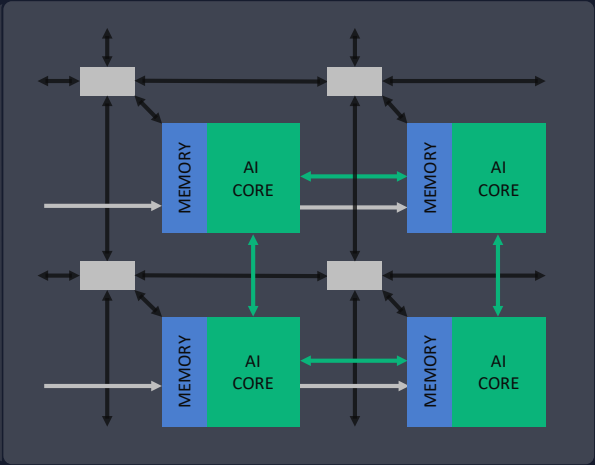
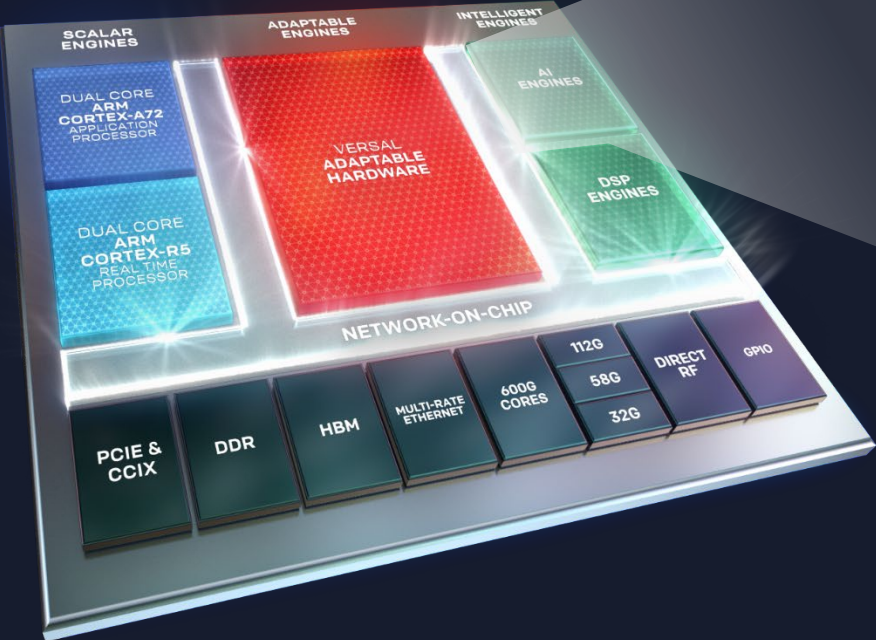


# Agenda



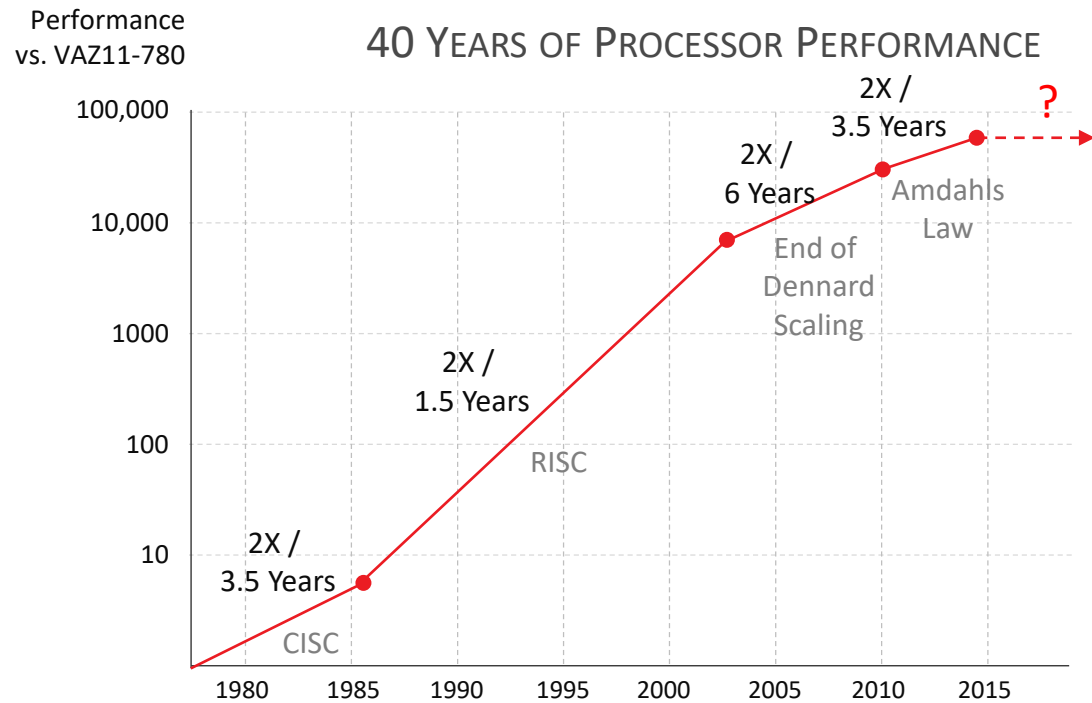
- > **Motivation for the first AI Engine architecture**
- > **AI Engine architecture in more detail**
- > **Programming**
- > **Some Applications**
- > **Motivation for the second generation AI Engine, AI – ML**
- > **AI-ML in more detail**
- > **Some Applications**

# Motivation for AI Engine

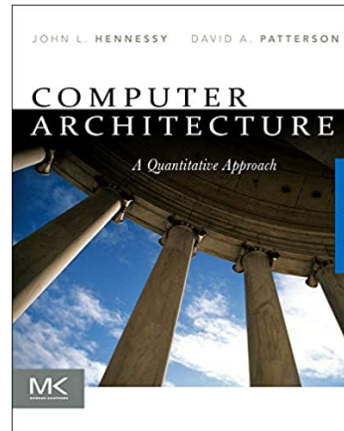


# Computer Architecture: choices, a new golden age!

Performance is not scaling: end of Moore's Law



Source: John Hennessy and David Patterson, Computer Architecture: A Quantitative Approach, 6/e 2018



2017 Turing award winners, talk:

“A New Golden Age for Computer Architecture: Domain-Specific Hardware/Software Co-Design, Enhanced Security, Open Instruction Sets, and Agile Chip Development,”



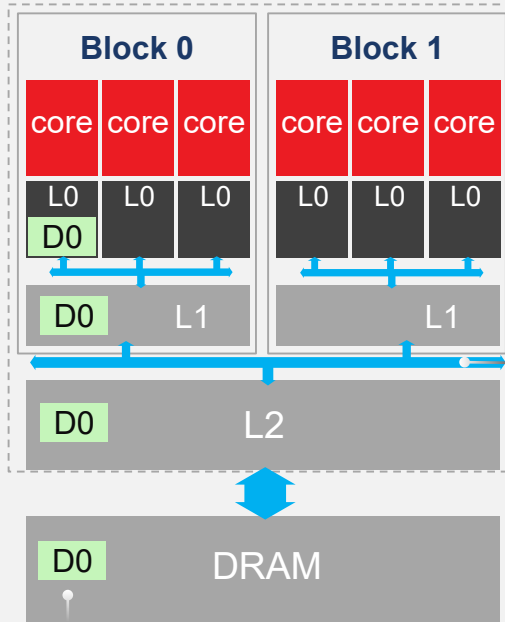
<https://www.acm.org/hennessy-patterson-turing-lecture>

# Original design rationale for AI engine

- > **Trend towards heterogeneous architectures**
- > **Can a novel array of dedicated, Xilinx designed processors, complement Programmable Logic?**
- > **Can these processors provide cost and power benefit for the customers in the domains of:**
  - >> Wireless applications
  - >> General signal processing and image processing applications
  - >> Emerging Machine Learning applications
- > **Can you raise the programming abstraction to programming environments:**
  - >> C/C++
  - >> Machine Learning Frameworks

# AI Engine: Multi-Core Compute with dedicated memory

## Traditional Multi-core (cache-based architecture)



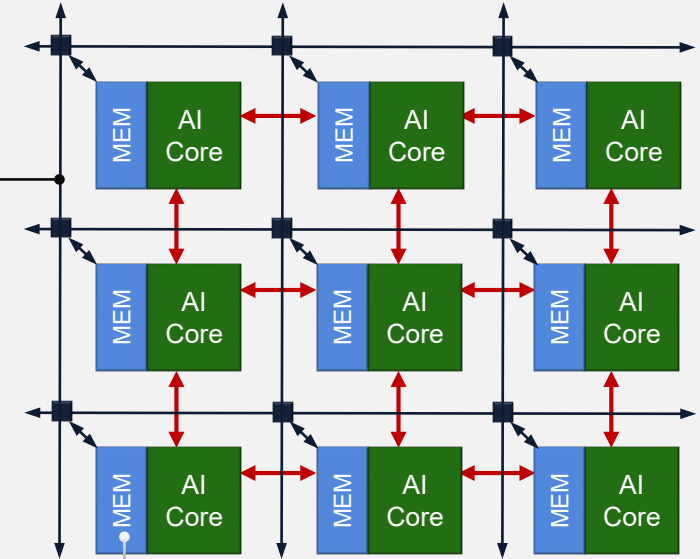
Fixed, shared Interconnect

- Blocking limits compute
- Timing not deterministic

Data Replicated

- Robs bandwidth
- Reduces capacity

## AI Engine Array (intelligent engine)





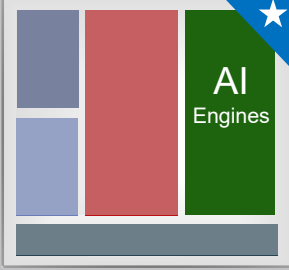

Dedicated Interconnect

- Non-blocking
- Deterministic

Local, Distributed Memory

- No cache misses
- Higher bandwidth
- Less capacity required

# Delivering Compute Acceleration

	CPU (Sequential) 	GPU (Parallel) 	Versal 	Custom ASIC 
SW Programmable	✓	✓	✓	✓
HW Adaptable	—	—	✓	—
Workload Flexibility	✓	✓	✓	—
Throughput vs. Latency	—	—	✓	✓
Device / Power Efficiency	—	—	✓	✓



# Versal Architecture Overview



**Adaptable Engines**  
2X compute density



## Scalar Engines

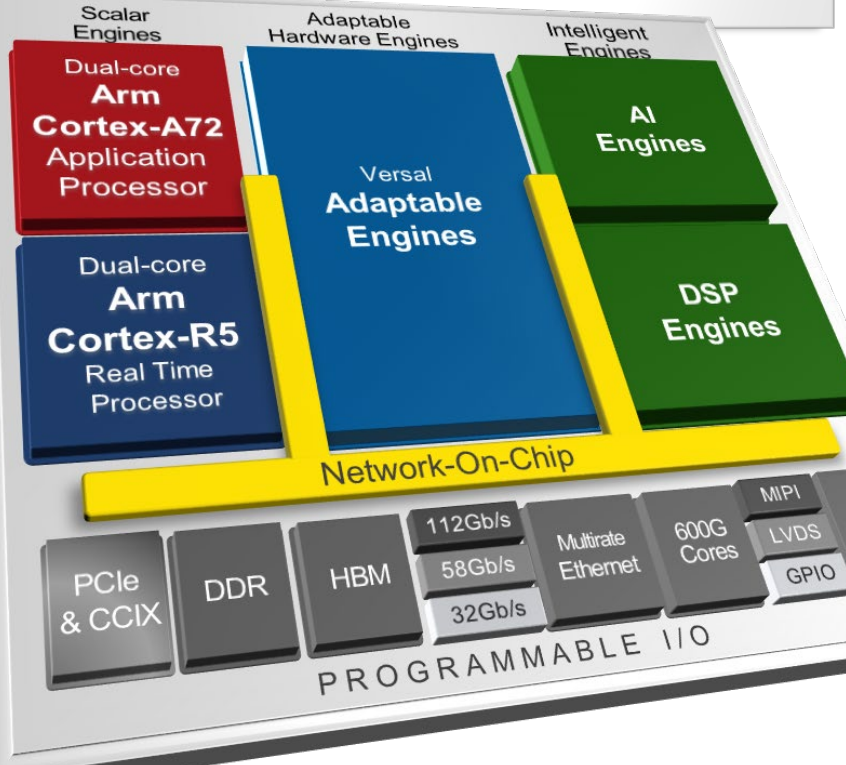
- Platform Control
- Edge Compute

## Protocol Engines

- Integrated 600G cores
- 4X encrypted bandwidth

## Programmable I/O

- Any interface or sensor
- Includes 4.2Gb/s MIPI



## AI Engines

- AI Compute
- Diverse DSP workloads

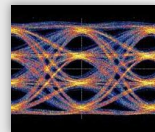


## Network-on-Chip

- Guaranteed Bandwidth
- Enables SW Programmability

## DDR Memory

- 3200-DDR4, 3200-LPDDR4
- 2X bandwidth/pin



## Transceivers

- Broad range, 25G → 112G
- 58G in mainstream devices



## PCIe & CCIX

- 2X PCIe & DMA bandwidth
- Cache-coherent interface to accelerators

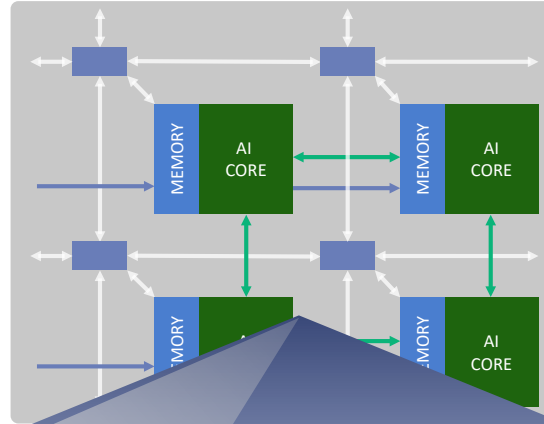


# Introducing the AI Engine

SW Programmable ➤

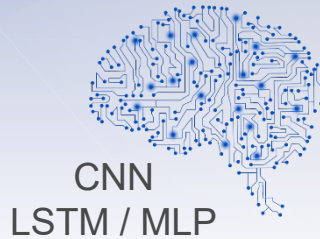
Deterministic ➤

Efficient ➤

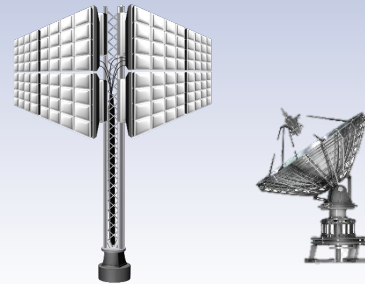


- 1GHz+ Multi-precision Vector Processor
- High bandwidth extensible memory
- Up to 400 AI Engines per device
- 8X Compute Density
- 40% Lower Power

