# Software Programmable AND Hardware Adaptable: Can You Have Your Cake and Eat It, Too?

**Manuel Uhm**
**Director, Silicon Marketing - Xilinx**
**Chair of the Board of Directors - Wireless Innovation Forum**

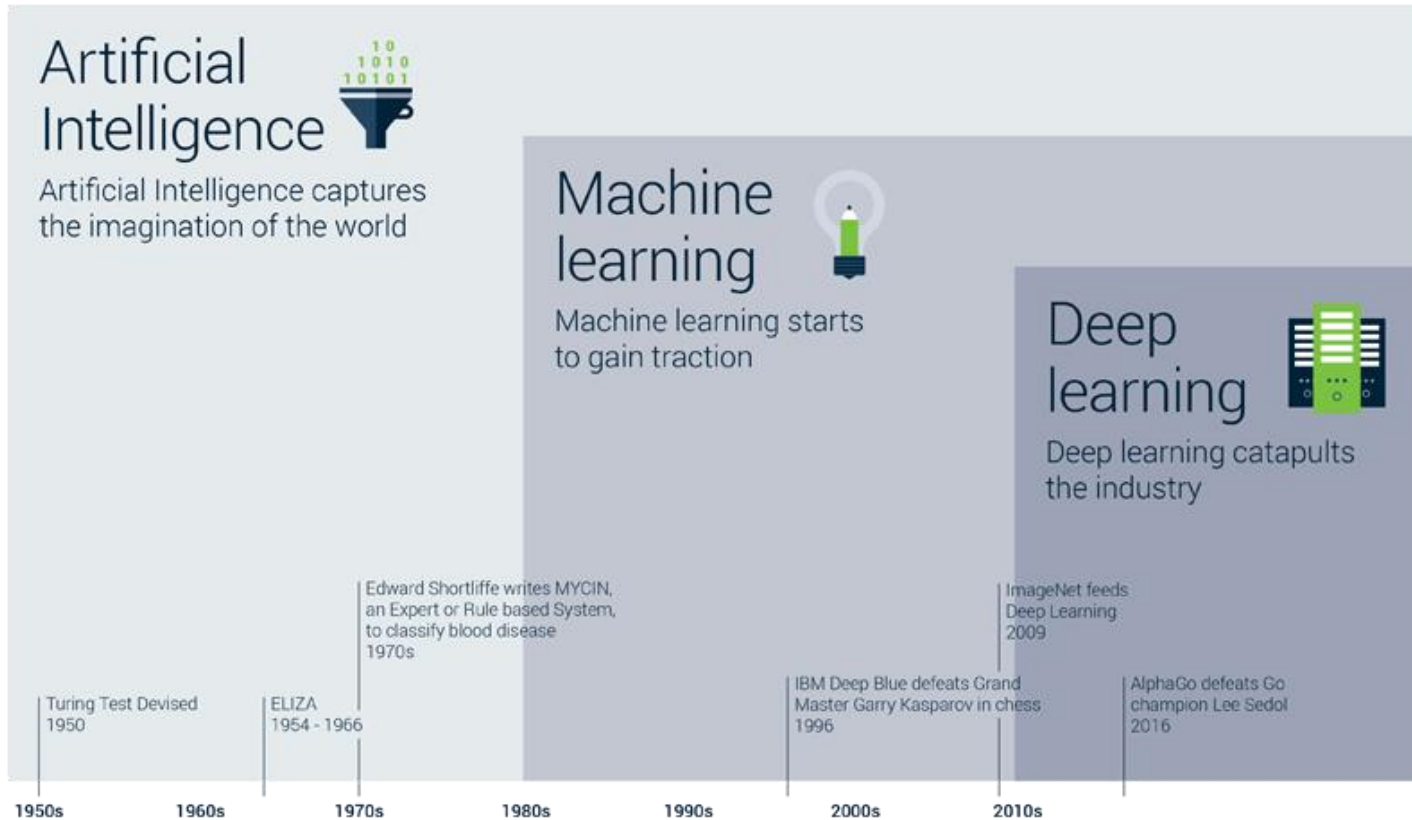**XILINX**®

# SDR Evolution



**Figure 1:** How successive generations of SDRs have come to dominate the radio industry and will continue to evolve.

Key semiconductor technology drivers:
- Moore's Law
- FPGAs
- RFICs
- Analog/Digital Integration

Source: Manuel Uhm, *Software-Defined Radio: To Infinity and Beyond*, Military Embedded Systems, October 2016

# AI Evolution



## Artificial Intelligence
Artificial Intelligence captures the imagination of the world

## Machine learning
Machine learning starts to gain traction

## Deep learning
Deep learning catapults the industry

Turing Test Devised
1950

ELIZA
1954 - 1966

Edward Shortliffe writes MYCIN, an Expert or Rule based System, to classify blood disease
1970s

IBM Deep Blue defeats Grand Master Garry Kasparov in chess
1996

ImageNet feeds Deep Learning
2009

AlphaGo defeats Go champion Lee Sedol
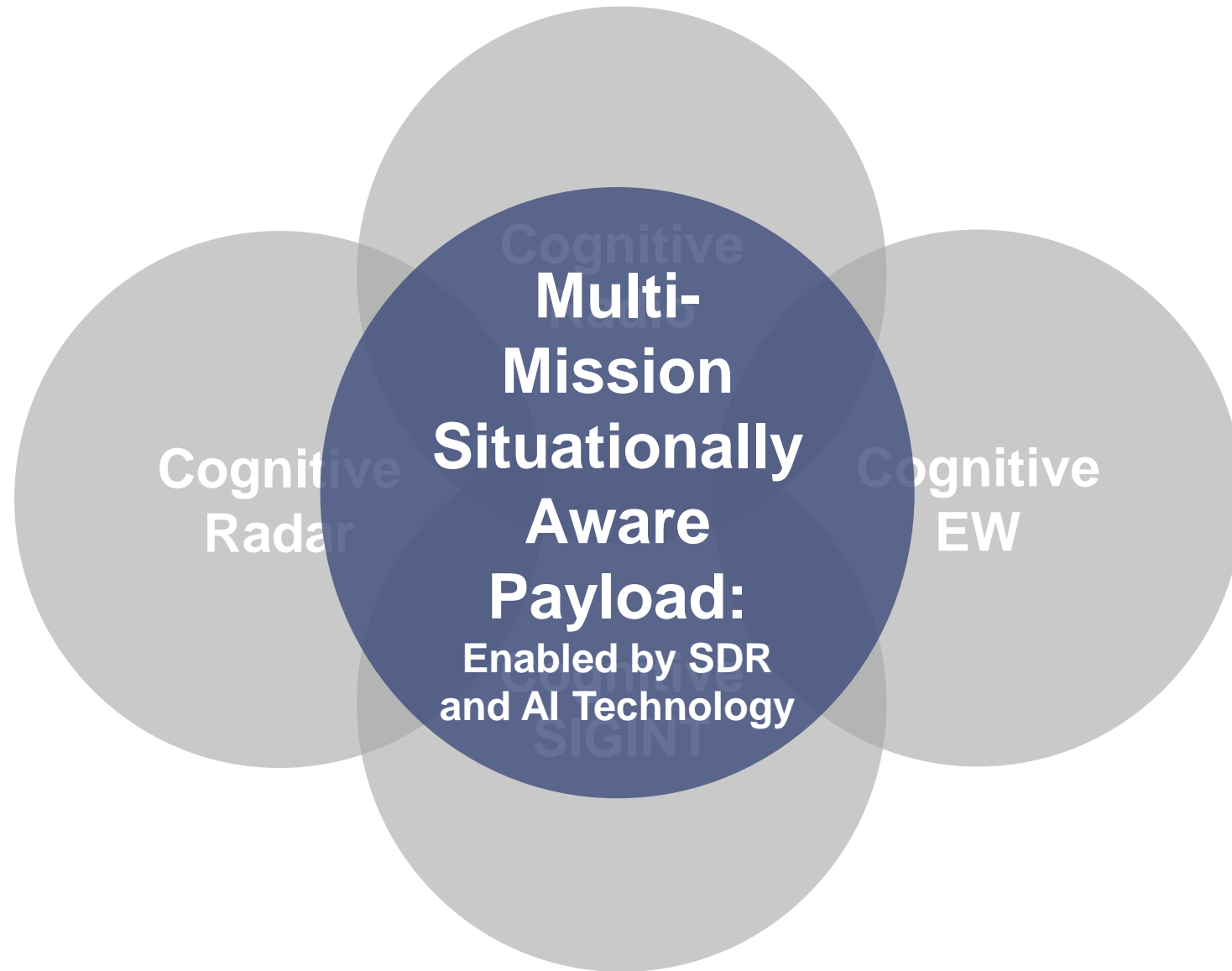2016