Artificial Intelligence



Signal Processing



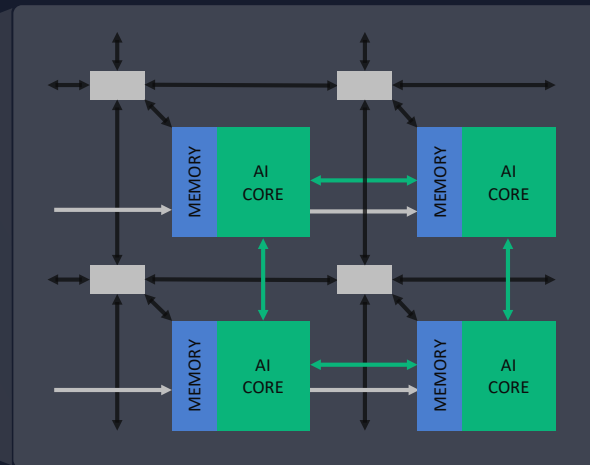
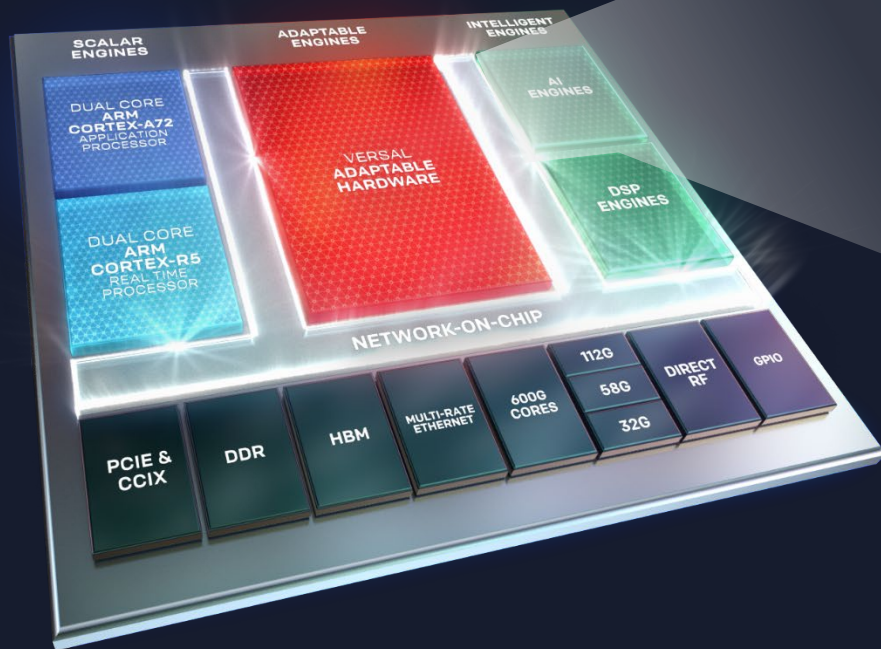
Computer Vision



## High performance multi-core processor system

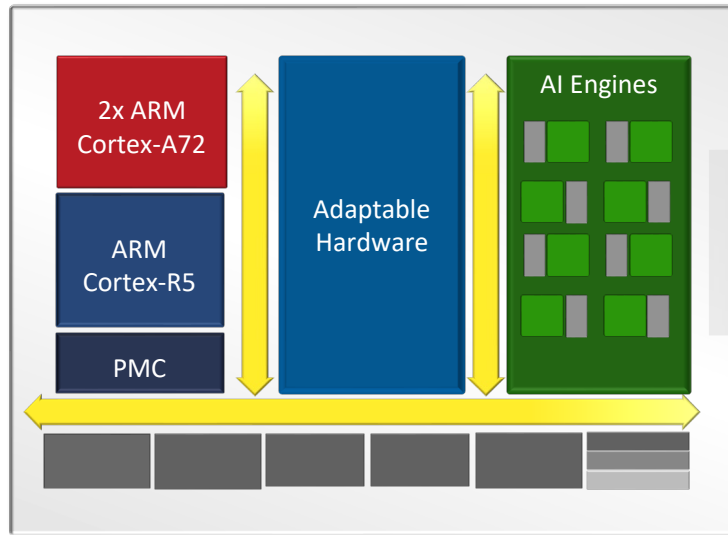
# AI Engine (first generation)

## Architecture, Programming & Applications

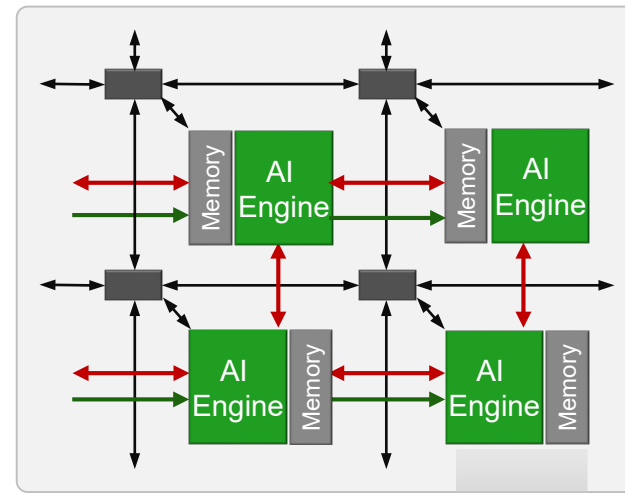


# AI Engine: Terminology

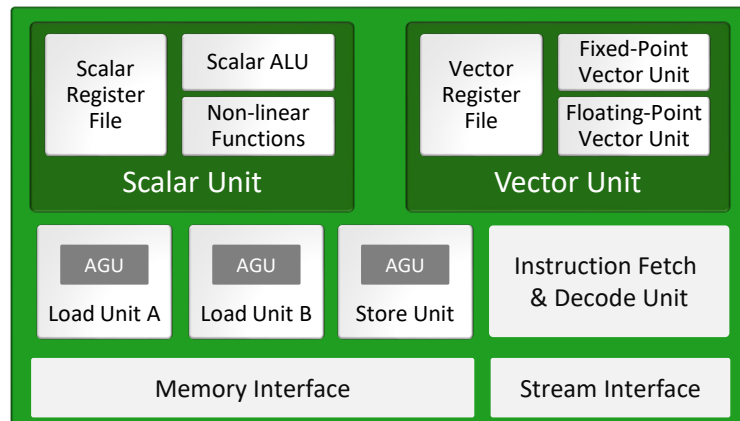
Versal AI



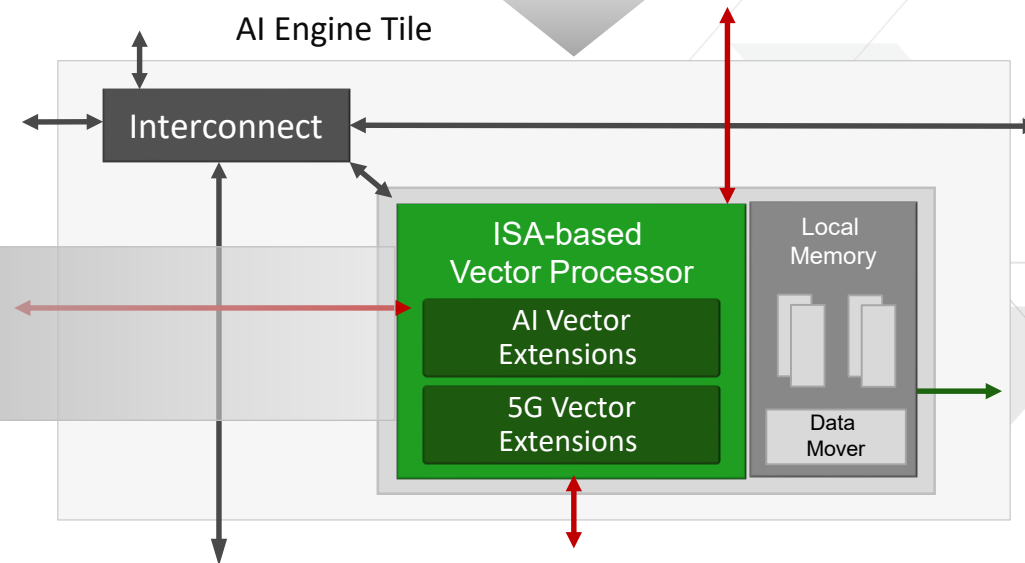
AI Engine Array



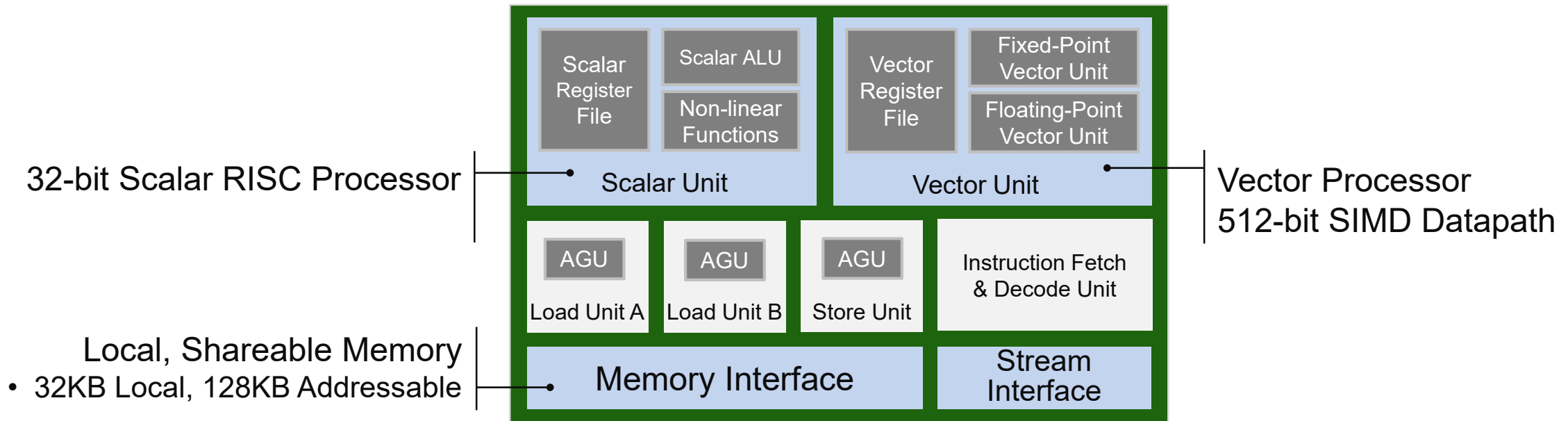
AI Engine Core



AI Engine Tile



# AI Engine: Processor Core



## Instruction Parallelism: VLIW

7+ operations / clock cycle

- 2 Vector Loads / 1 Mult / 1 Store
- 2 Scalar Ops / Stream Access

Highly Parallel

## Data Parallelism: SIMD

Multiple vector lanes

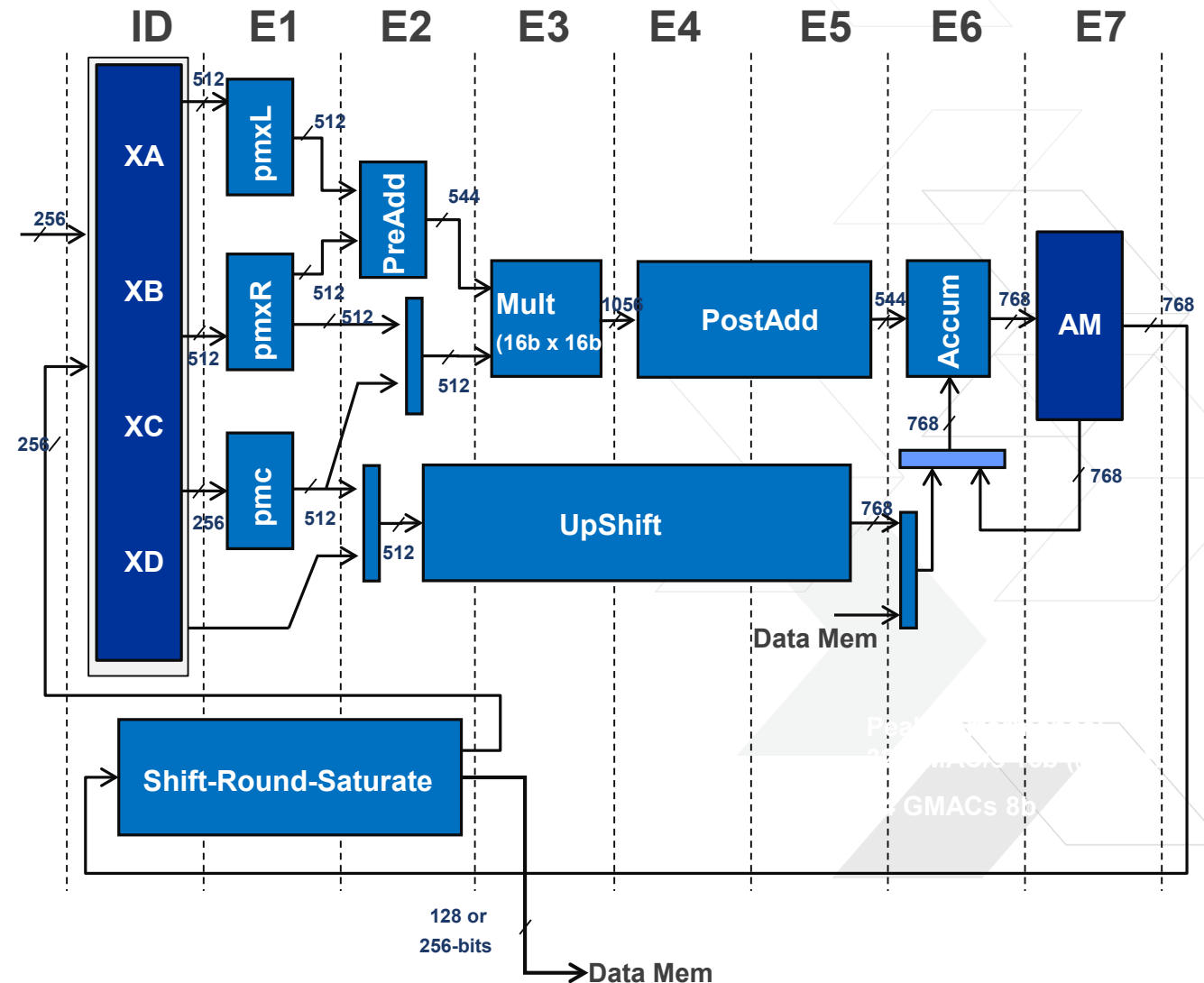
- Vector Datapath
- 8 / 16 / 32-bit & SPFP operands

Up to 128 MACs / Clock Cycle per Core (INT 8)

# AIE Core: Fixed-Point Vector Unit

## > Vector fixed-point/integer unit

- >> Multiple precision for complex and real operand
- >> Full permute unit (32-bit granularity)
- >> Pre-adder
- >> Multiplier
- >> Two-step post-adding
- >> Accumulator
- >> Shift, round, saturate



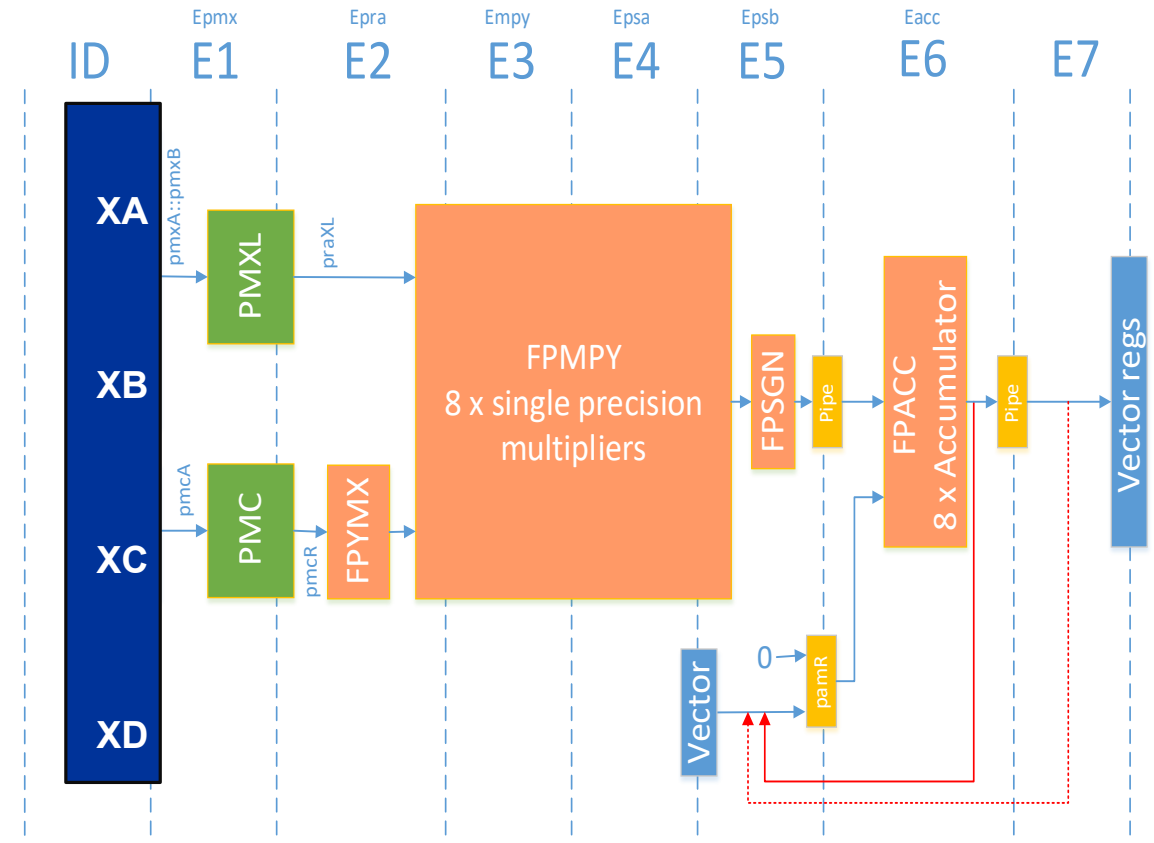
# AIE Core: Floating-point Vector Unit

## > Vector floating-point Unit

- >> Single precision
- >> 8 multiply-accumulate per cycle
- >> Sign change (FPSGN) is on per-lane basis

## > Exceptions

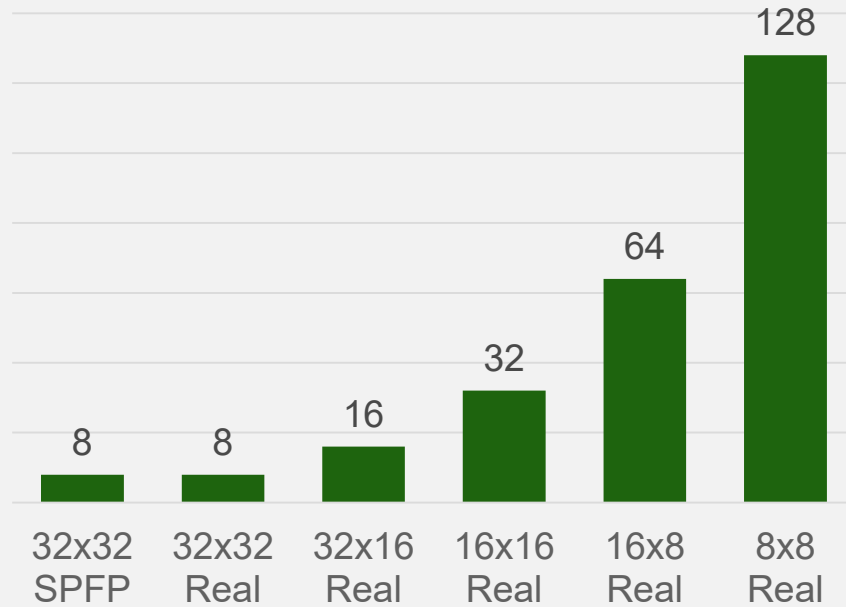
- >> Per-lane exceptions
- >> ZERO, INFINITY, INEXACT, HUGE\_INT
- >> TINY, HUGE, INVALID, DIVIDE\_BY\_ZERO
- >> The latter can be converted into events broadcasted to PL-AI interface and sent as interrupts to PS/PMC



# Multi-Precision Support

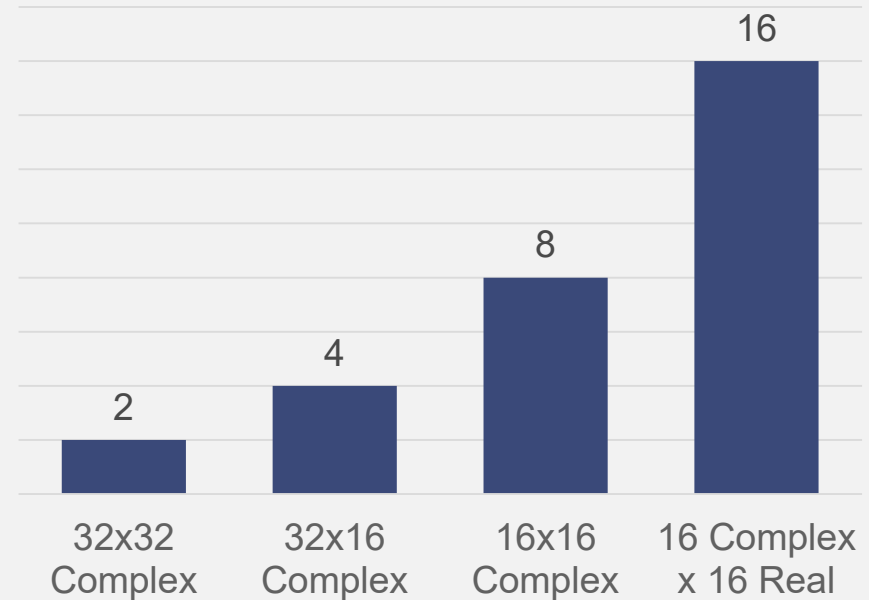
## AI Data Types

MACs / Cycle (per core)



## Signal Processing Data Types

MACs / Cycle (per core)

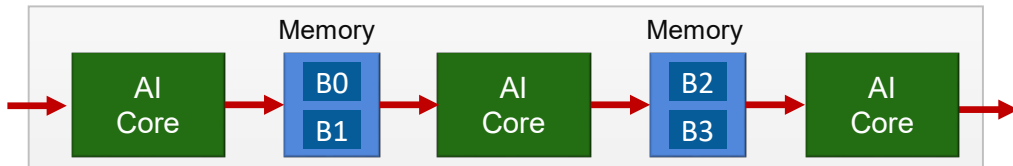




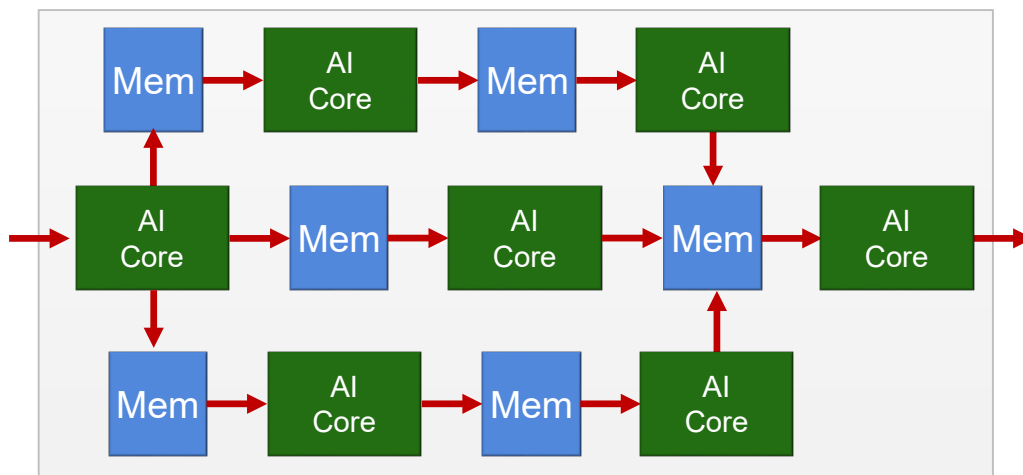
# Data Movement Architecture

## Memory Communication

Dataflow Pipeline



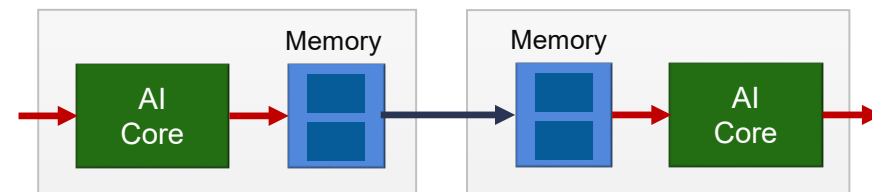
Dataflow Graph



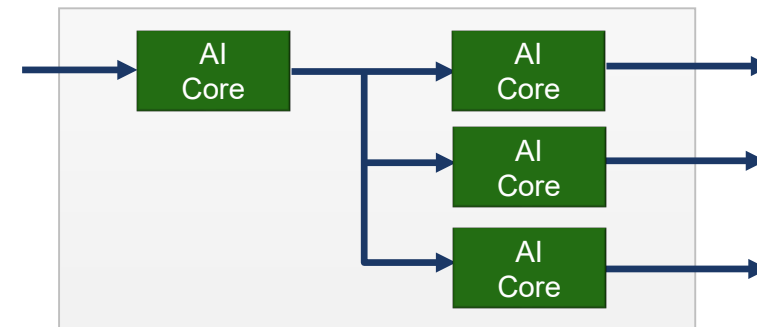
- Memory Interface
- Stream Interface
- Cascade Interface

## Streaming Communication

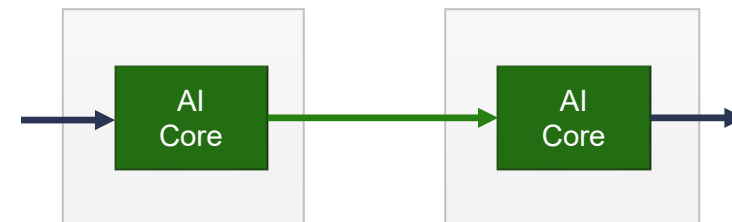
Non-Neighbor



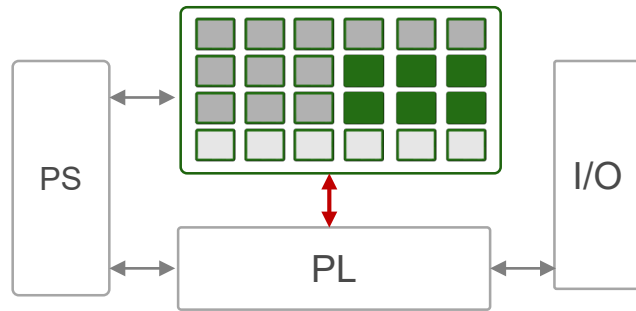
Streaming Multicast



Cascade Streaming



# AI Engine Integration with Versal

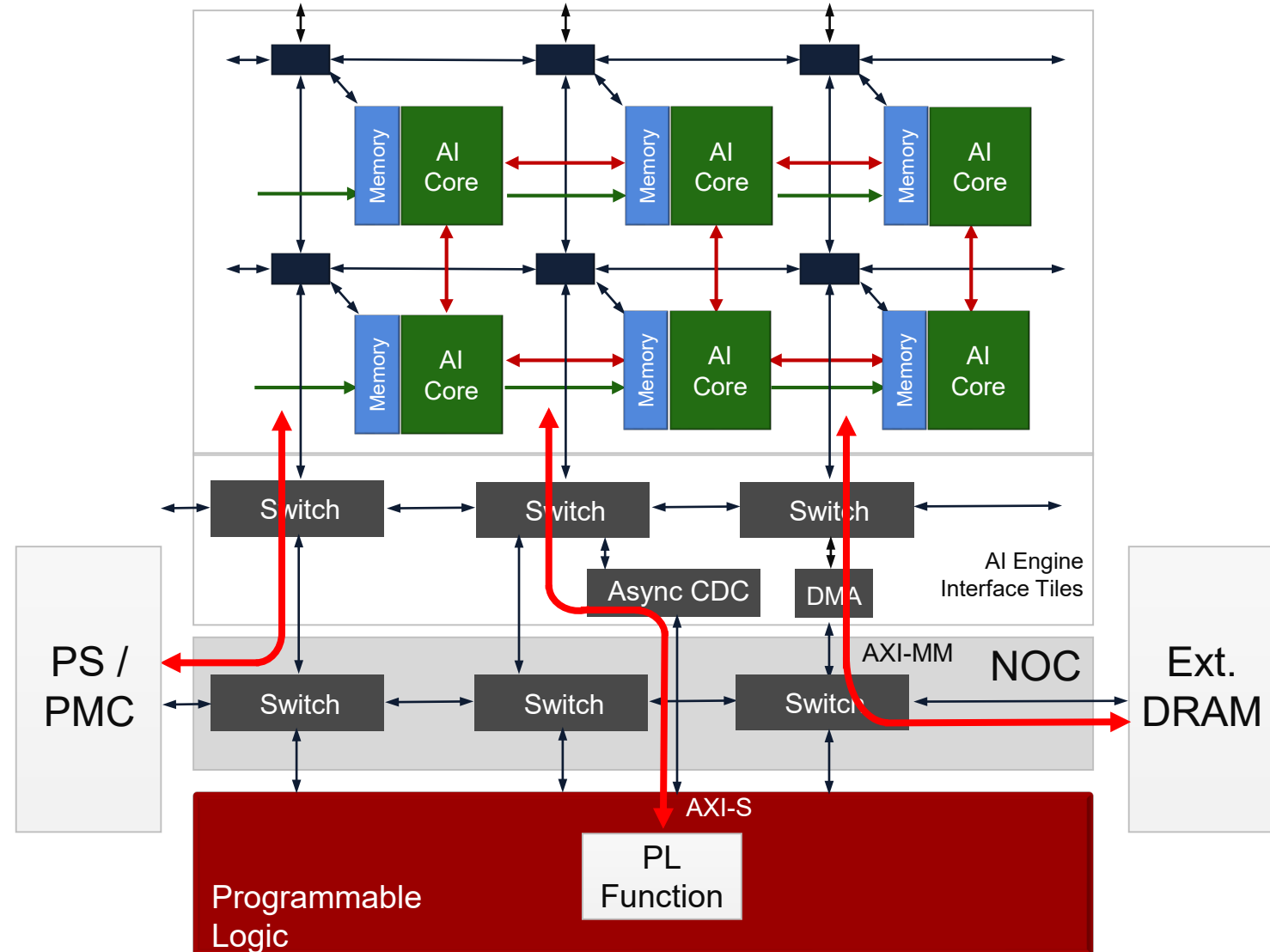


## > TB/s of Interface Bandwidth

- >> AI Engine to Programmable Logic
- >> AI Engine to NOC

## > Leveraging NOC connectivity

- >> PS manages Config / Debug / Trace
- >> AI Engine to DRAM (no PL req'd)