| 1950s | 1960s | 1970s | 1980s | 1990s | 2000s | 2010s |

**Key semiconductor technology drivers:**
- Moore's Law
- GPUs
- FPGAs
- ASICs

Source: Verhaert, *2019 Perspective on Artificial Intelligence Evolution*

XILINX

# SDR & AI Payload Convergence

Cognitive Radar

Cognitive SIGINT

Cognitive EW

**Multi-Mission Situationally Aware Payload:**
Enabled by SDR and AI Technology

**XILINX**

# End of the Line for Processor Performance?

**DENNARD SCALING**

Power Density Rises

**MOORE'S LAW**

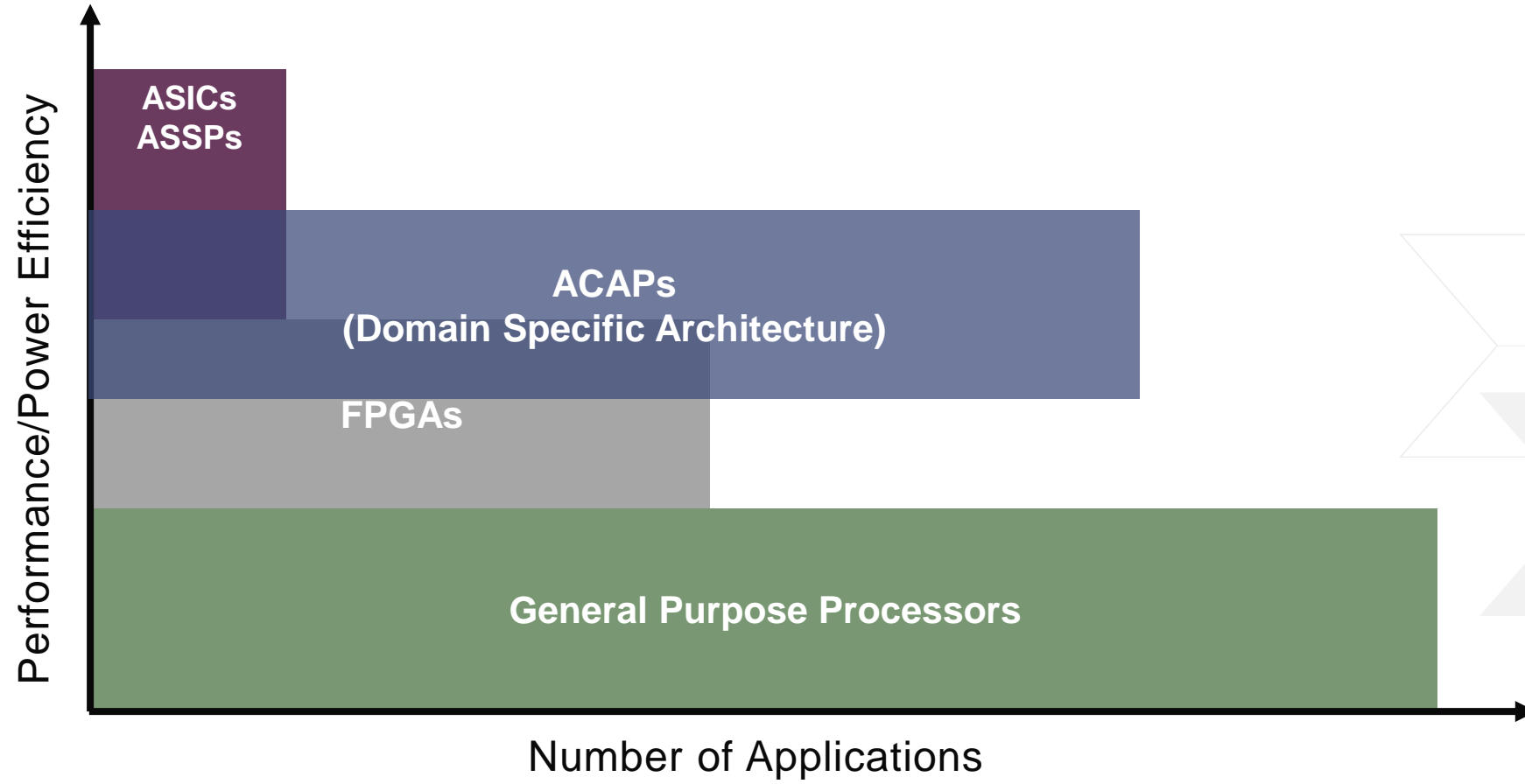End of "PPA" Improvement

**AMDAHL'S LAW**

Multicore Hits Limit



40 Years of Processor Performance

Performance vs. VA11-780

- CISC 2x / 3.5yrs (22%/yr)
- RISC 2x / 1.5yrs (52%/yr)
- End of Dennard Scaling Multicore 2x / 3.5yrs (23%/yr)
- Amdahl's Law 2x / 6yrs (12%/yr)
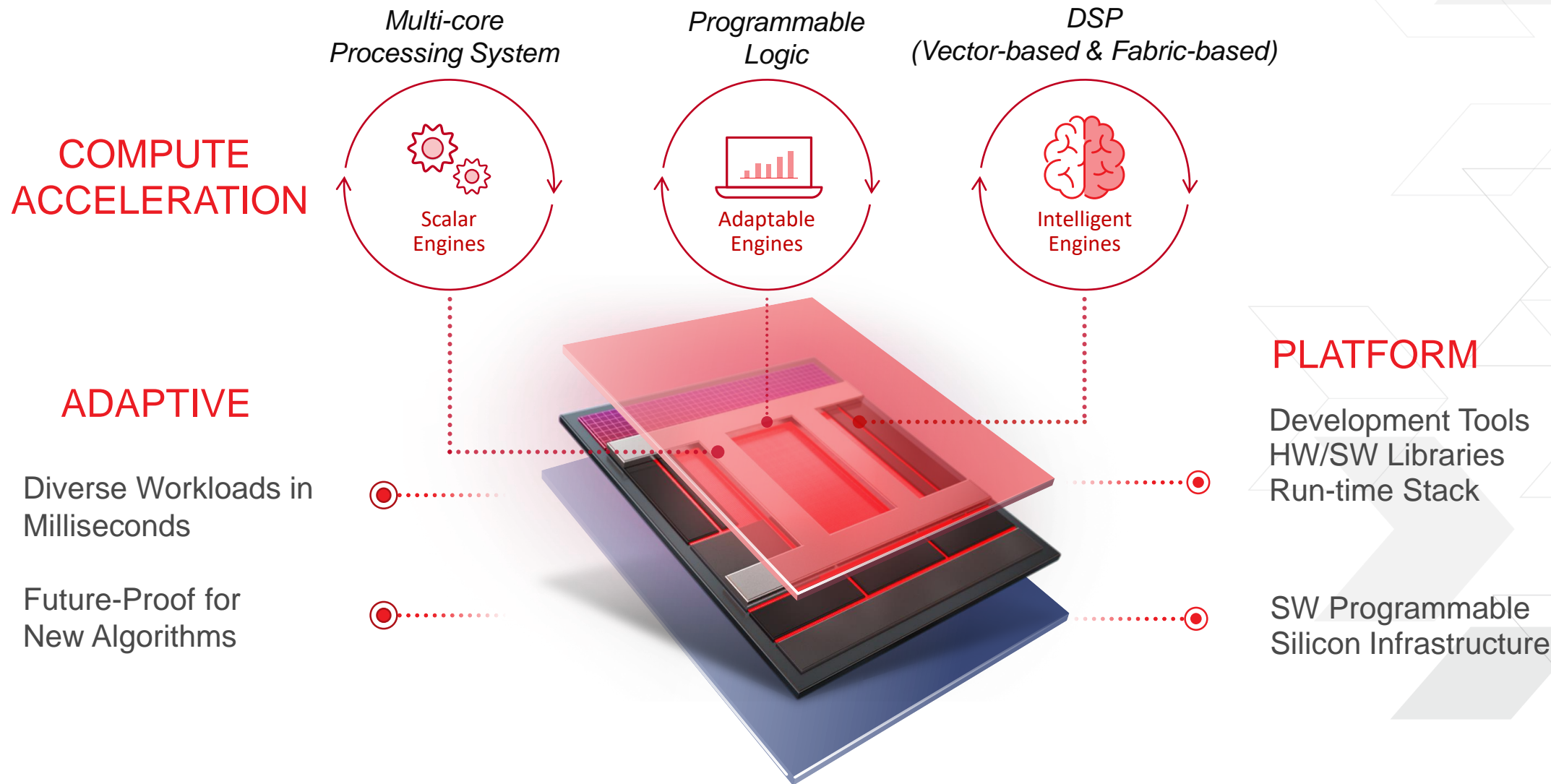- End of the line? 2x / 20yrs (3%/yr)

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

## Moving Forward: Domain-Specific Architectures (DSAs)

XILINX

# Evolving Processor Landscape

# The Adaptive Compute Acceleration Platform



**Multi-core Processing System**

Scalar Engines

**Programmable Logic**

Adaptable Engines

**DSP (Vector-based & Fabric-based)**

Intelligent Engines

COMPUTE ACCELERATION

ADAPTIVE

Diverse Workloads in Milliseconds

Future-Proof for New Algorithms

PLATFORM

Development Tools
HW/SW Libraries
Run-time Stack

SW Programmable
Silicon Infrastructure

Enabling Data Scientists, SW Developers, HW Developers

XILINX

# It's About the Platform



## ADAPTIVE

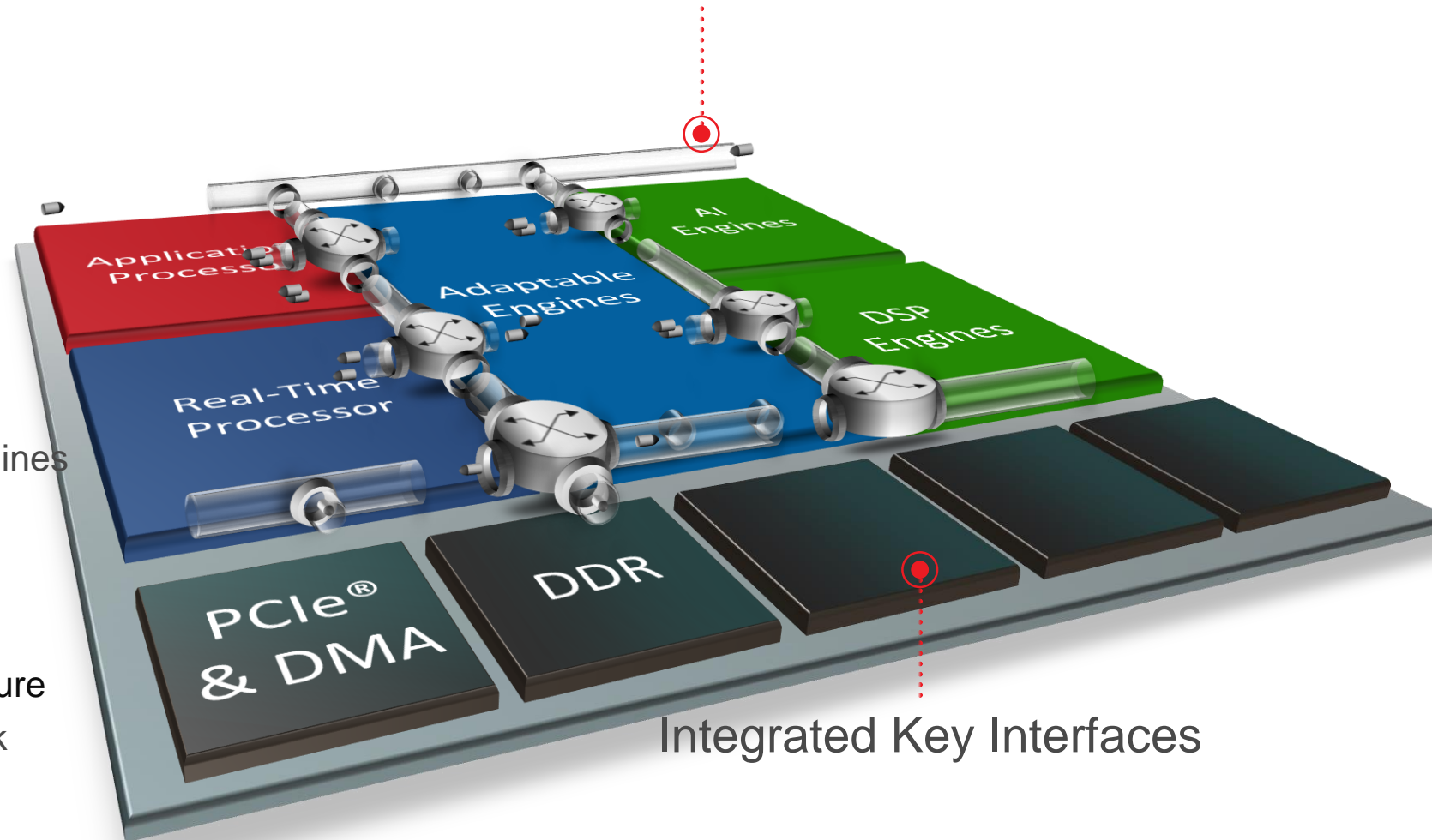> Adaptable to Diverse Workloads
> Future-proof Algorithms