# AI Engine Delivers High Compute Efficiency

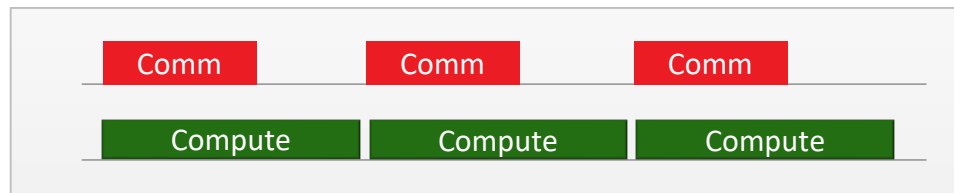
## > Adaptable, 'non-blocking' interconnect

- >> Flexible data movement architecture
- >> Avoids interconnect "bottlenecks"

## > Adaptable memory hierarchy

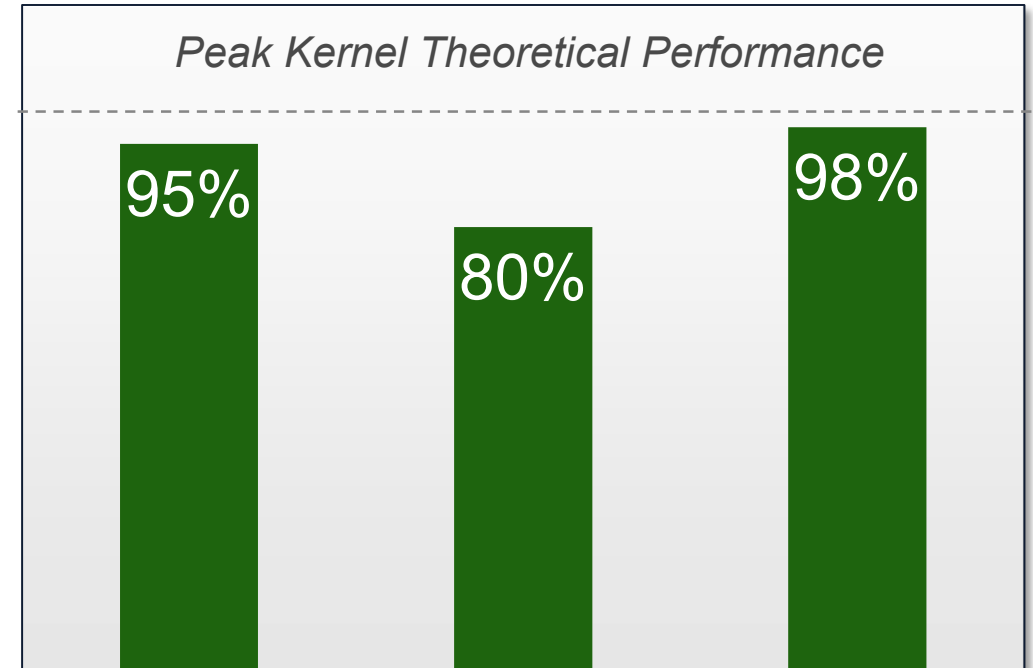
- >> Local, distributed, shareable = extreme bandwidth
- >> No cache misses or data replication
- >> Extend to PL memory (BRAM, URAM)

## > Transfer data while AI Engine Computes



Overlap Compute and Communication

## Vector Processor Efficiency



ML Convolutions

Block-based  
Matrix Multiplication  
(32×64) × (64×32)

FFT

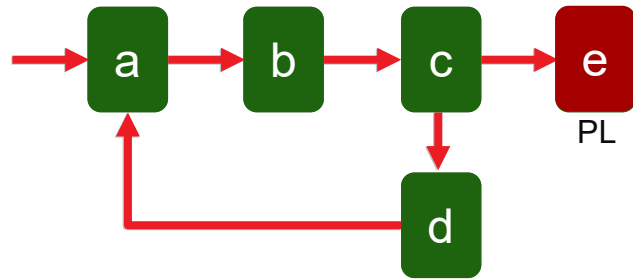
1024-pt  
FFT/iFFT

DPD

Volterra-based  
forward-path DPD

# AI Engine Programming Experience: Dataflow Model

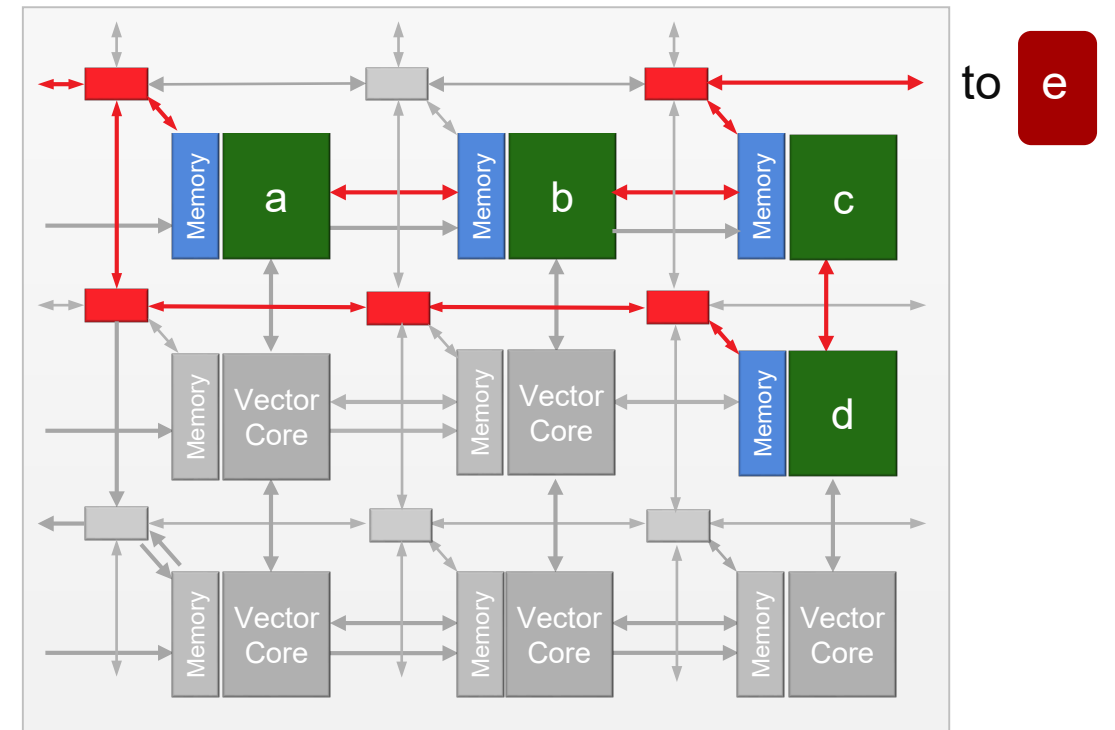
1 User defines dataflow logic



2 User describes dataflow graph using C/C++ APIs

3 Compiler transparently manages placement & interconnect

Physical Mapping to AI Engines



# All Developers Can Build and Deploy on All Platforms



Build



Embedded  
Developers



Enterprise  
Application Developers



Enterprise Infrastructure  
Developers



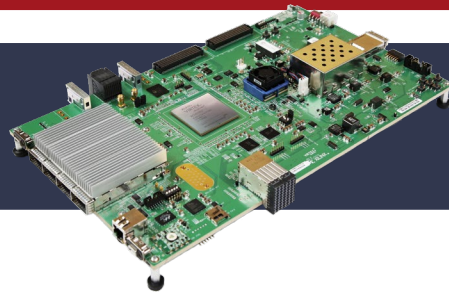
Data & AI  
Scientists



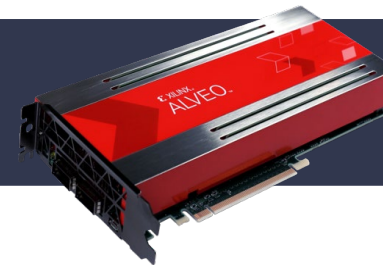
Deploy



Zynq-7000



Zynq UltraScale+ MPSoC



Alveo



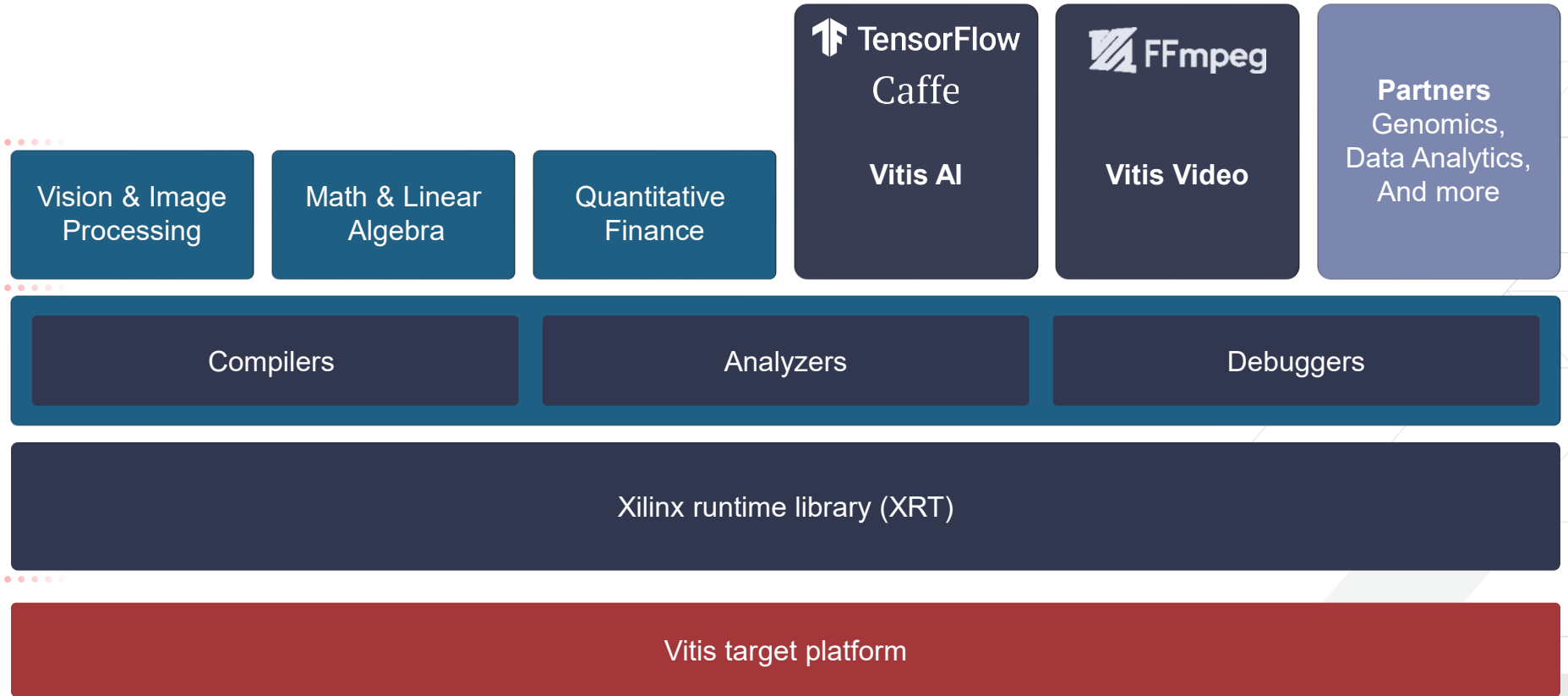
Versal ACAPs

# Vitis Unified Software Platform

Domain-specific development environments

Vitis accelerated libraries

Vitis core development kit



Edge Deployment

On-Premise Deployment

Cloud Deployment

# Vitis libraries, including AIE Libraries

## Domain-Specific Libraries



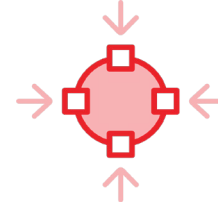
Vision & Image



Quantitative Finance



Data Analytics & Database

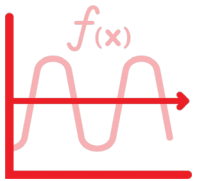


Data Compression

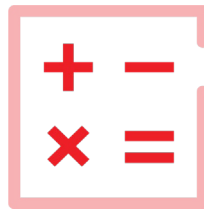


Data Security

## Common Libraries



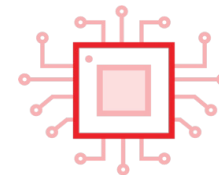
Math



Linear Algebra



Statistics



DSP



Data Management

500+ functions across multiple libraries for performance-optimized out-of-the-box acceleration



# AIE Designs for massive MIMO and DFE Platforms

A set of highly efficient and scalable AIE designs of product quality

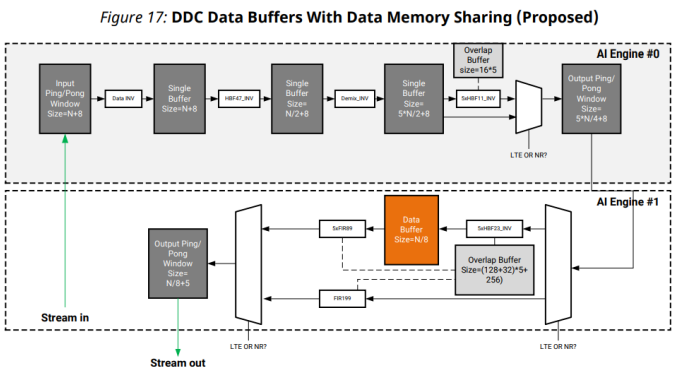
>> Advanced design techniques published on Xilinx.com as application notes with source code

Application Note: Versal™ AI Core Devices




**Digital Down-conversion Chain Implementation on AI Engine**

XAPP1351 (v1.0) February 15, 2021



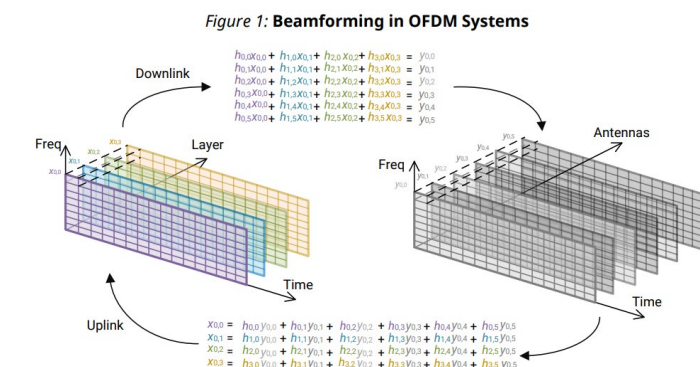
XAPP 1351  
Multi-carrier Multi-rate Filters