## COMPUTE ACCELERATION

> Heterogeneous architecture
> Scalar, Adaptable, and Intelligent Engines

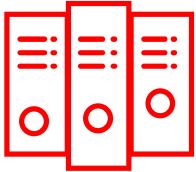## PLATFORM

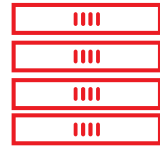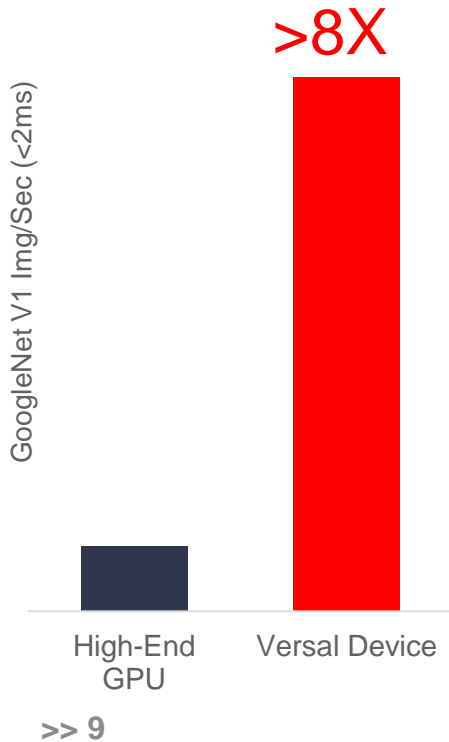> SW Programmable Silicon Infrastructure
> Integrates with Dev Tools & SW Stack

Programmable Network on Chip
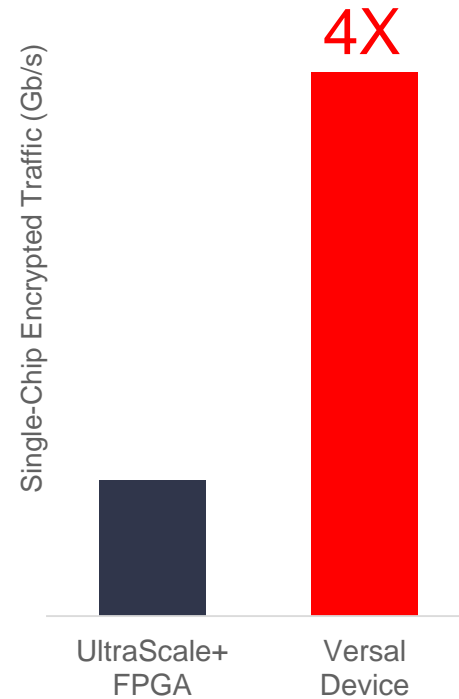
Integrated Key Interfaces

XILINX

# Breakthrough Performance for Cloud, Network, and Edge

**Cloud Compute**
Breakthrough AI Inference

**Networking**
Multi-terabit Throughput

**5G Wireless**
Compute for Massive MIMO

**Edge Compute**
AI Inference at Low Power



>8X

GoogleNet V1 Img/Sec (<2ms)

High-End GPU | Versal Device

4X

Single-Chip Encrypted Traffic (Gb/s)

UltraScale+ FPGA | Versal Device

5X

Int 16x16 DSP Compute (TeraMAC/ sec)

UltraScale+ RFSoC | Versal Device

15X

ResNet50 img/sec (batch=1)

UltraScale+ MPSoC | Versal Device

>> 9

XILINX CONFIDENTIAL

**&#10157; XILINX.**

# Hardware: The Foundation for the Software Stack

User Application
C, C++, Python

Application-Specific Frameworks

OS • Runtime • Middleware

SW Dev Tools • Vivado® Design Suite

| Embedded SW | Libraries or RTL | C Libraries |
|---|---|---|
| Scalar Engines | Adaptable Engines | Intelligent Engines |

Versal™ Platform
(PCIe®, DMA, DDR, Programmable NoC)

User Application

Application-Specific Frameworks

Runtime

SW Dev Tools

XILINX

# Versal ACAP: A Platform for SW *and* HW Developers

## Fully Software Programmable
with Hardware Design Path

| | |
|---|---|
| User Application<br>C, C++, Python | |
| Frameworks | Software Platform |
| Runtime | |
| OS • Drivers | |
| IP • Libraries | Hardware Platform |
| **Evaluation & Deployment Boards** | |
| **Versal™ ACAP Device & Integrated Shell** | |

Vitis™ Unified Software Platform

Vivado® Design Suite

# Possible Platform Example: Multi-Mission Situationally Aware UAV Payload with Versal ACAP



**UAV Platform**

| Multi-Mission Applications: Comms, Radar, SIGINT, EW |
|:---:|

| Frameworks<br>TensorFlow          PYTORCH |
|:---:|

| Xilinx Runtime (XRT) |
|:---:|

| VxWorks | INTEGRITY SECURITY SERVICES | QNX SOFTWARE SYSTEMS | xfopenCV | DSPlib | ML Overlay |
|:---:|:---:|:---:|:---:|:---:|:---:|

| Scalar Engines | Adaptable Engines | AI Engines |
|:---:|:---:|:---:|

| Versal ACAP Eval Board |
|:---:|

| VERSAL ACAP |
|:---:|

XILINX

# Key to ACAP is Mapping to Embedded IP & Accelerators
*Application Example: Video Security (Surveillance)*

**FPGA**



- 100G MAC
- 100G MAC
- Video
  - Scaling
  - Compression
- Video Transcoding
- AI Inference
- Host Interface (PCIe®)
- Enqueu/Dequeu

**ACAP**



Scalar Engines | Adaptable Engines | Intelligent Engines

- Application Processor
- Real-Time Processor
- SCALING
- COMPRESSION
- ACCELERATOR (KERNEL)
- ELERATOR (KERNEL)
- ACCELERATOR (KERNEL)
- ELERATOR (KERNEL)
- MACHINE LEARNING
- VIDEO TRANSCODING
- Programmable Network on Chip
- PCIe & CCIX (w/DMA)
- DDR4
- SerDes
- Multirate Ethernet
- MIPI
- LVDS
- GPIO

> All workloads are in logic, no room for additional differentiation
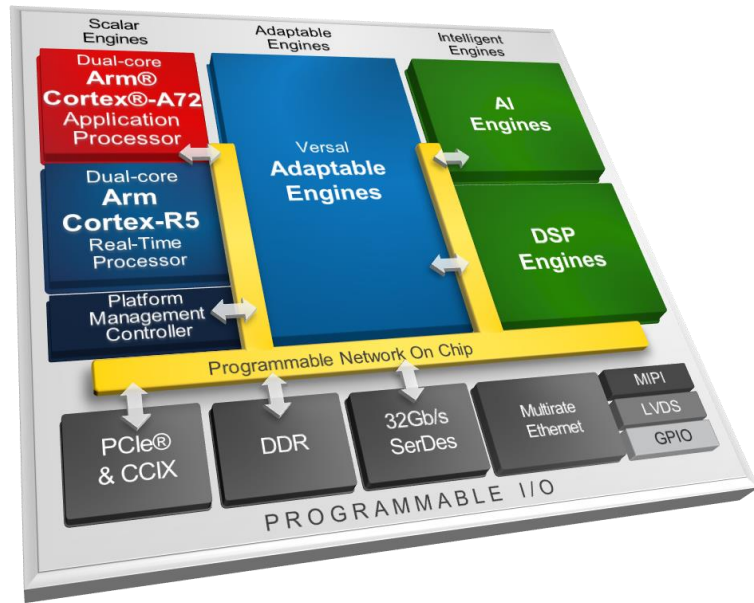
> Need for the arbitration between workloads and memory

> Workloads are mapped to the right engines

> High bandwidth, guaranteed QoS via Programmable NoC

> Power efficiency and greater performance

XILINX

# Part of a Comprehensive Product Portfolio

| | 28nm | 20nm | 16nm | 7nm |
|---|---|---|---|---|
| **Cloud Applications** | | | | XILINX VERSAL™ AI Core |
| **RF Applications** | ZYNQ® Zynq-7000 | | ZYNQ® RFSoC | XILINX VERSAL™ AI RF |
| **Edge Applications** | ZYNQ® Zynq-7000S | | ZYNQ® UltraSCALE+ | XILINX VERSAL™ AI Edge |
| **Broad Application** | VIRTEX®7 KINTEX®7 ARTIX®7 SPARTAN®7 | VIRTEX® UltraSCALE KINTEX® UltraSCALE | VIRTEX® UltraSCALE+ HBM VIRTEX® UltraSCALE+ KINTEX® UltraSCALE+ | XILINX VERSAL™ HBM XILINX VERSAL™ Premium XILINX VERSAL™ Prime |

XILINX

# Versal™ AI Core Series



Scalar Engines · Adaptable Engines · Intelligent Engines

Dual-core Arm® Cortex®-A72 Application Processor
Versal Adaptable Engines
AI Engines
Dual-core Arm Cortex-R5 Real-Time Processor
DSP Engines
Platform Management Controller
Programmable Network On Chip
PCIe® & CCIX
DDR
32Gb/s SerDes
Multirate Ethernet
MIPI
LVDS
GPIO
PROGRAMMABLE I/O

## Breakthrough AI Inference Throughput

> Portfolio's highest throughput for low latency inference

> Optimized for cloud, networking, and autonomous applications

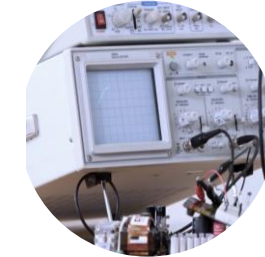> For highest dynamic range of AI and workload acceleration

Data Center Compute
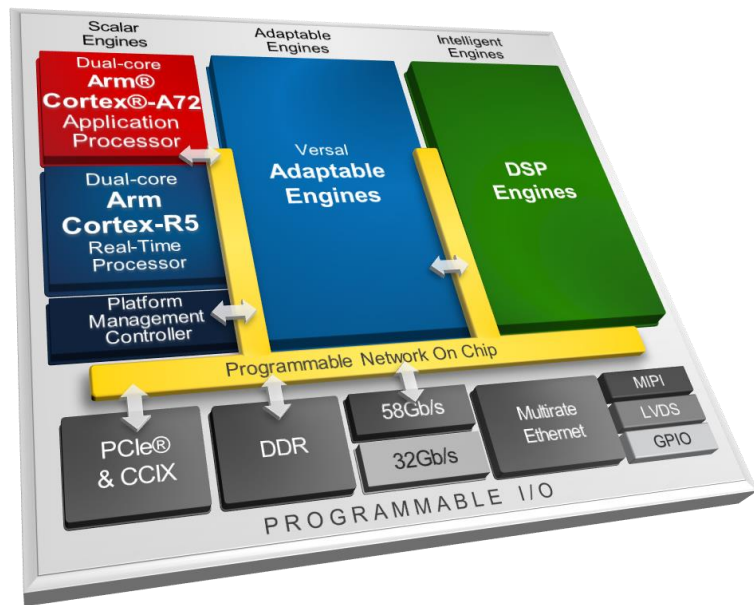
5G Radio & Beamforming