Application Note: Versal™ AI Core Devices



**Beamforming Implementation on AI Engine**

XAPP1352 (v1.0) January 11, 2021



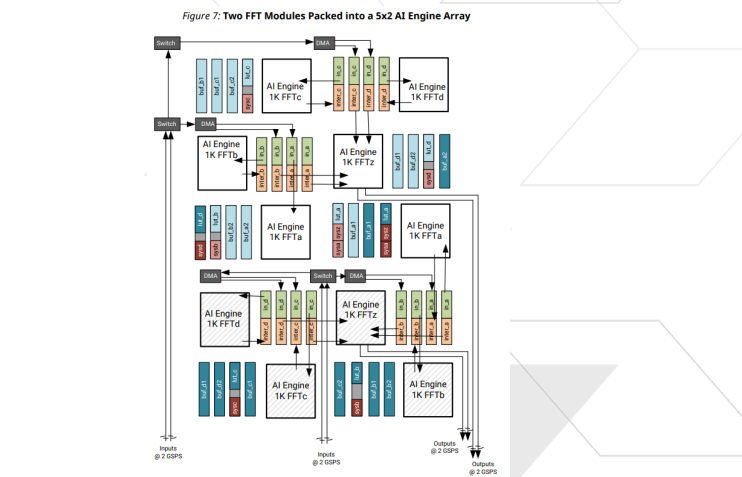
XAPP 1352  
Scalable beamforming

Application Note: Versal™ AI Core Devices



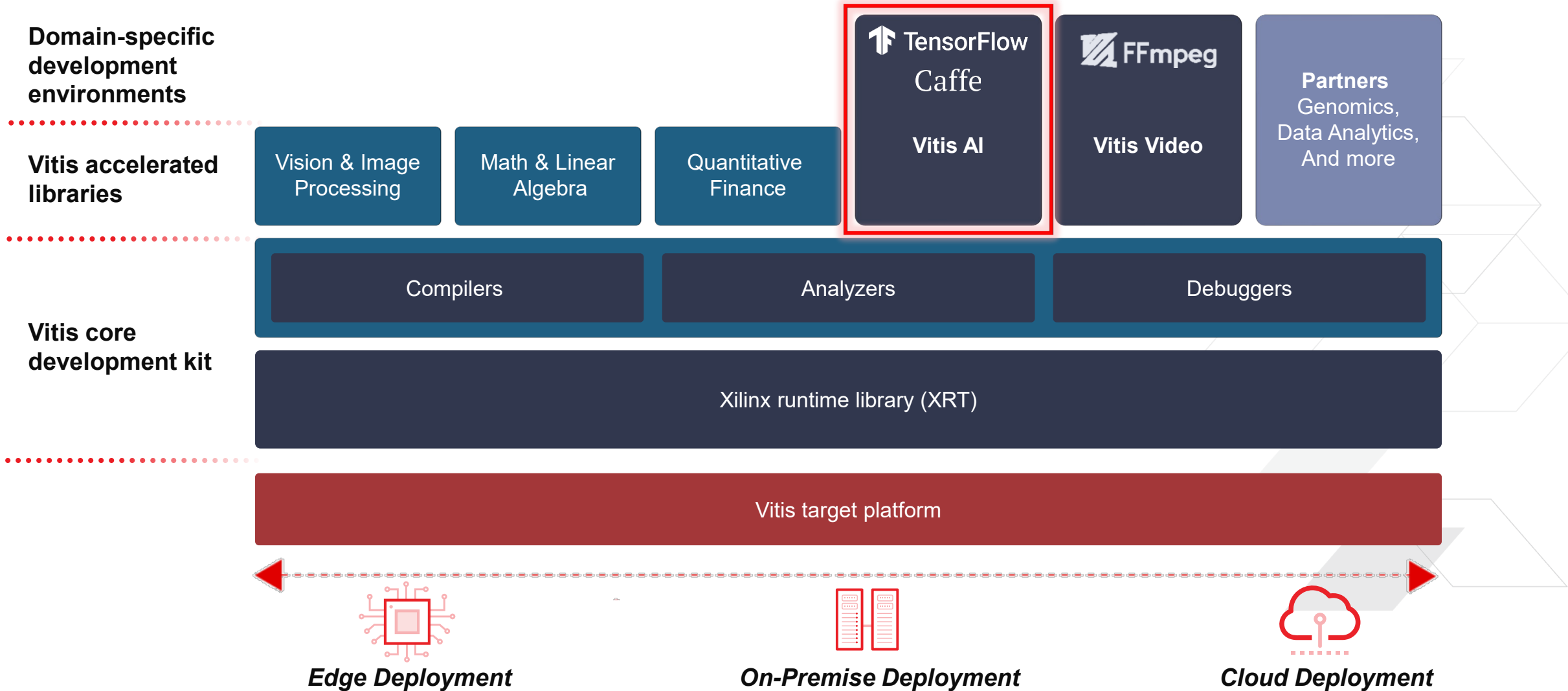
**Block-by-Block Configurable Fast Fourier Transform Implementation on AI Engine**

XAPP1356 (v1.0) January 11, 2021



XAPP 1356  
Configurable FFT/IFFT

# Vitis Unified Software Platform

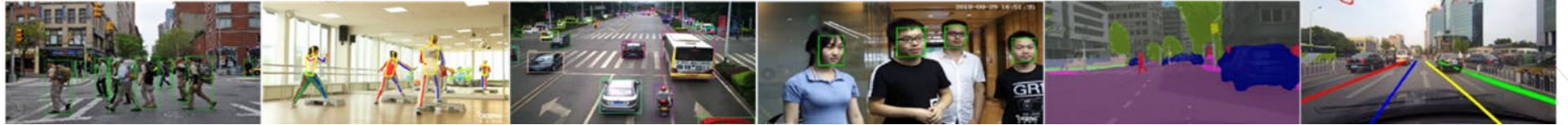


# Vitis AI: Deep Learning Acceleration

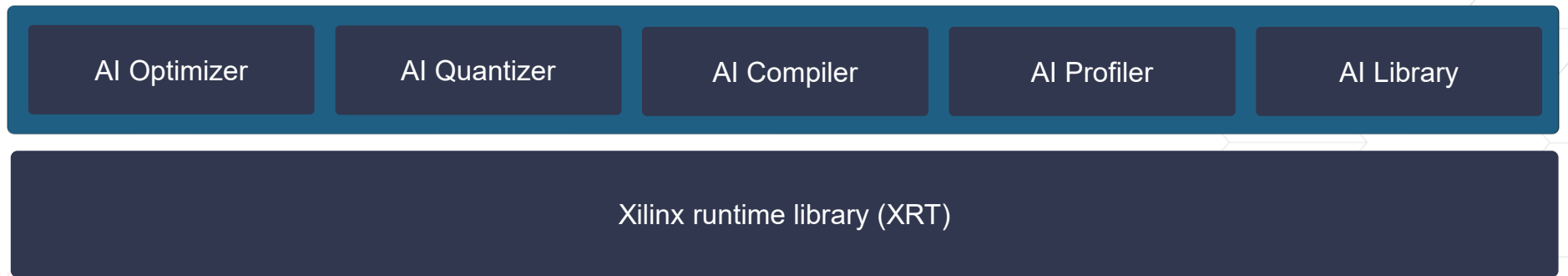
Frameworks



Vitis AI models



Vitis AI development kit



Deep Learning Processing Unit (DPU)



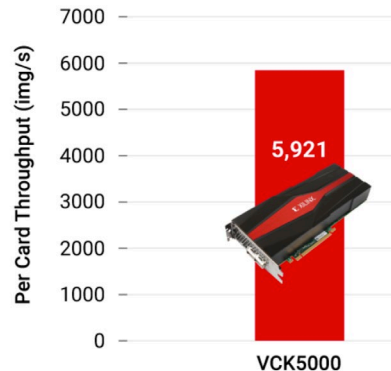
Edge Deployment

On-Premise Deployment

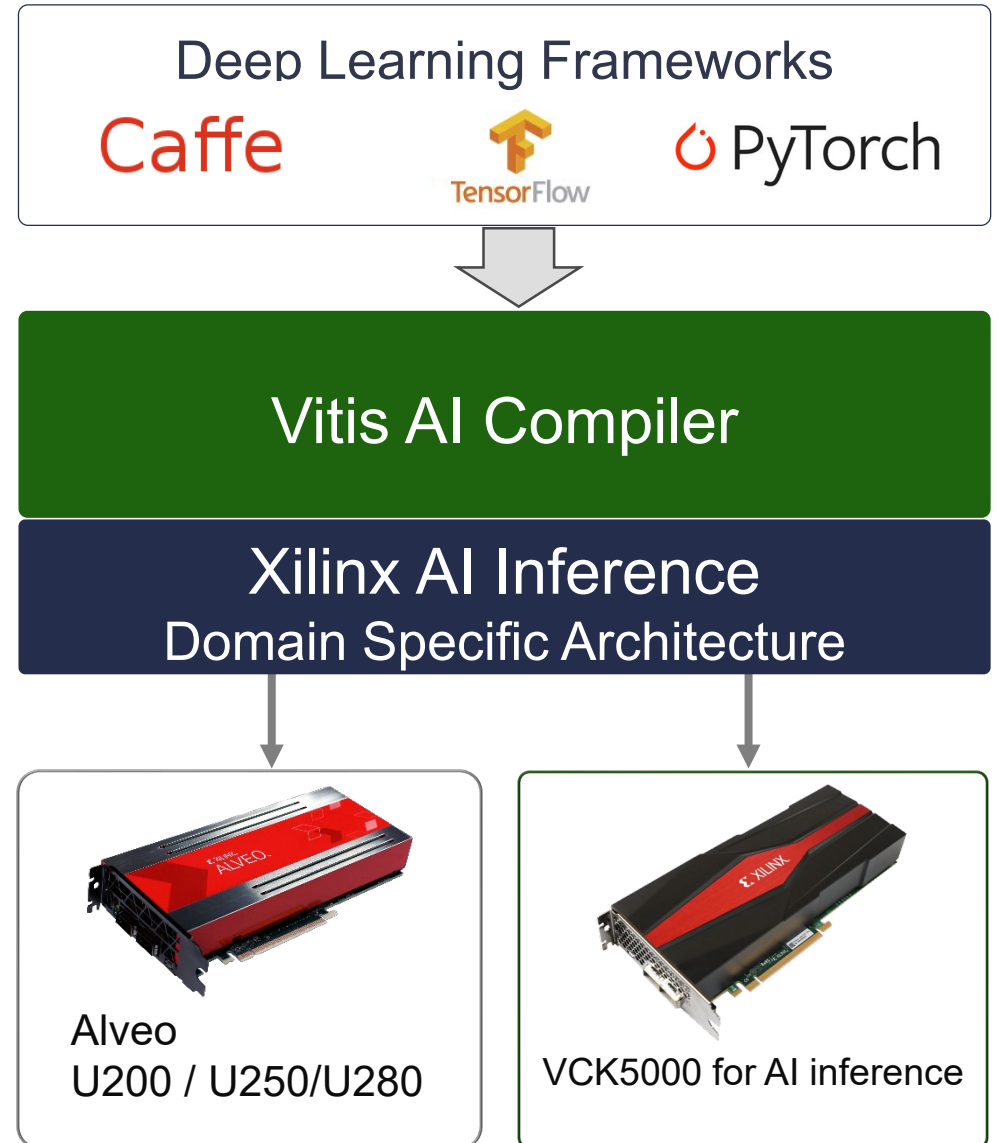
Cloud Deployment

# Vitis AI: Inference in the Data Center

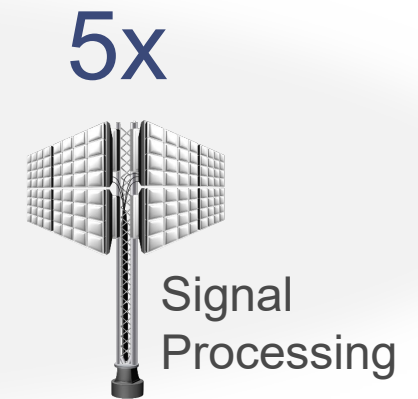
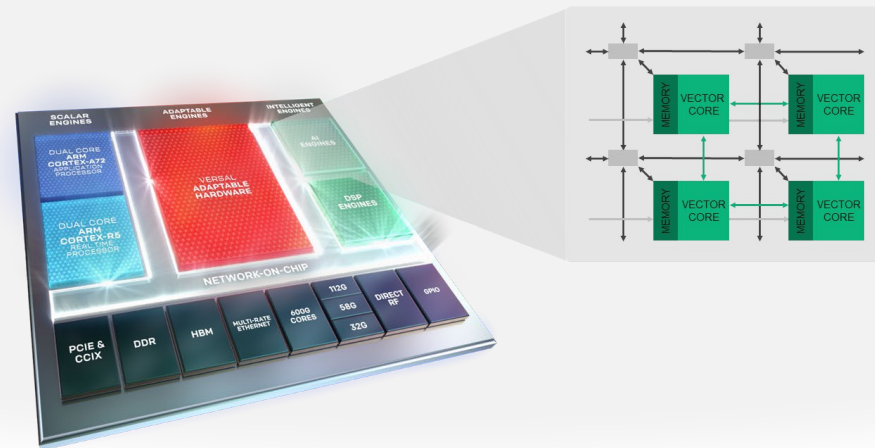
- 1 User works in Framework of choice
  - Develop & train custom network
  - User provides trained model
- 2 Xilinx Vitis AI Compiler implements network
  - Targets AI Inference Domain Specific Architecture
  - Quantize, merge layers, prune
  - Compile to AI Engines
- 3 Scalable across hardware targets
  - AI engine performance benefit 8x -10x



ML commons v1.0 Data Center Closed Division Server Resnet-50



# AI Engine: Accelerating AI Inference & Signal Processing



## Software Programmable

- Frameworks & C/C++
- SW Compile, Debug & Deploy

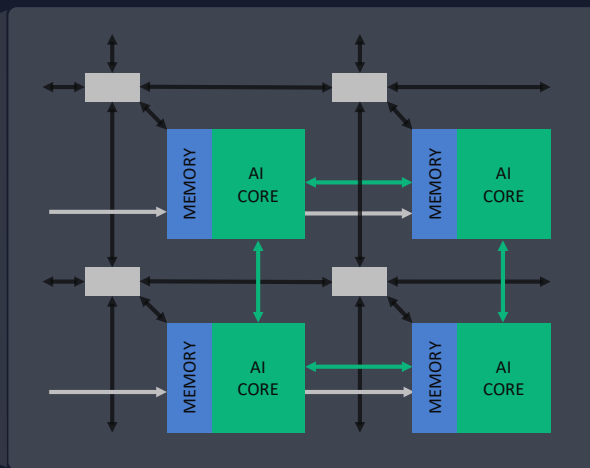
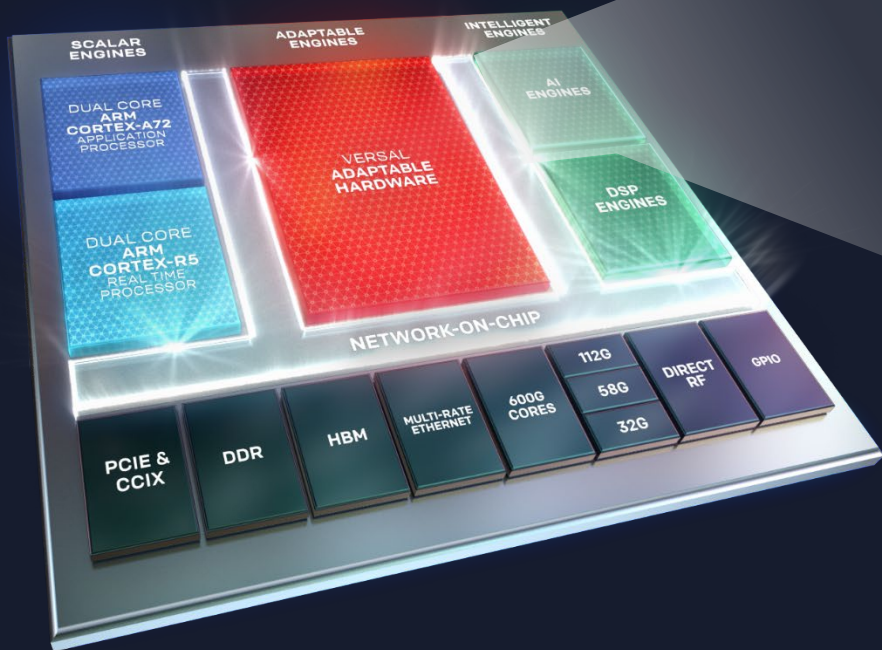
## Deterministic