ADAS, AD Prototyping
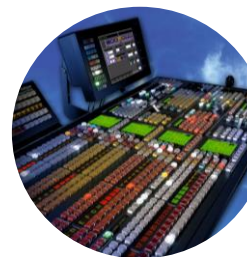
Cable Access
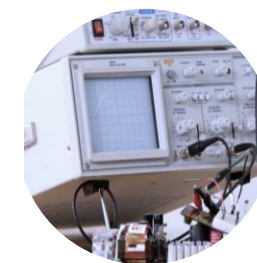
A&D

Wireless Test Equipment

# Versal Prime Series



**Broad Applicability across Multiple Markets**

> Mid-range series in the Versal™ portfolio

> Optimized for connectivity

> For inline acceleration and diverse workloads

Nx100G Ethernet & OTN Networking
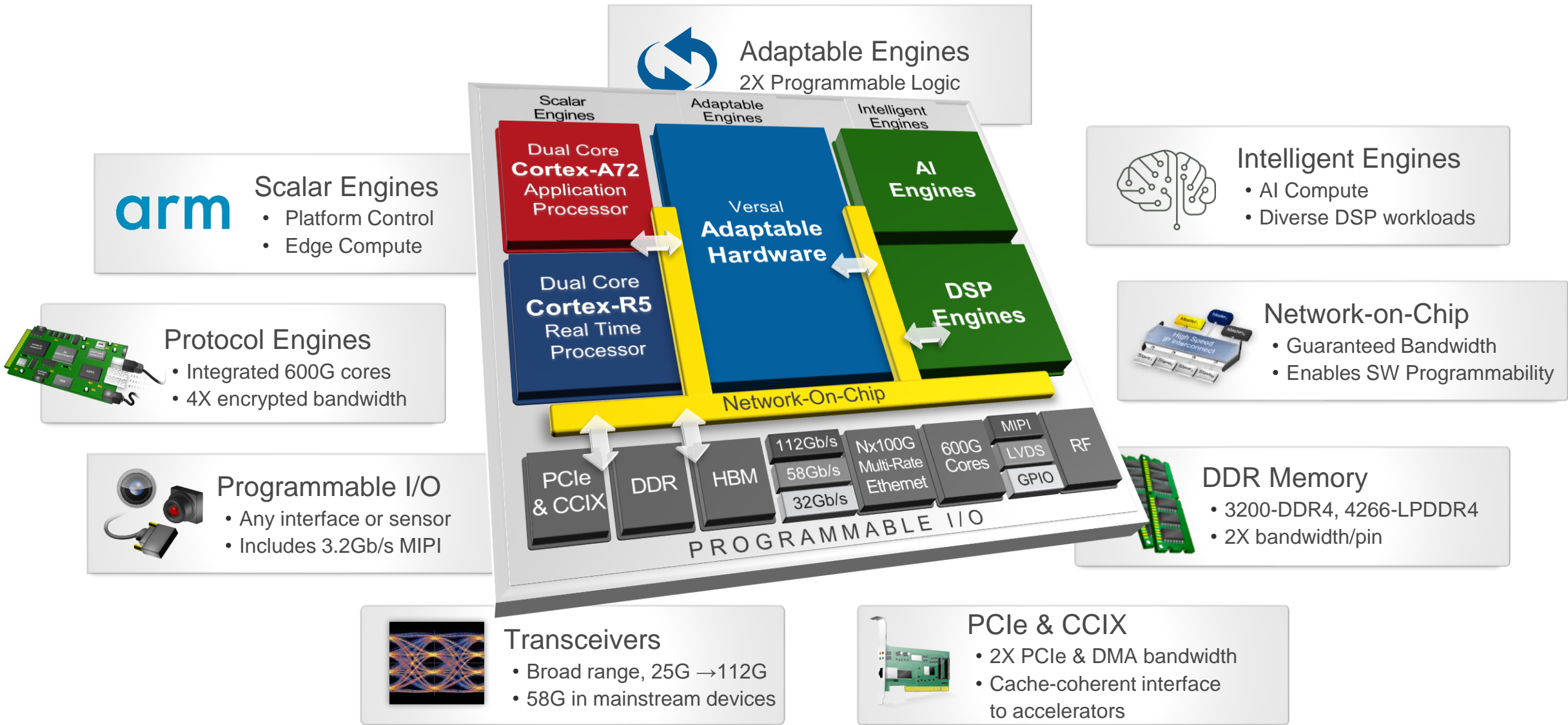
Data Center Network & Storage

Avionics Control

Broadcast Switches

Network Test Equipment

Medical Imaging

XILINX

# Versal Architecture Overview



**Scalar Engines**
- Platform Control
- Edge Compute

**Protocol Engines**
- Integrated 600G cores
- 4X encrypted bandwidth

**Programmable I/O**
- Any interface or sensor
- Includes 3.2Gb/s MIPI

**Adaptable Engines**
2X Programmable Logic

**Intelligent Engines**
- AI Compute
- Diverse DSP workloads

**Network-on-Chip**
- Guaranteed Bandwidth
- Enables SW Programmability

**DDR Memory**
- 3200-DDR4, 4266-LPDDR4
- 2X bandwidth/pin

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

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

Within diagram:
Scalar Engines — Dual Core Cortex-A72 Application Processor; Dual Core Cortex-R5 Real Time Processor
Adaptable Engines — Versal Adaptable Hardware
Intelligent Engines — AI Engines; DSP Engines
Network-On-Chip
PROGRAMMABLE I/O — PCIe & CCIX, DDR, HBM, 112Gb/s, 58Gb/s, 32Gb/s, Nx100G Multi-Rate Ethernet, 600G Cores, MIPI, LVDS, GPIO, RF

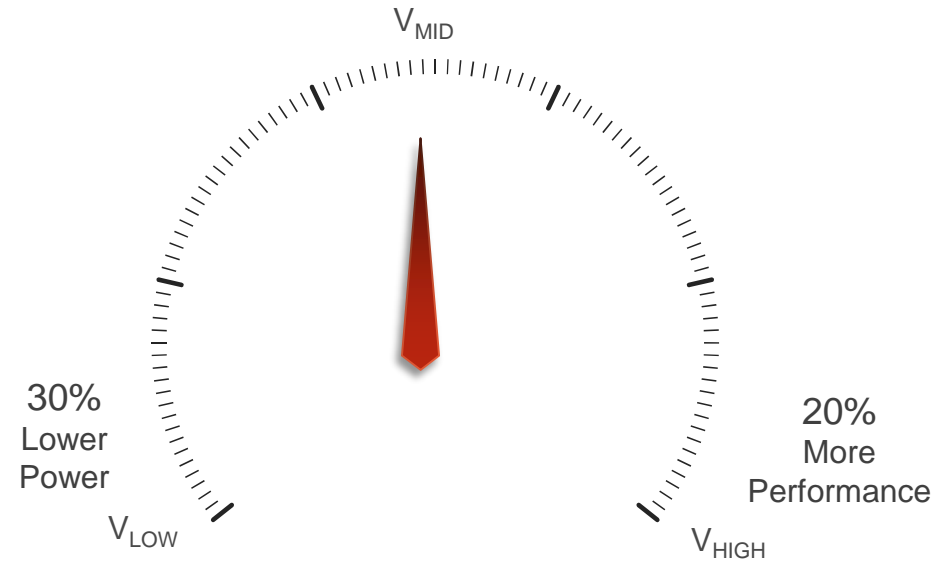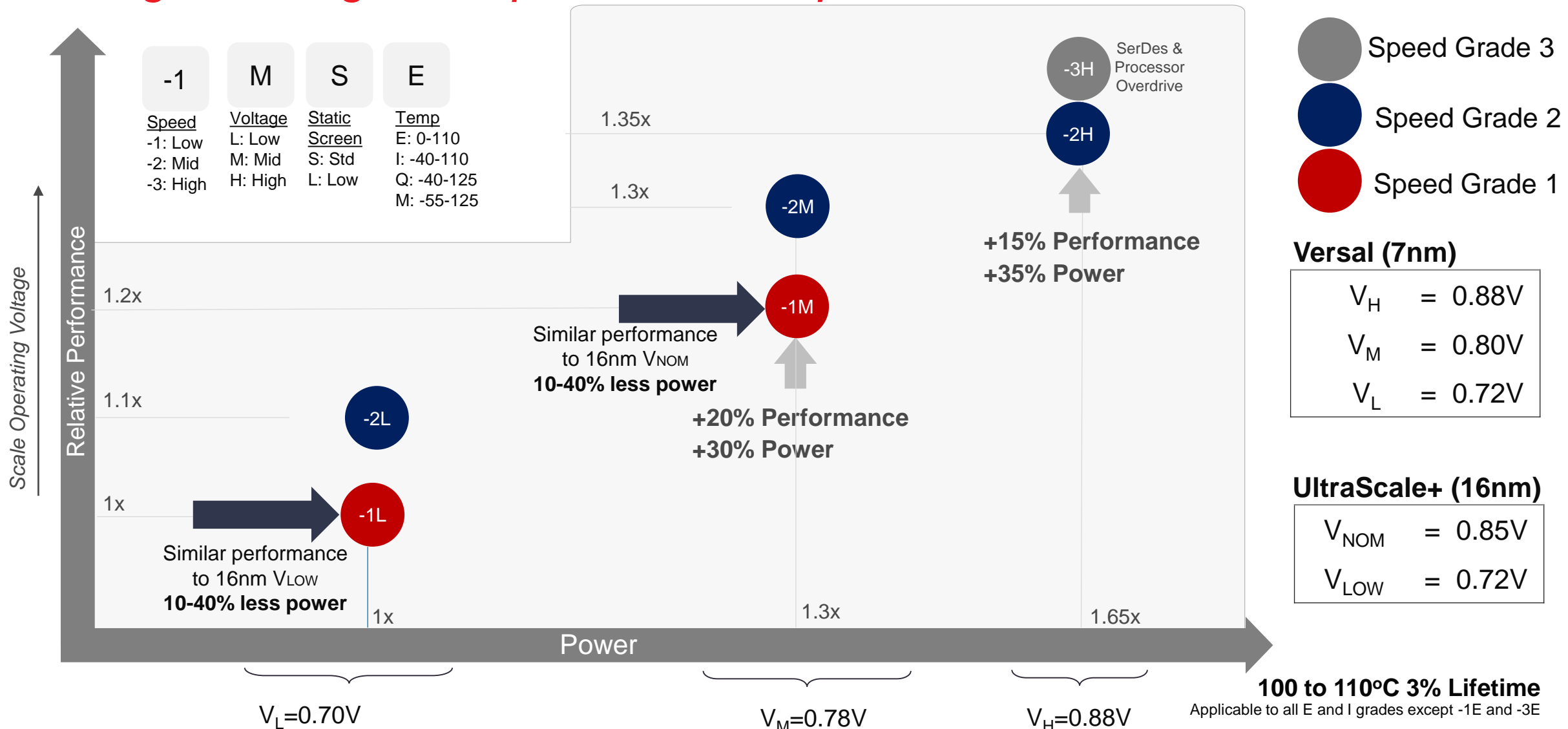XILINX