- Max throughput w/ low latency
- Real-time inference leadership

## Efficient

- Up to 8X compute density
- At ~40% lower power

# Motivation for second generation: AI - ML





# AI/ML Edge Inference: Use-cases and Requirements



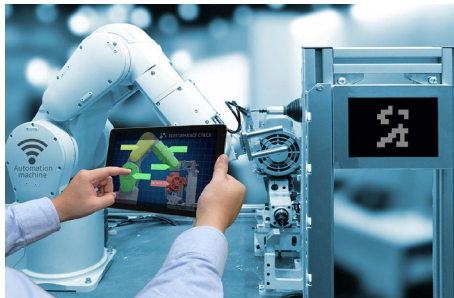
**Automotive**  
ADAS  
Automated Driving



**Industrial Robotics**  
Vision guidance  
Environmental awareness



**Medical**  
Ultrasound  
Endoscopy, Surgical Robotics



**Industrial IoT**  
Industrial PCs  
Smart Grid Controllers



**Vision & Smart City**  
Machine Vision Camera  
Edge AI Box



**Aerospace & Defense**  
Unmanned Aerial Vehicles  
Multi-Mission Payload Systems

**Challenge: Delivering massive AI compute at low-latency and low-power**



# Versal™ AI Edge: Architecture Overview



## ADAPTABLE TO MULTIPLE WORKLOADS

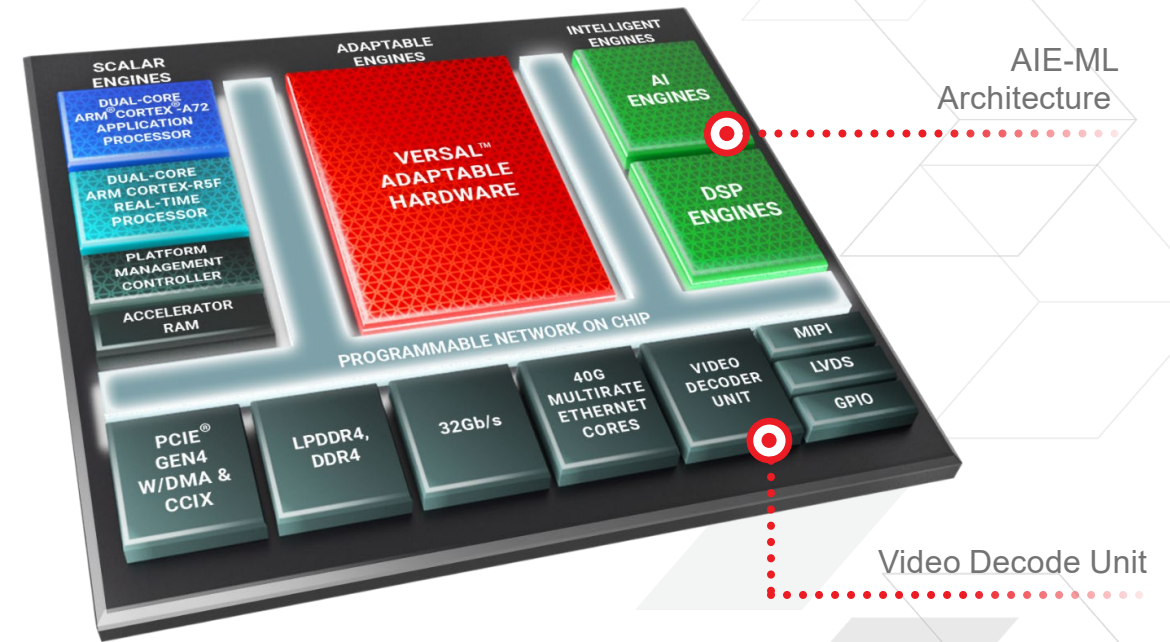
- ▶ Versal AI Core
- ▶ Versal Premium
- ▶ Versal AI Edge

## COMPUTE ACCELERATION

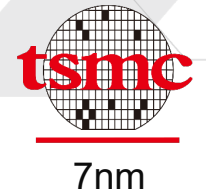
- ▶ Scalar Engines
- ▶ Adaptable Engines
- ▶ Intelligent Engines

## PLATFORM

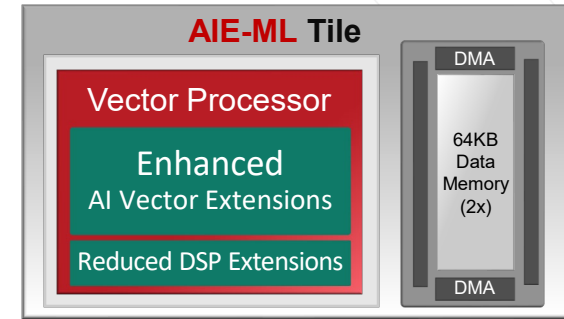
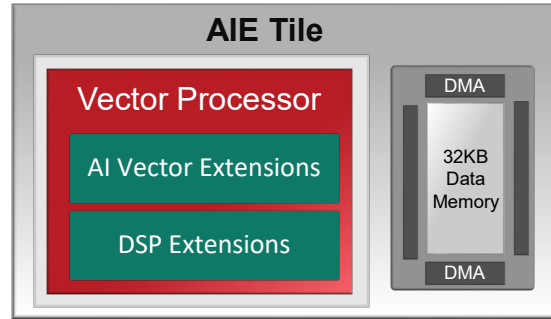
- ▶ Software programmable NoC
- ▶ Platform Management Controller
- ▶ Dedicated Interfaces (e.g., PCIe, DDR)



Technology Node: TSMC 7nm FinFET



# Xilinx AIE Tile: Architecture Features



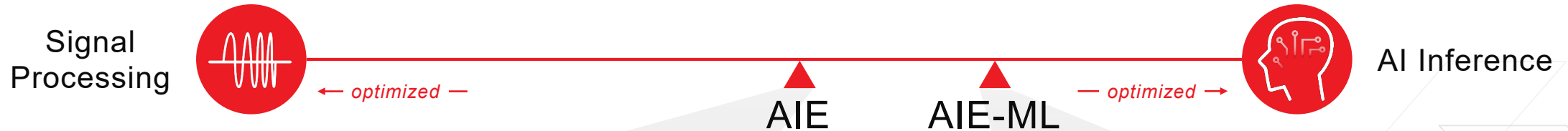
		AIE Tile		AIE-ML Tile
Target Markets		Wireless 5G, AI/ML inference, A&D		AI/ML inference
Compute (Mults / tile)	BFLOAT16	—	▶ <b>New</b> ▶	128
	INT8	128	▶ <b>2X</b> ▶	256
	INT4	—	▶ <b>New</b> ▶	512
Tile Local Data Memory		32 KB	▶ <b>2X</b> ▶	64 KB
AIE Array Interconnect B/W		1X	▶ <b>1X</b> ▶	1X
Compression and Sparsity		No		Yes
Scratchpad On-Chip Memory		PL uRAM		AIE Memory (512KB/tile)

BEAMFORMING,  
RADAR PROCESSING,  
ML INFERENCE

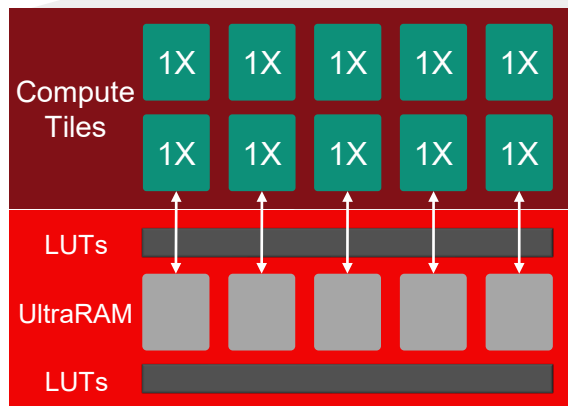
ML INFERENCE  
(CNN, RNN, MLP)



# Intelligent Engines Optimized for Any AI Application



## AI Engine Architecture

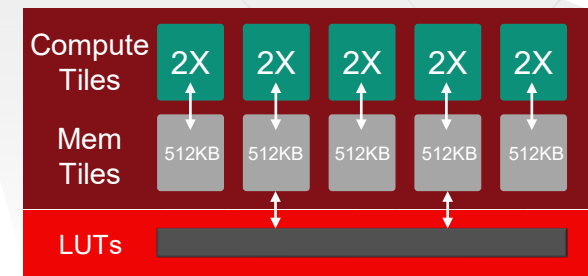


- ▶ Optimized for signal processing AND ML
- ▶ Flexibility for high performance DSP applications
- ▶ Native support for INT8, INT16, FP32

AIE	OPS / Tile	AIE-ML
256	INT4	1024
256	INT8	512
64	INT16	128
16	BFLOAT16	256
16	INT32	
16	FP32	42*
<b>KB / Tile</b>		
32	Data Memory	64
16	Program Memory	16

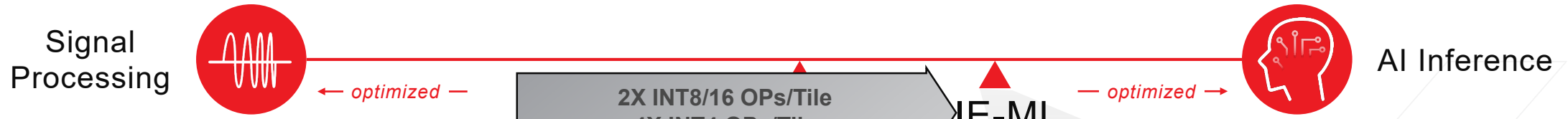
\*Via software emulation

## AIE-ML Architecture

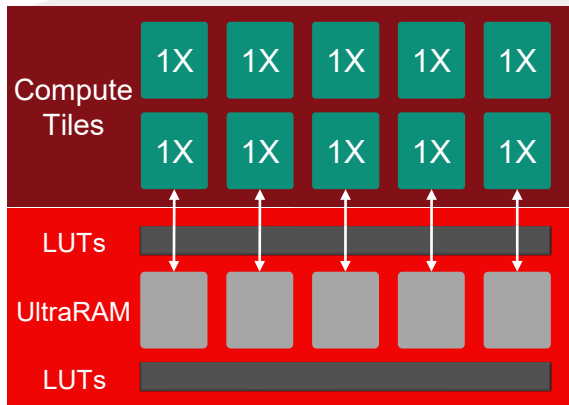


- ▶ Optimized for ML Inference Applications
- ▶ Maximum AI/ML compute with reduced footprint
- ▶ Native support for INT4, INT8, INT16, bfloat16
- ▶ Fine grained sparsity HW optimization
- ▶ Enhanced FFT & complex math support

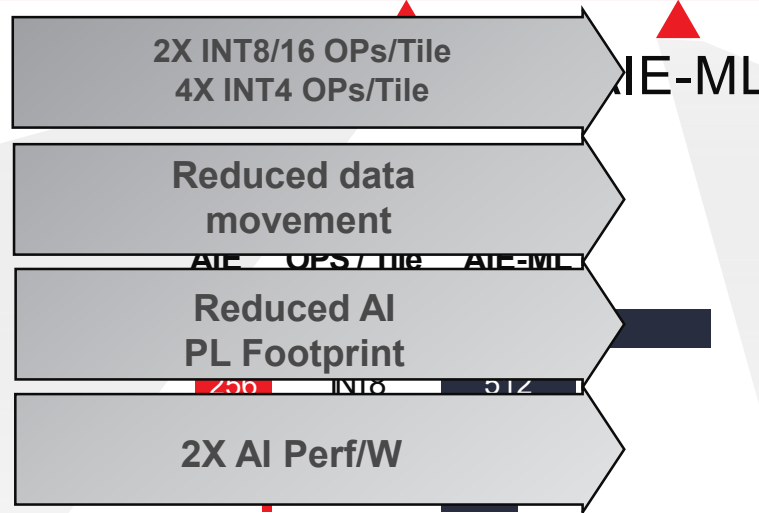
# Intelligent Engines Optimized for Any AI Application



## AI Engine Architecture



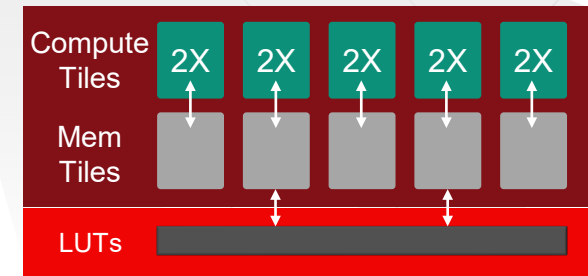
- ▶ Optimized for signal processing AND ML
- ▶ Flexibility for high performance DSP applications
- ▶ Native support for INT8, INT16, FP32



16	INT32	
16	FP32	42*
<b>KB / Tile</b>		
32	Data Memory	64
16	Program Memory	16

\*Via software emulation

## AIE-ML Architecture



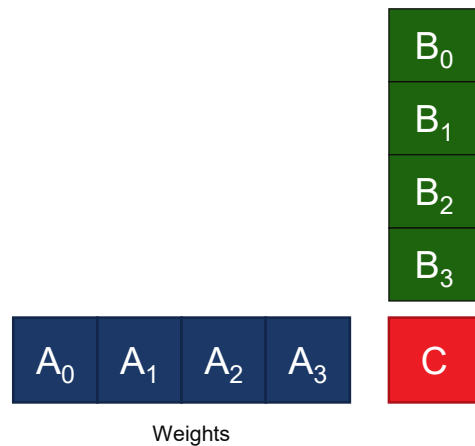
- ▶ Optimized for ML Inference Applications
- ▶ Maximum AI/ML compute with reduced footprint
- ▶ Native support for INT4, INT8, INT16, bfloat16
- ▶ Fine grained sparsity HW optimization
- ▶ Enhanced FFT & complex math support

# Flexible Dataflow: AIE-ML Array Examples

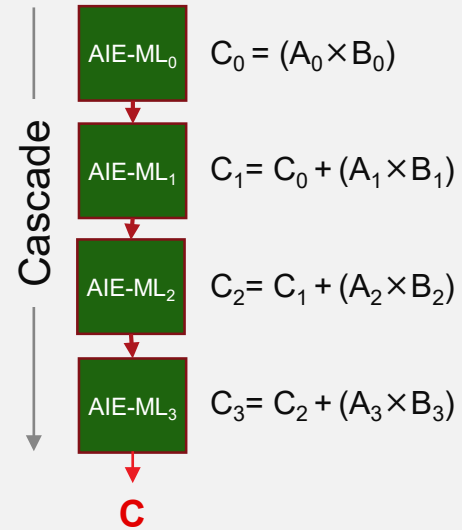
## DATA REDUCTION

$$C = (A_0 \times B_0) + (A_1 \times B_1) + (A_2 \times B_2) + (A_3 \times B_3)$$