# Tune for Power & Performance in Versal ACAP

> Three operating voltages to choose from

> Balance power/performance for target app

> Equivalent to 3 speed grades in one device

| -1 | M | S | E |
|---|---|---|---|

| Speed | Voltage | Static Screen | Temp |
|---|---|---|---|
| -1: Low | L: Low | S: Std | E: 0-110 |
| -2: Mid | M: Mid | L: Low | I: -40-110 |
| -3: High | H: High | | Q: -40-125 |
| | | | M: -55-125 |

$V_{MID}$

30%
Lower
Power

20%
More
Performance

$V_{LOW}$     $V_{HIGH}$

XILINX

# Granular Control of Power vs. Performance
## *Voltage Scaling with Speed Grade Options*



**Legend boxes:**

| -1 | M | S | E |
|---|---|---|---|
| Speed | Voltage | Static Screen | Temp |
| -1: Low | L: Low | S: Std | E: 0-110 |
| -2: Mid | M: Mid | L: Low | I: -40-110 |
| -3: High | H: High | | Q: -40-125 |
| | | | M: -55-125 |

Speed Grade 3
Speed Grade 2
Speed Grade 1

-3H  SerDes & Processor Overdrive

-2H

**+15% Performance**
**+35% Power**

-2M

-1M

Similar performance to 16nm V$_{NOM}$
**10-40% less power**

**+20% Performance**
**+30% Power**

-2L

-1L

Similar performance to 16nm V$_{LOW}$
**10-40% less power**

Relative Performance / Scale Operating Voltage (y-axis): 1x, 1.1x, 1.2x, 1.3x, 1.35x

Power (x-axis): 1x, 1.3x, 1.65x

V$_L$=0.70V

V$_M$=0.78V

V$_H$=0.88V

**Versal (7nm)**

| V$_H$ | = 0.88V |
|---|---|
| V$_M$ | = 0.80V |
| V$_L$ | = 0.72V |

**UltraScale+ (16nm)**

| V$_{NOM}$ | = 0.85V |
|---|---|
| V$_{LOW}$ | = 0.72V |

**100 to 110°C 3% Lifetime**
Applicable to all E and I grades except -1E and -3E

XILINX

# New Intelligent Engines

*Massive AI Inference Throughput and Wireless Compute*

## Up to 1.3GHz VLIW / SIMD vector processors

> Versatile core for ML and other advanced DSP workloads

## Massive array of interconnected cores

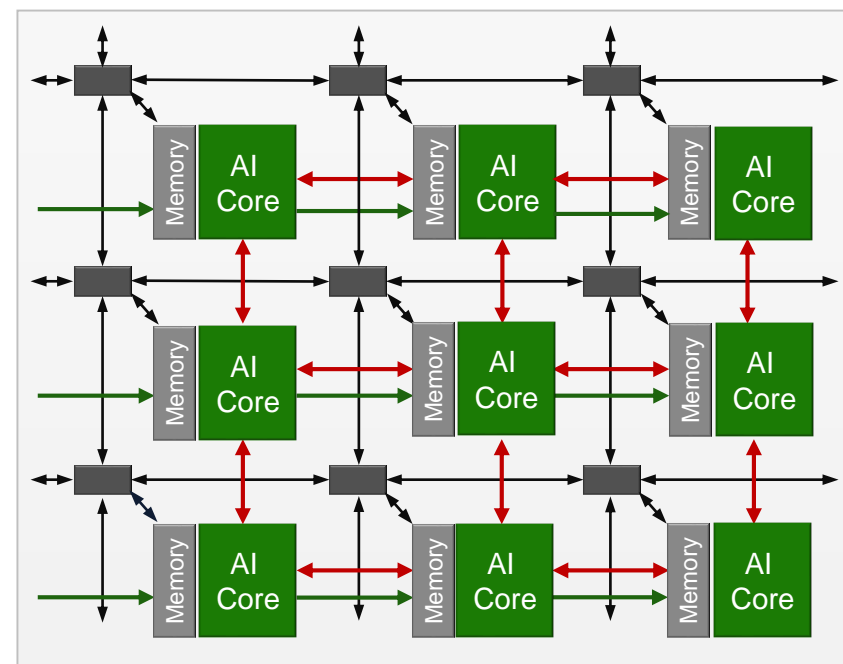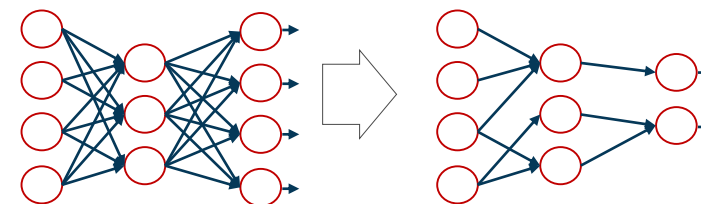> Instantiate multiple tiles (10s to 100s) for scalable compute

## Terabytes/sec of interface bandwidth to other engines

> Direct, massive throughput to adaptable HW engines

> Implement core application with AI for "Whole App Acceleration"

## SW programmable for any developer

> C programmable, compile in minutes
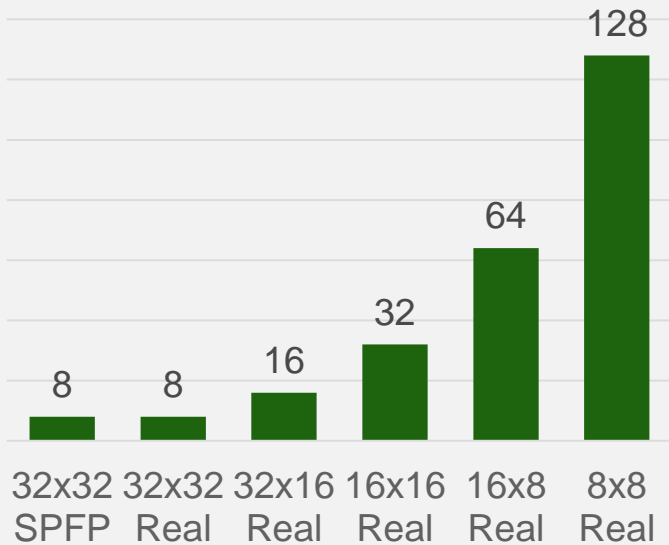
> Library-based design for ML framework developers

ML Inference and Optimizations
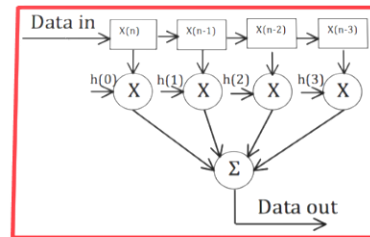
XILINX

# AI Engine: Multi-Precision Math Support

## Real Data Types

### MACs / Cycle (per core)



| Value | Data Type |
|---|---|
| 8 | 32x32 SPFP |
| 8 | 32x32 Real |
| 16 | 32x16 Real |
| 32 | 16x16 Real |
| 64 | 16x8 Real |
| 128 | 8x8 Real |

## Optimized For:



$$\begin{bmatrix} 0 & 2 & 5 & 2 \\ .4 & 0 & 0 & 0 \\ 0 & .5 & 0 & 0 \\ 0 & 0 & .6 & 0 \end{bmatrix} \begin{bmatrix} 5 \\ 2 \\ 2 \\ 2 \end{bmatrix} = \begin{bmatrix} 18 \\ 2 \\ 1 \\ 1.2 \end{bmatrix}$$

**Linear Algebra**

Matrix-Matrix Mult

Matrix-Vector Mult



**Convolution**

FIR Filters

2-D Filters

$$F(x) = \sum_{n=0}^{N-1} f(n)e^{-j2\pi\left(x\frac{n}{N}\right)}$$

$$f(n) = \frac{1}{N}\sum_{n=0}^{N-1} F(x)e^{j2\pi\left(x\frac{n}{N}\right)}$$

**Transforms**

FFTs/IFFTs

DCT, etc

## Complex Data Types

### MACs / Cycle (per core)



| Value | Data Type |
|---|---|
| 2 | 32x32 Complex |
| 4 | 32x16 Complex |
| 8 | 16x16 Complex |
| 16 | 16 Complex x 16 Real |