COMPUTE SINGLE OUTPUT MATRIX

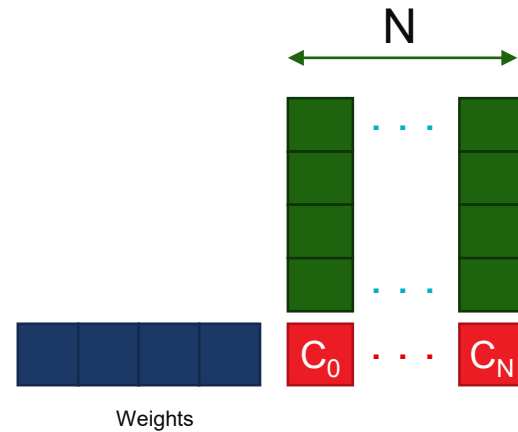


AIE-ML IMPLEMENTATION

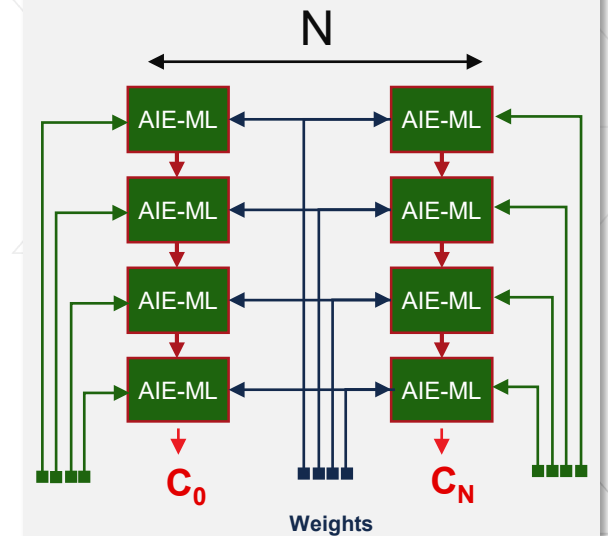


## DATA MULTICASTING

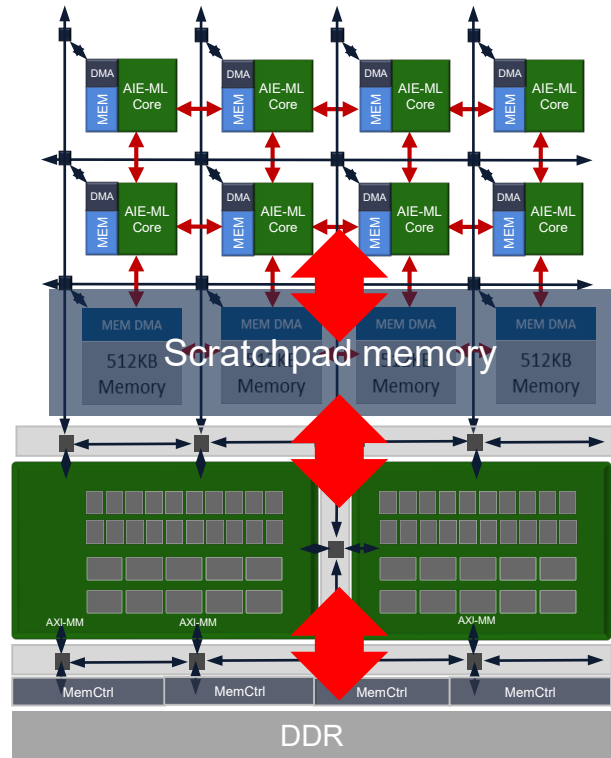
COMPUTE MULTIPLE OUTPUT MATRICES



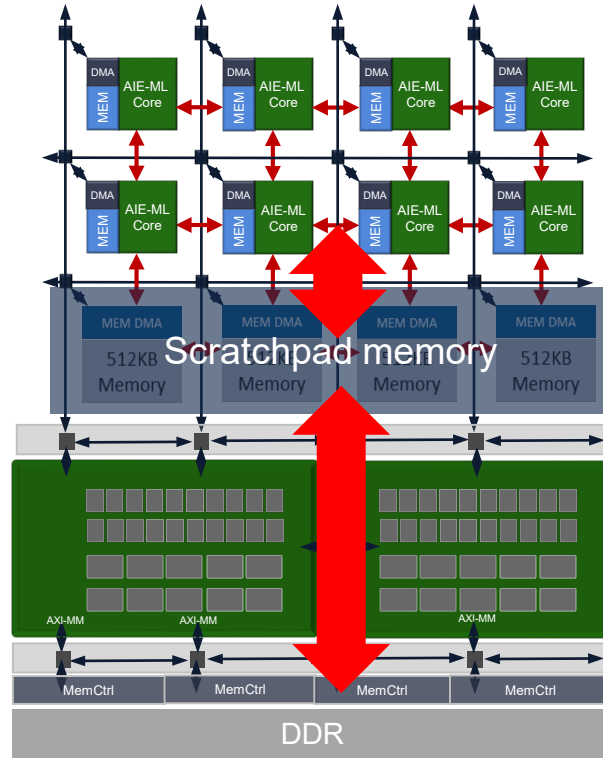
AIE-ML IMPLEMENTATION



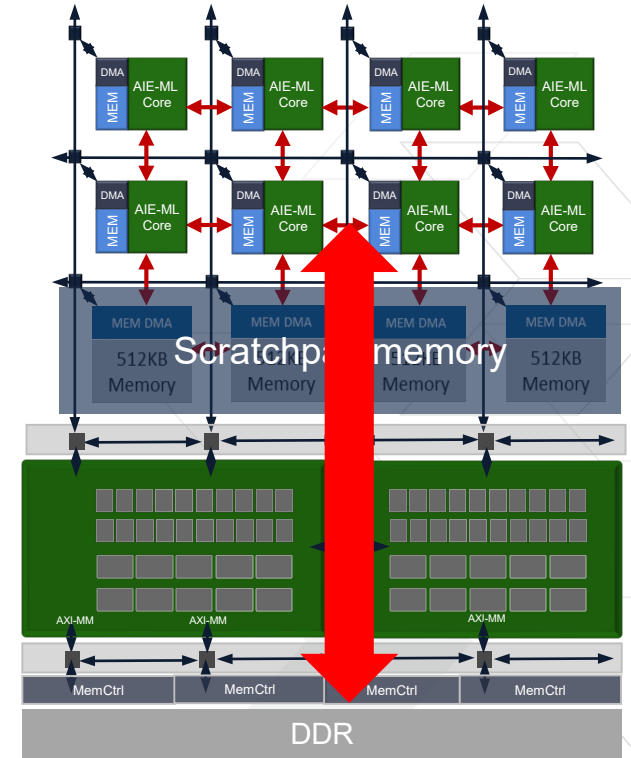
# Flexible Dataflow: Device-level Examples



**DDR ↔ Scratchpad  
using Programmable Logic**



**DDR ↔ Scratchpad  
using NoC**



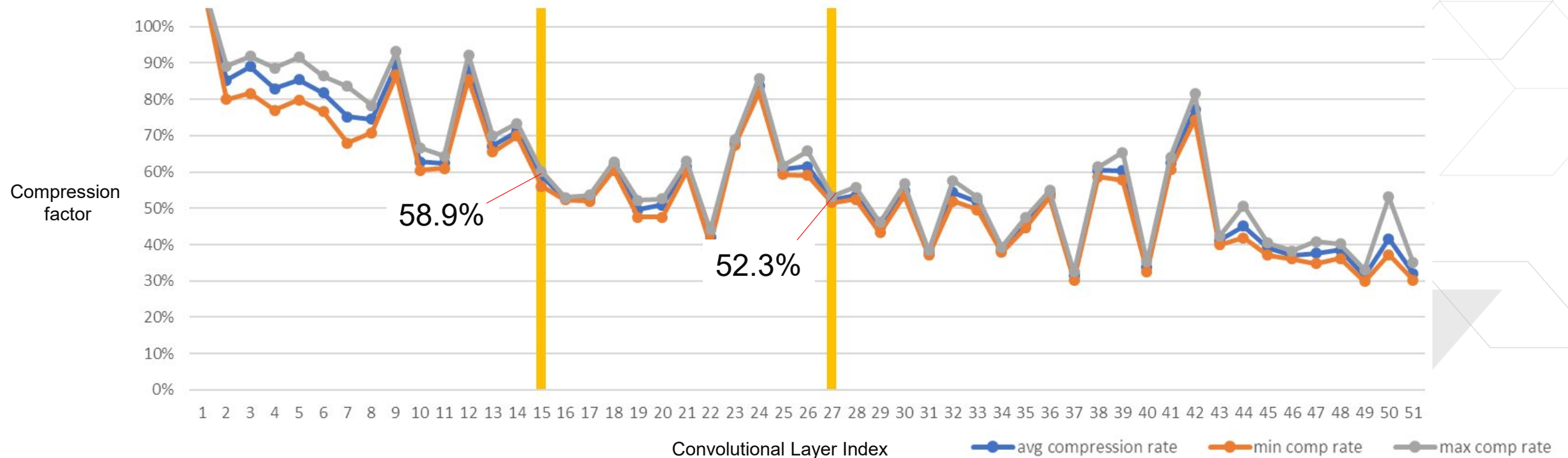
**DDR ↔ AIE-ML Tiles  
using NoC**

# Compression to DDR: Acceleration of activation Spill/Restore

## > Reduction in external memory bandwidth using compression

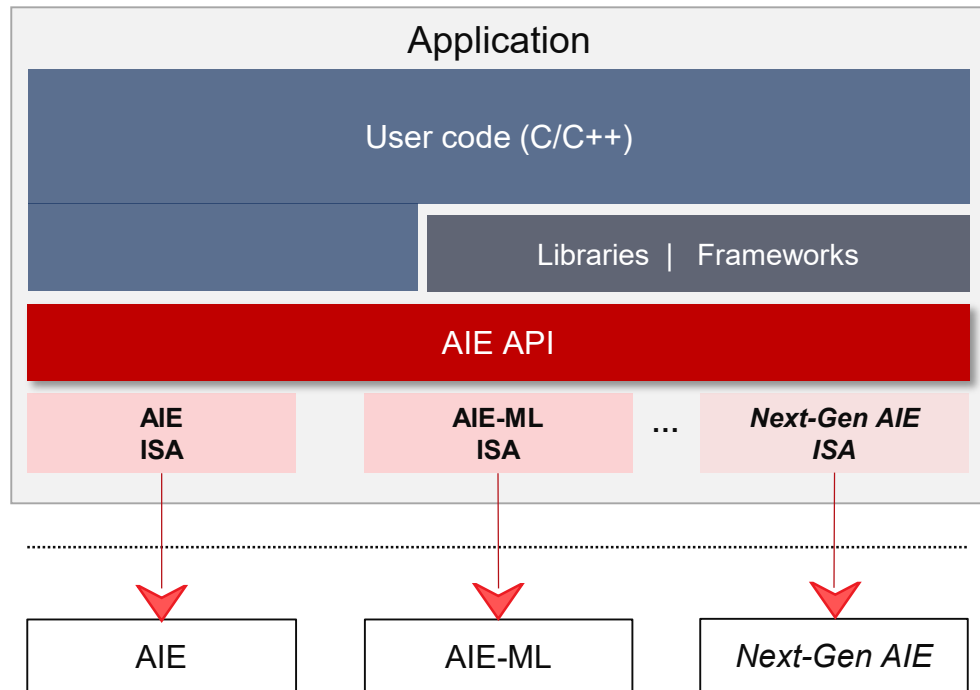
- >> Quantized int8 weights, int8 activations
- >> HD images require tiling, with spill/restore to DDR
  - Without compression ~56GB/s
  - With compression ~29GB/s

HD ResNet-50 Intermediate Feature Map Compression Factor





# AIE Architecture: Source Code Portability



- > **Backward compatible**
- > **Evolving**  
New devices can add new functionality
- > **Efficient**  
Applications can be optimized to maximize compute efficiency