# AI Engine Tile

> AI Engine core
  >> 512b SIMD vector units
    – Both fixed and floating point
    – 16KB program memory
  >> 32b scalar RISC processor
  >> 256-bit load (x2) and store units with individual AGUs
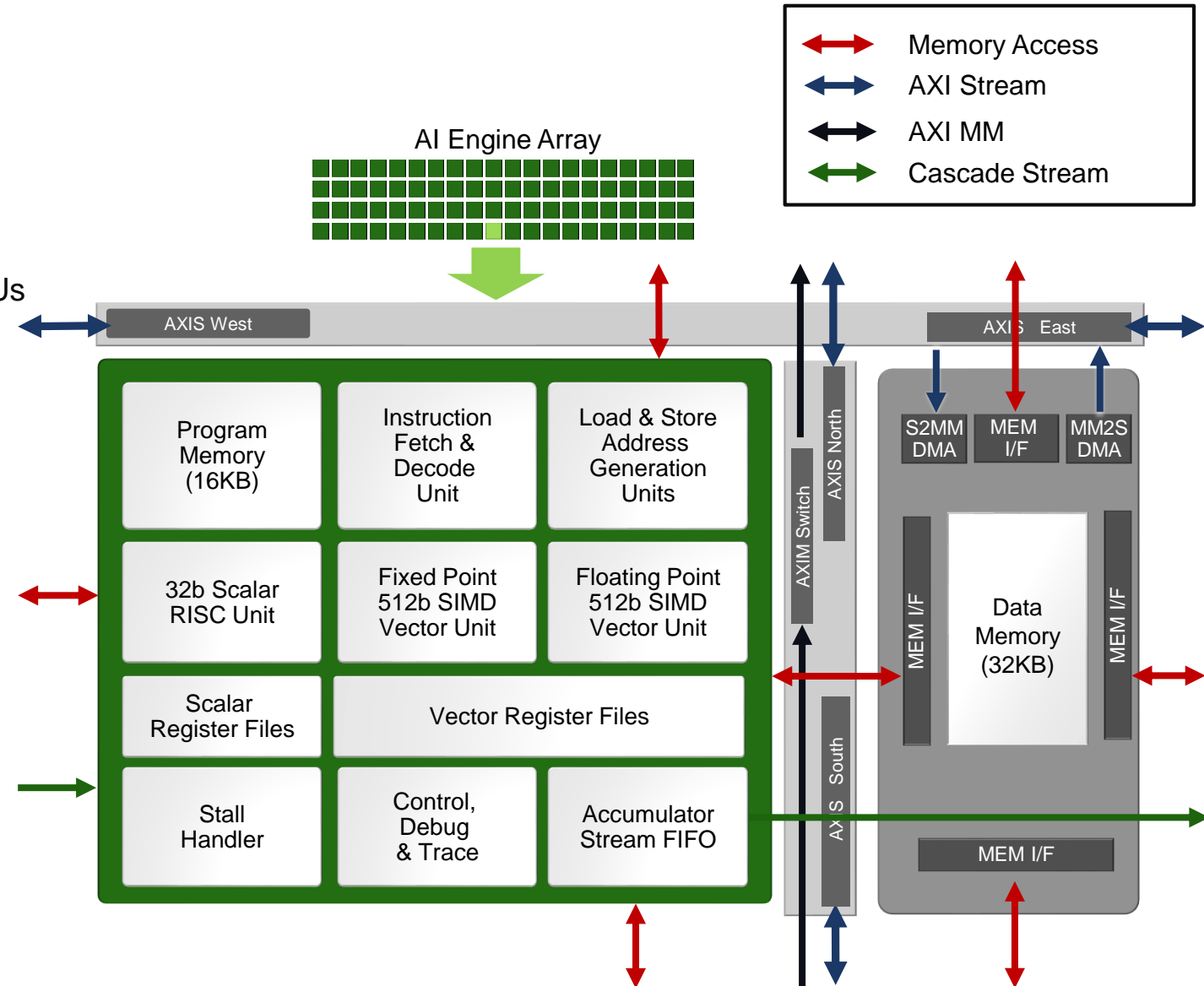
> 128KB direct core memory access
  >> 32KB local
  >> 32KB north, south, east & west

> Streaming interconnects
  >> AXI Memory Mapped (AXI-MM) switch
    – Configuration, control and debug
  >> AXI-Stream crossbar switch
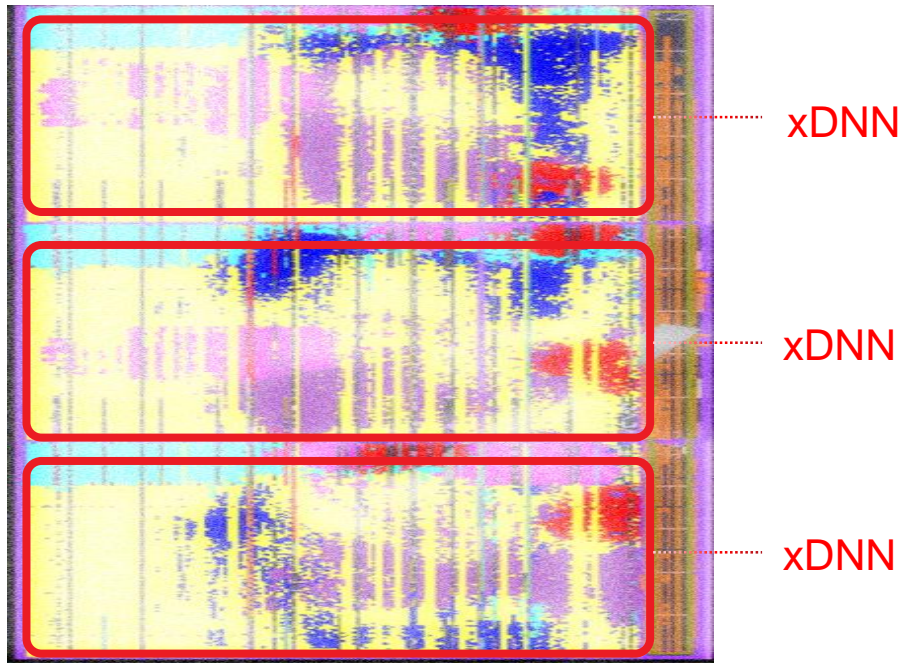    – Routing N/S/E & west around the array

> Debug/Trace/Profile functionality
  >> Debug using memory-mapped AXI4 i/f
  >> Connect to PMC via JTAG or HSDP



AI Engine Array

Legend:
- Memory Access
- AXI Stream
- AXI MM
- Cascade Stream

AXIS West | AXIS East

Program Memory (16KB) | Instruction Fetch & Decode Unit | Load & Store Address Generation Units

32b Scalar RISC Unit | Fixed Point 512b SIMD Vector Unit | Floating Point 512b SIMD Vector Unit

Scalar Register Files | Vector Register Files

Stall Handler | Control, Debug & Trace | Accumulator Stream FIFO

AXIM Switch | AXIS North | AXIS South

S2MM DMA | MEM I/F | MM2S DMA

MEM I/F | Data Memory (32KB) | MEM I/F

MEM I/F

XILINX

# Leveraging AI Engines for Compute-Intensive Applications

## Virtex® UltraScale+™ VU9P

**100%** of Logic for ML Acceleration (20 TOPS INT8)



xDNN

xDNN

xDNN

## Versal™ AI Core Series

**>5X** More Compute w/AI Engines (133 TOPS INT8)



AI Engine

Optimized Libraries

ML Overlay

70% of Device for Application Differentiation

XILINX

# DSP Engines

*Versatility and Granular Control of Datapath*

## Enhanced Compute architecture

> Greater than 1GHz of performance
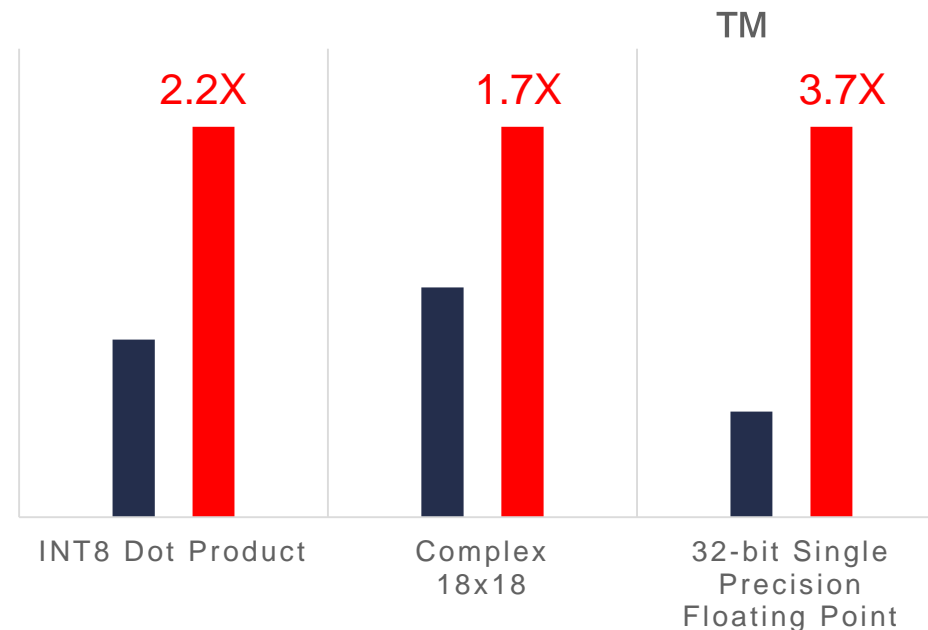
## Versatility for Wireless, ML, HPC, and more

> Integrated FP32, FP16 floating point, INT24 (HPC)

> Integrated complex 18x18 operation (wireless, cable access)

> Double the performance in INT8 operation (AI inference)

## Code Portability for UltraScale+™ 16nm designs

> Support for legacy IP and LogiCORE™ libraries

> Compatibility with SysGen, Model Composer, HLS tools
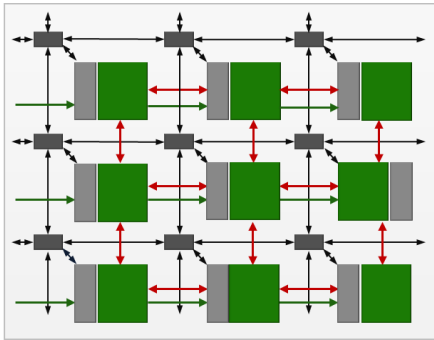
### Performance Improvement

■ UltraScale+ 16nm  ■ Versal  7nm

TM

| INT8 Dot Product | Complex 18x18 | 32-bit Single Precision Floating Point |
| --- | --- | --- |
| 2.2X | 1.7X | 3.7X |

**XILINX**
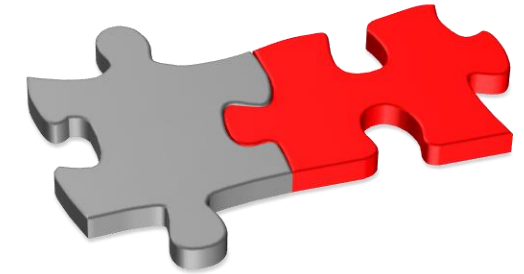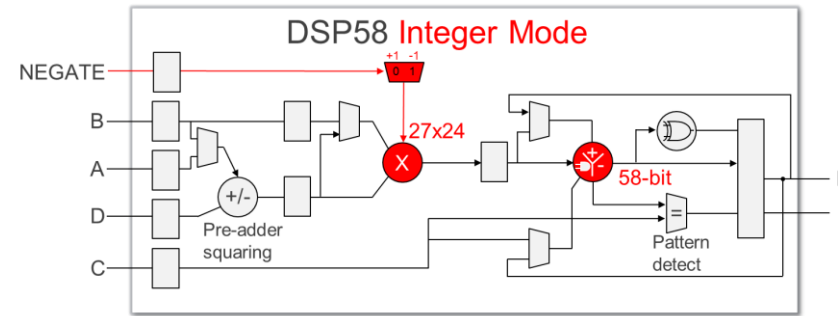
# AI and DSP Engines

## AI Engine 2D Array
VLIW and SIMD Architecture
C/C++ Programmable



## DSP Engine
Additional Features
RTL Entry





## Why AI Engine?

> Massive compute performance

> S/W programmable (C/C++)

> Fast compile – increase productivity

## Why DSP Engine?

> Existing RTL/HLS IP usage

> Additional features not available in AI Engine, e.g., 58-bit logic unit
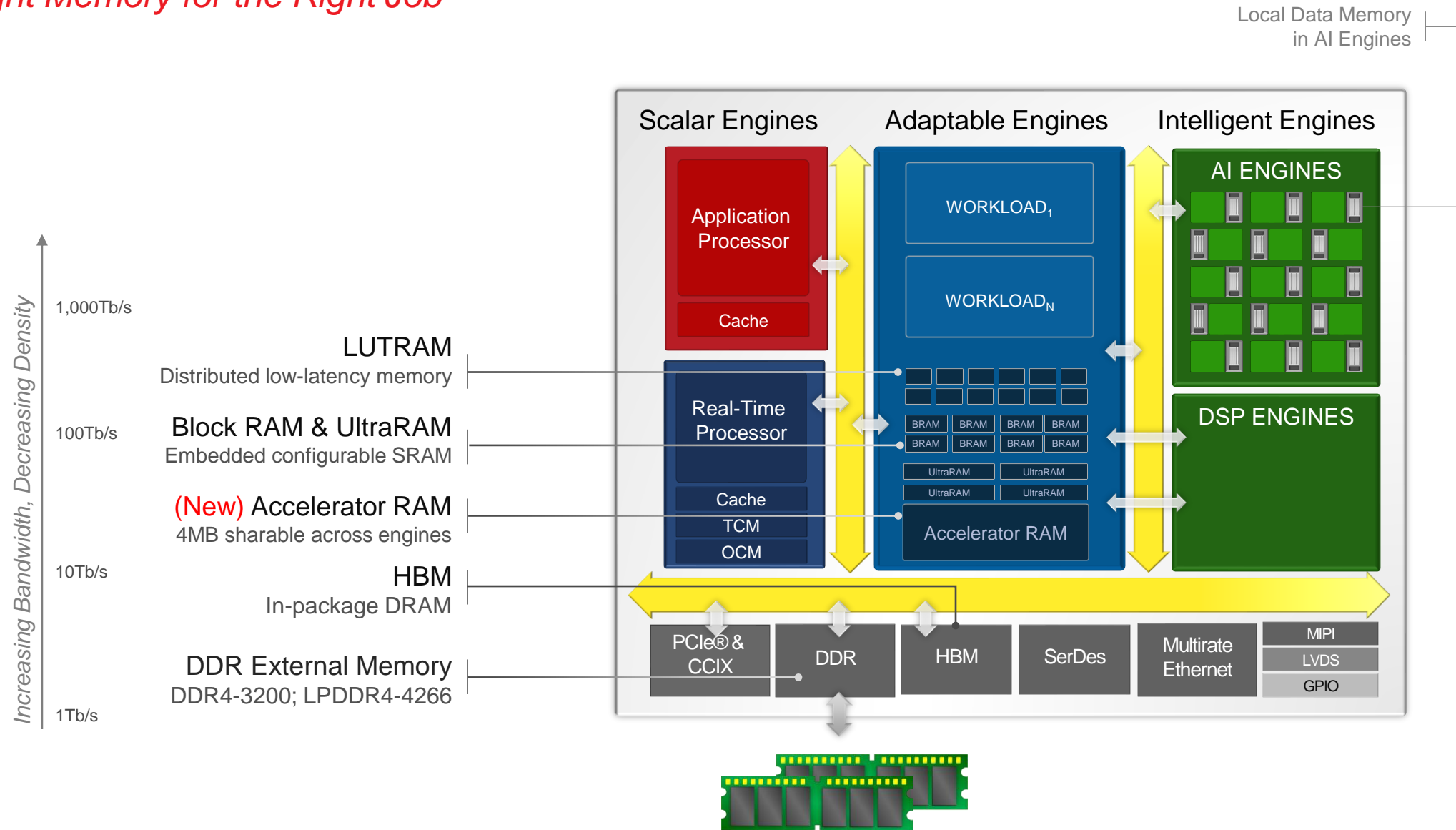
> Pre/Post processing to/from AI Engine

## Why both?
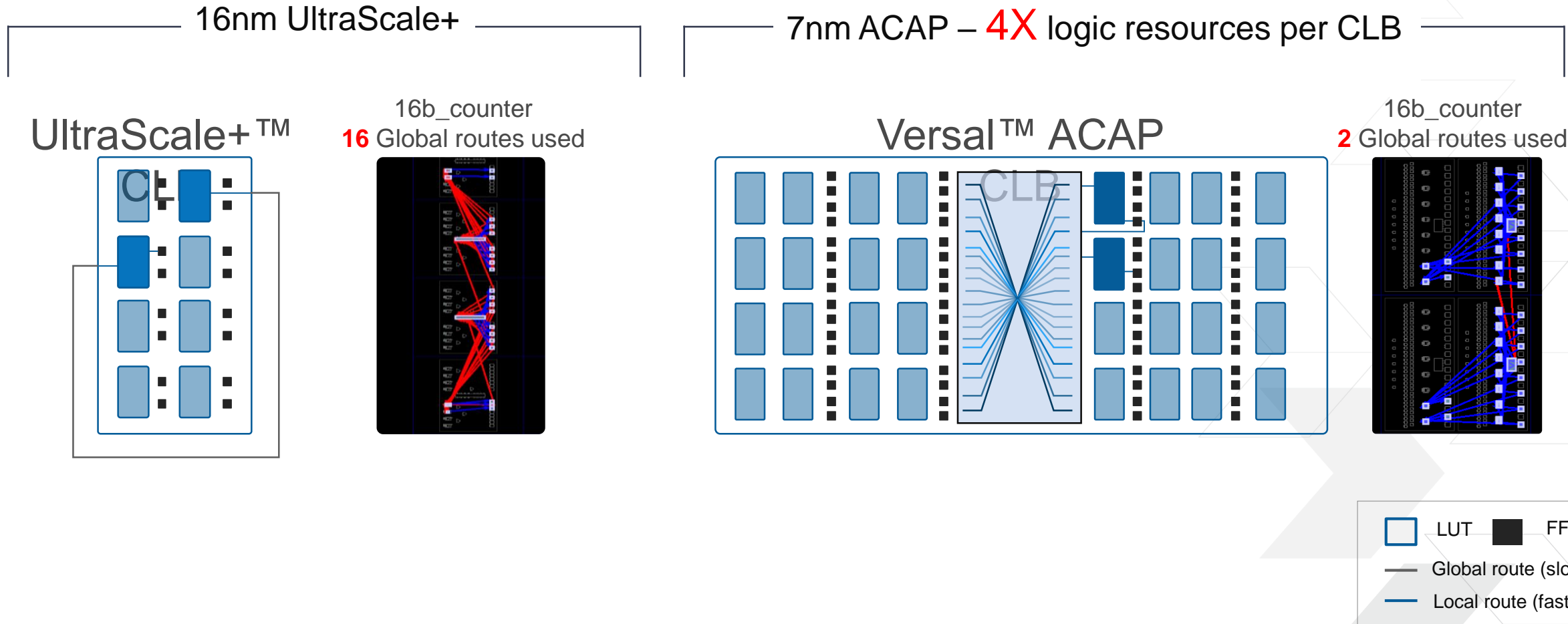
> AI Engine for efficiency

> PL & DSP for flexibility

**Versal™ ACAPs Accelerate the Complete Application**

&#8721; XILINX.

# Adaptable Cache-less Memory Hierarchy

*The Right Memory for the Right Job*



Local Data Memory in AI Engines

**Scalar Engines** — Application Processor — Cache — Real-Time Processor — Cache — TCM — OCM

**Adaptable Engines** — WORKLOAD$_1$ — WORKLOAD$_N$ — BRAM — UltraRAM — Accelerator RAM

**Intelligent Engines** — AI ENGINES — DSP ENGINES

PCIe® & CCIX — DDR — HBM — SerDes — Multirate Ethernet — MIPI — LVDS — GPIO

*Increasing Bandwidth, Decreasing Density*

1,000Tb/s
100Tb/s
10Tb/s
1Tb/s

**LUTRAM**
Distributed low-latency memory

**Block RAM & UltraRAM**
Embedded configurable SRAM

**(New) Accelerator RAM**
4MB sharable across engines

**HBM**
In-package DRAM

**DDR External Memory**
DDR4-3200; LPDDR4-4266

**XILINX**

# Re-Architected Hardware Logic for 4X Compute Density

16nm UltraScale+

7nm ACAP – **4X** logic resources per CLB

UltraScale+™

16b_counter
**16** Global routes used

Versal™ ACAP

16b_counter
**2** Global routes used

CLB

CLB

☐ LUT  ■ FF

— Global route (slow)
— Local route (fast)

**New CLB Interconnect Reduces Need for Global Interconnect**

⚡ **XILINX**

# NoC for Ease of Use, Guaranteed Bandwidth, and Power Efficiency
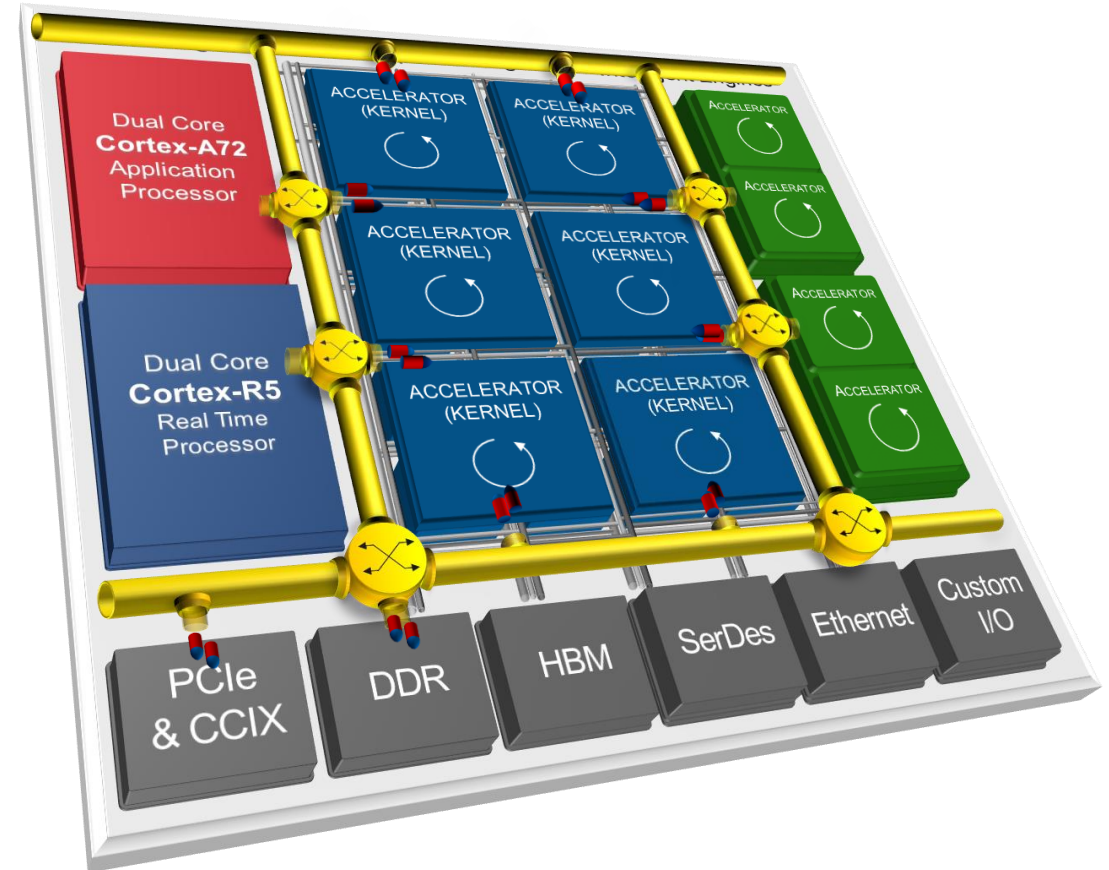
## High bandwidth terabit network-on-chip