Portable common high-level API across the AIE architecture

# Smart World Application: Use-Case Mapping

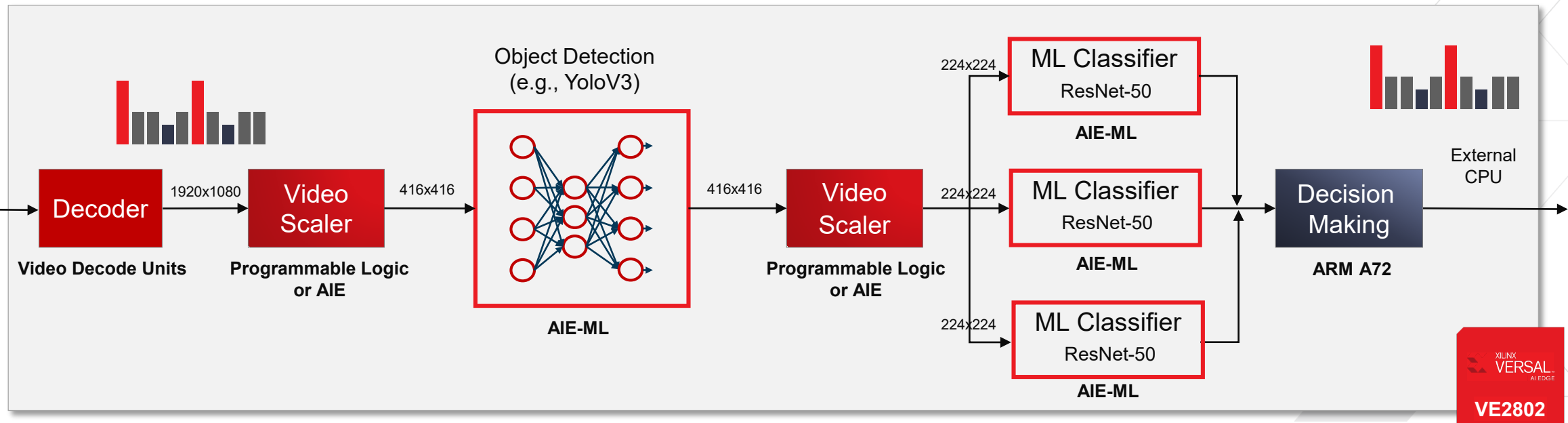
Smart Retail



Smart Hospital

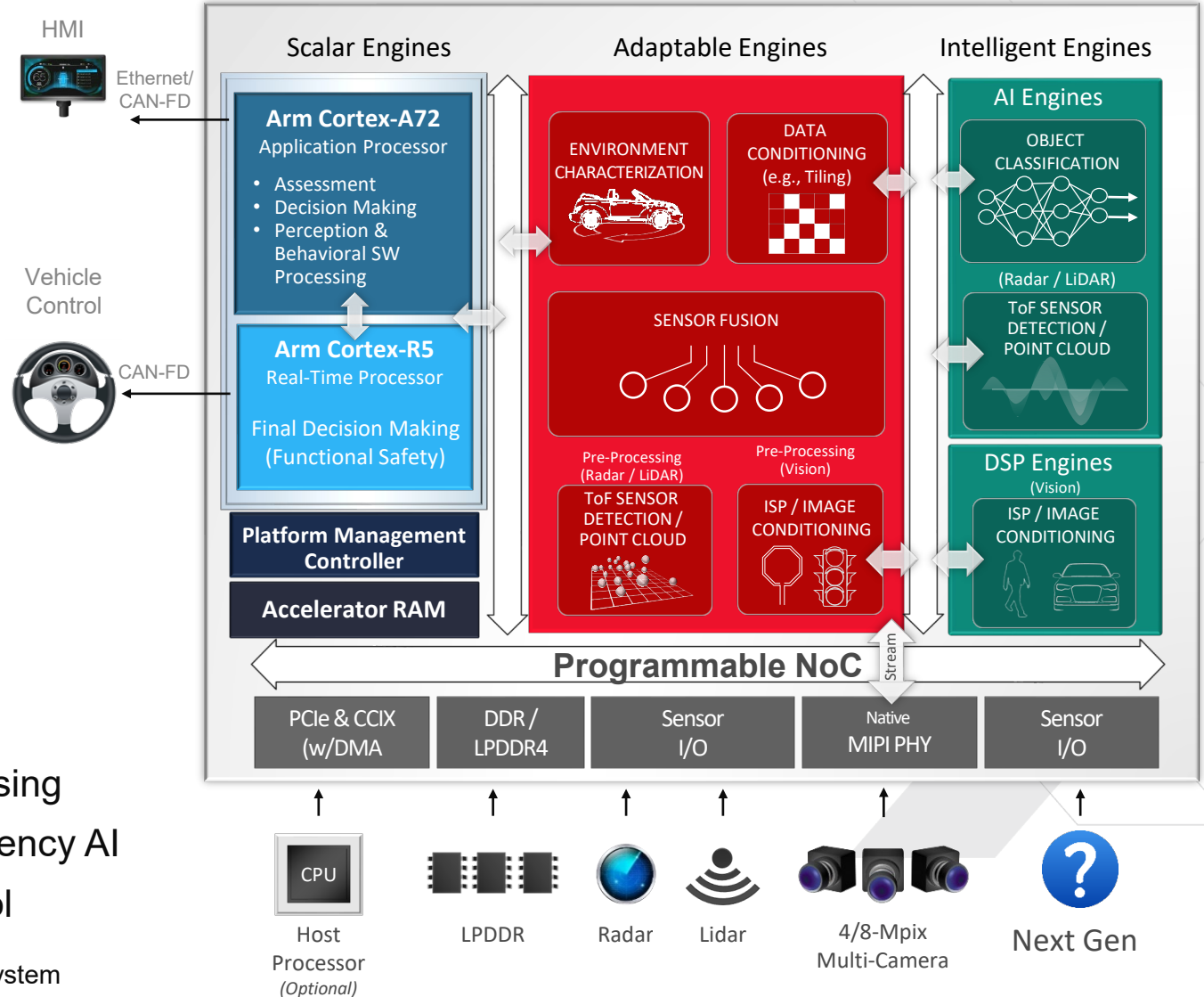
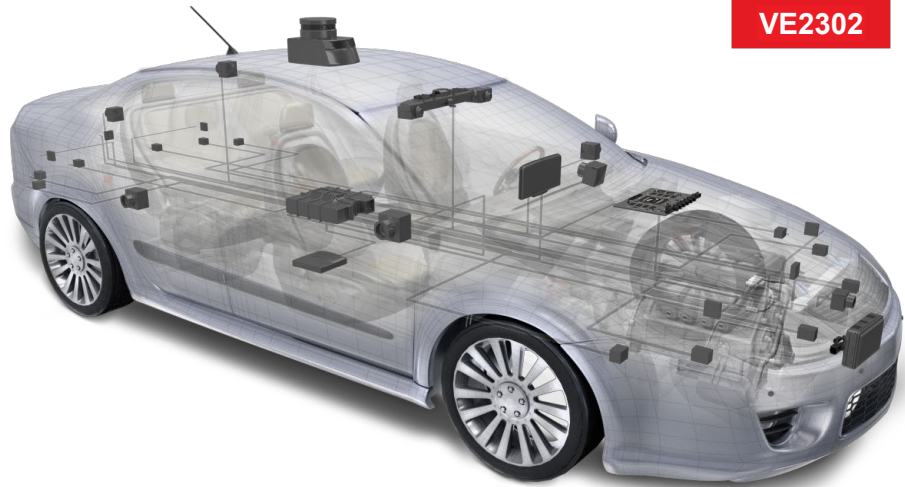


Smart City



Accelerating the Whole Application: From Video Decode to Decision Making

# Versal AI Edge in ADAS and Automated Driving

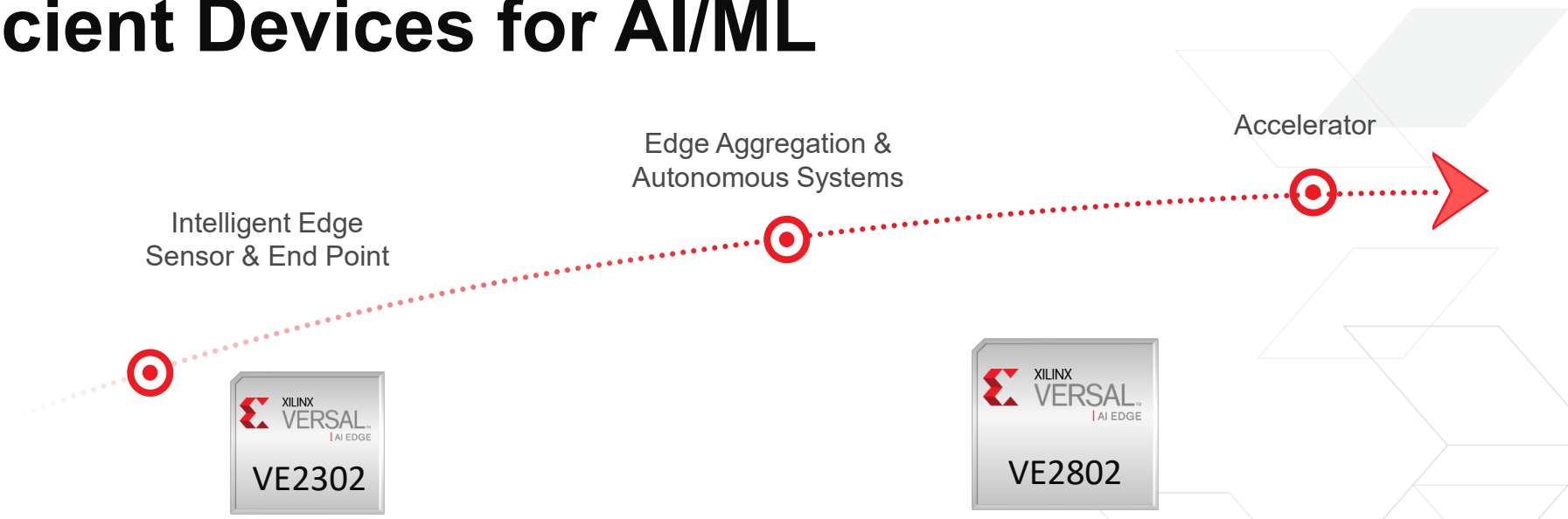


## Accelerating the Whole Application From Sensor to AI to Real-Time Control

- ▶ Adaptable Engines for sensor fusion and pre-processing
- ▶ Intelligent Engines for signal conditioning and low latency AI
- ▶ Scalar engine for decision making and vehicle control

1: Demonstrates capabilities of architecture, not representing a single chip AD system

# Range of efficient Devices for AI/ML

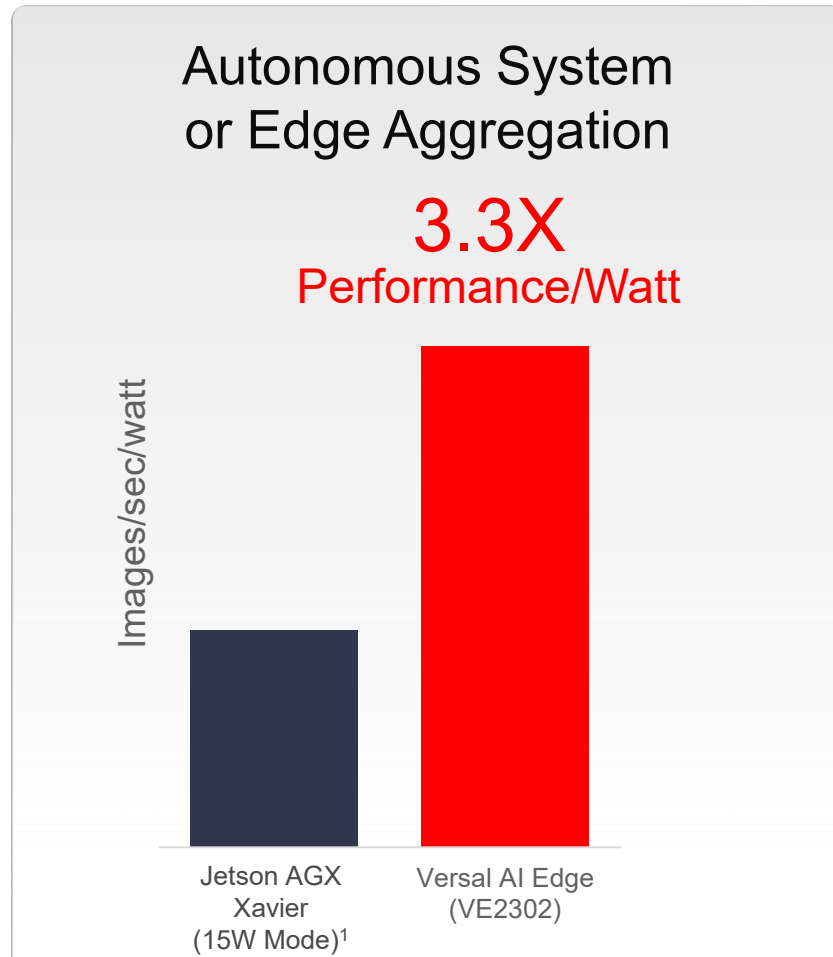


Engines	AI Compute (INT8x4) <sup>1</sup>	<b>67 TOPS</b>	<b>479 TOPS</b>
	AI Compute (INT8) <sup>1</sup>	<b>31 TOPS</b>	<b>228 TOPS</b>
	AIE-ML Tiles	34	304
	Adaptable Engines	150K LUTs	521K LUTs
	Processing Subsystem	Dual-Core Arm® Cortex®-A72 Application Processing Unit / Dual-Core Arm Cortex-R5F Real-Time Processing Unit	
RAM	Accelerator RAM (4MB)	✓	-
	Total Memory	172Mb	575Mb
	32G Transceivers	8	32
	PCIe®	✓	✓ (PCIe gen5 w/ DMA)
	Video Decode Unit (VDU)	-	✓
	Power <sup>2</sup>	<b>15-20W</b>	<b>75W</b>

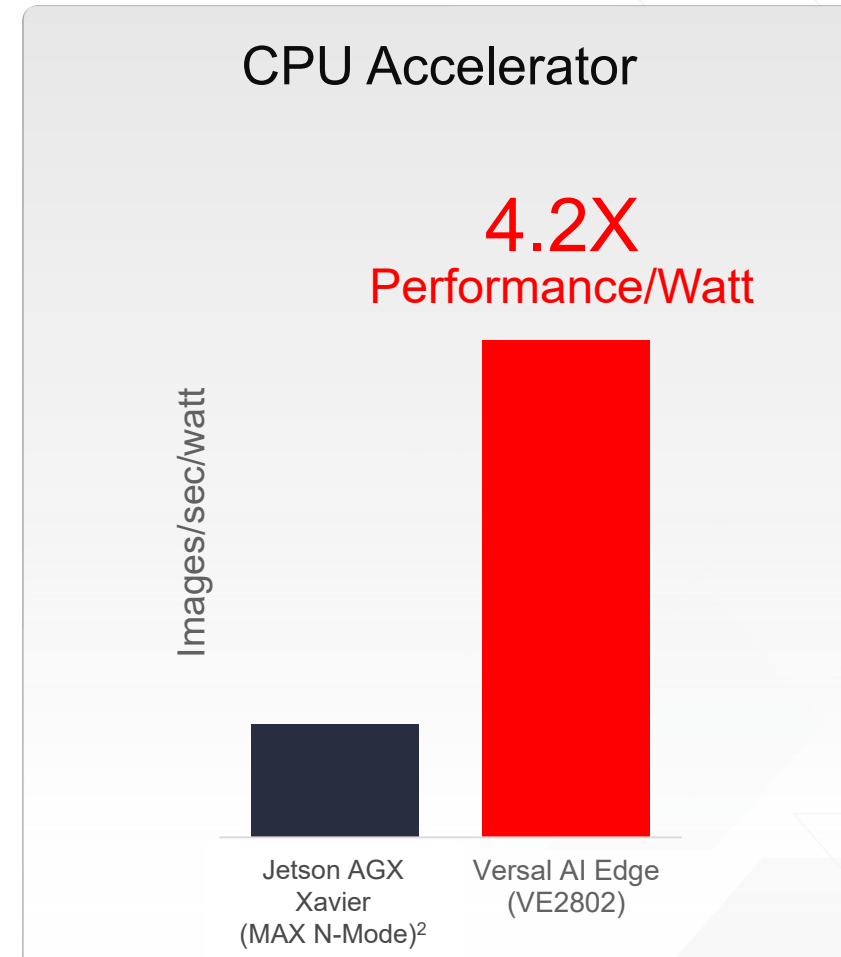
1: Total AI compute includes AI Engines, DSP Engines, and Adaptable Engines

2: Power Projections

# Performance/Watt Results



ResNet50 224x224, batch=1

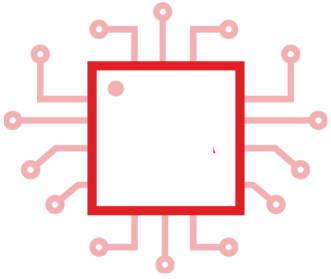


ResNet50 224x224, batch=1

1: Jetson AGX Xavier: <https://developer.nvidia.com/embedded/jetson-agx-xavier-dl-inference-benchmarks>.

2: Jetson AGX Xavier MAX N-Mode and VE2802 represent the highest performing device configuration in their respective portfolios  
Jetson AGX Xavier device power estimated by subtracting published memory and IO power from total module power

# Conclusion for 2 generations of AI Engines



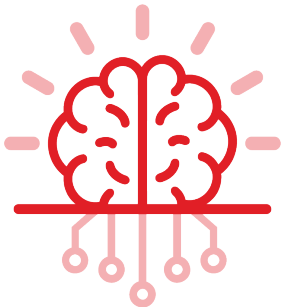
## Xilinx Optimized Devices for Wireless and AI/ML Acceleration

- ▶ Built using Xilinx 7nm Versal™ Architecture
- ▶ First Generation: Versal AI Core series with AIE
- ▶ Second Generation: Versal AI Edge series with AI-ML



## Whole Application Acceleration

- ▶ Wireless Beamforming and Digital Frontend, Video processing with ML
- ▶ Autonomous driving, Smart City for Edge Server Accelerator



## Multi-generation continuous improvements

- ▶ Major improvements in TOPs/W over multiple generations
- ▶ ML-optimized features (e.g., optimized datatypes and memory hierarchy)



---

**Thank You**