> Memory mapped access to all resources

> Built-in arbitration between engines and memory

## High Bandwidth, Low Latency, Low power

> Guaranteed QoS

## Eases Kernel Placement

> Easily swap kernels at NoC port boundaries

> Simplifies connectivity between kernels

XILINX

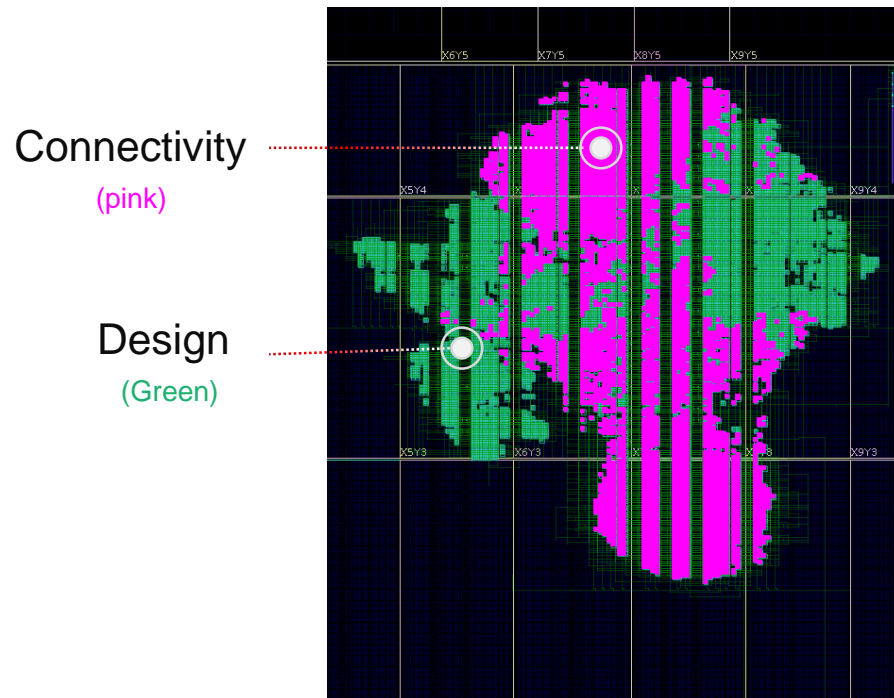# Programmable NoC vs. Logic Utilization using UltraScale+

*Simple Test Case – 4 AXI Traffic Generators Connected to 4 Block RAM*

## UltraScale+™ FPGA

Logic Resources for
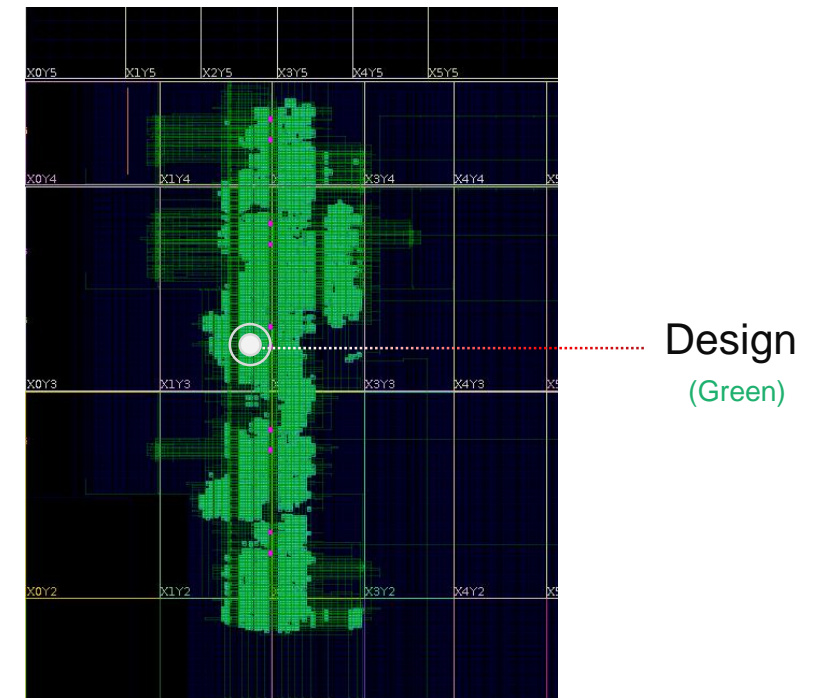Connectivity (SmartConnect)

42,012 FF

37,015 LUTs

## Versal™ ACAP

Logic Resources for
Connectivity

0 FF

0 LUTs



Connectivity
(pink)

Design
(Green)

Design Size
Reduction

~60% FF
~50% LUT

Design
(Green)

 XILINX.

# Introducing the "Integrated Shell"

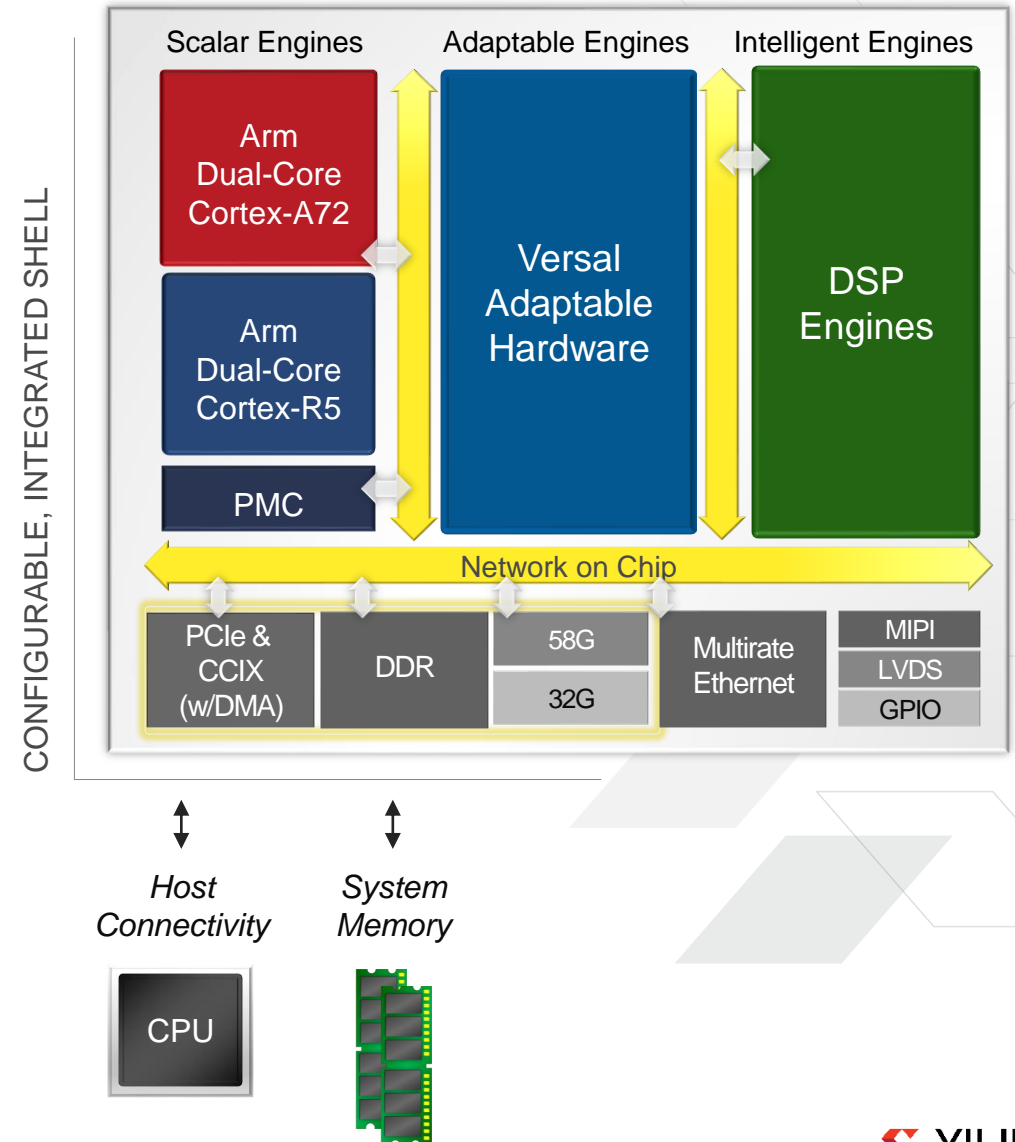**'Shell': Pre-Built Core Infrastructure & System Connectivity**

> External host interface

> Memory subsystem

> Basic interfaces (e.g., JTAG, USB, GbE)

**Key Architectural Elements of the Shell**

> Platform Management Controller (PMC)

> Integrated host interfaces: PCIe & CCIX, DMA

> Scalable Memory Subsystem: DDR4 & LPRDDR4

> Network-on-Chip for connectivity and arbitration

**Greater Performance, Device Utilization, and Productivity**

> More of the platform available for application's workload(s)

> Target application runs faster with less device congestion

> Turn-key, pre-engineered timing closure – no debug



CONFIGURABLE, INTEGRATED SHELL

Scalar Engines | Adaptable Engines | Intelligent Engines

Arm Dual-Core Cortex-A72

Arm Dual-Core Cortex-R5

PMC

Versal Adaptable Hardware

DSP Engines

Network on Chip

PCIe & CCIX (w/DMA) | DDR | 58G / 32G | Multirate Ethernet | MIPI / LVDS / GPIO

Host Connectivity

System Memory

CPU

XILINX

# Summary

> Heterogeneous processing is required for the future as there is no single processor type that is ideally suited for all algorithms and applications

> This can be a challenge for design development and productivity

> ACAPs are a response to this new reality

> Embracing multiple levels of abstraction through a unified platform allows developers to use the tools and languages that they are familiar with while simplifying debugging

*Visit https://www.xilinx.com/products/silicon-devices/acap/versal.html for datasheets, whitepapers, and product tables.*



Xilinx VC1902 Versal ACAP with 400 AI Engines.
First shipment June 2019.

XILINX

# Building the Adaptable, Intelligent World

XILINX.