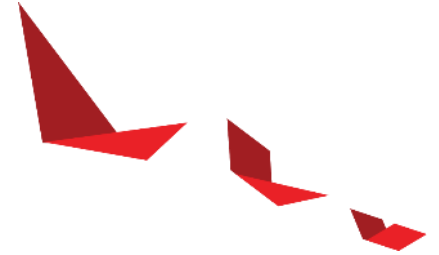




# Next generation compute efficiency with Xilinx FPGAs and the new Versal ACAP

Cathal McCabe  
Xilinx University Program, Xilinx Ireland  
18 February 2020

# Overview

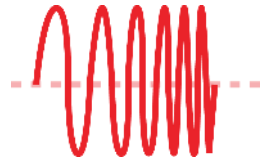


- ▶ Requirements for next generation compute systems
- ▶ The technology conundrum
- ▶ FPGA technology evolution
- ▶ Current and next generation Xilinx technologies

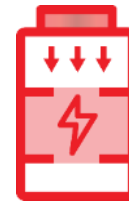
# What are your requirements for next generation HEP systems?



Performance



Data rates



Power



Cost



Compute density



Machine Learning

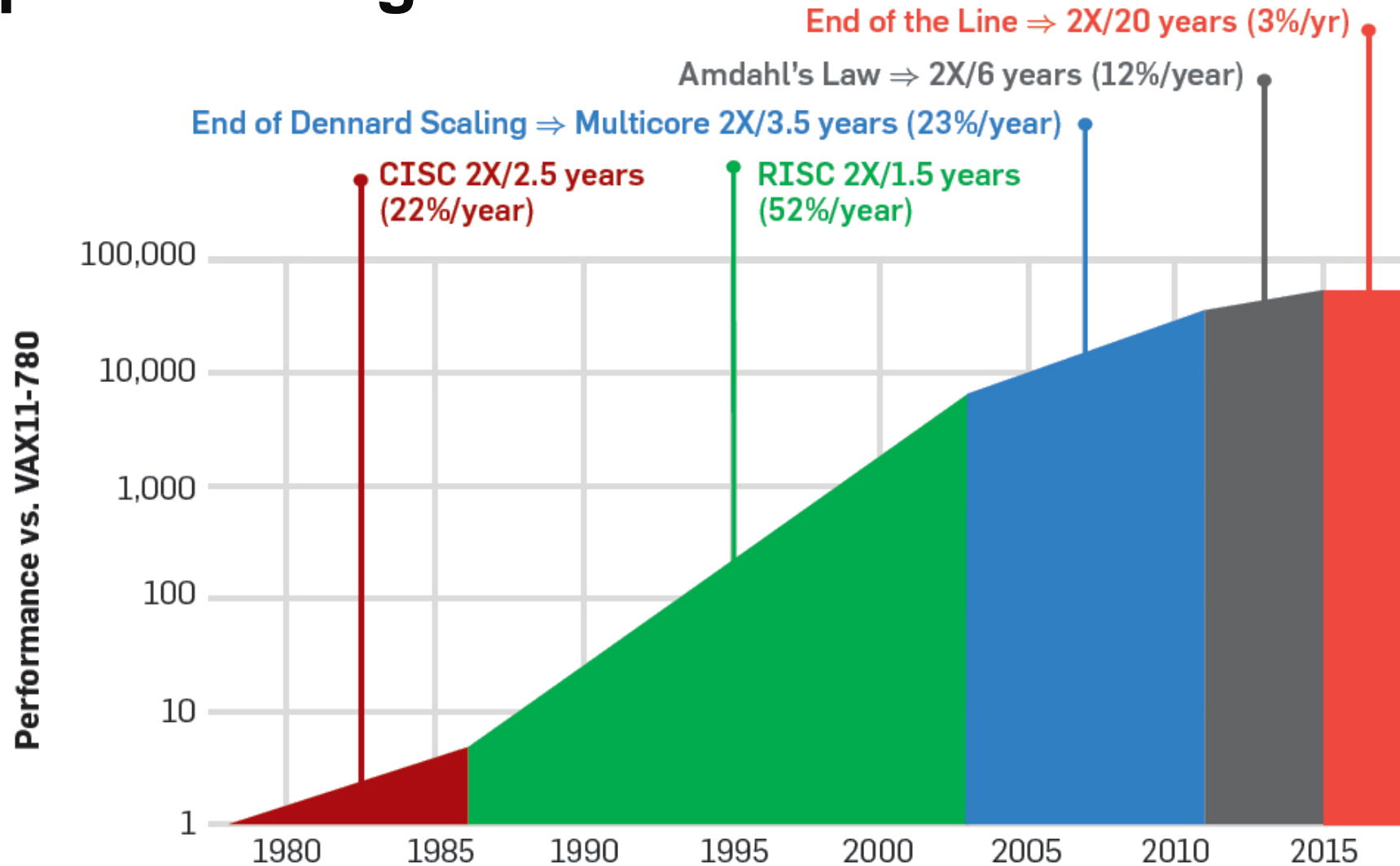


Adaptability



Cloud scalability

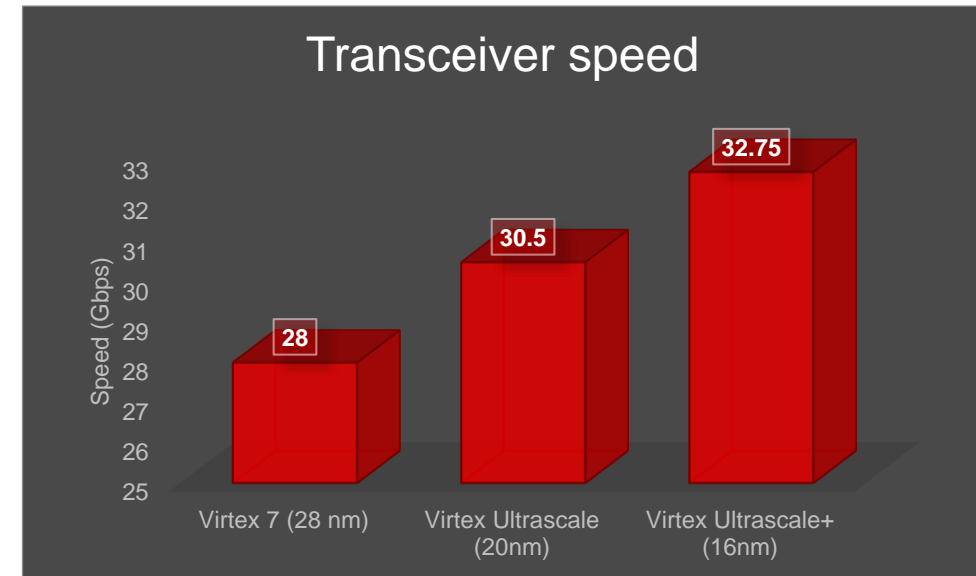
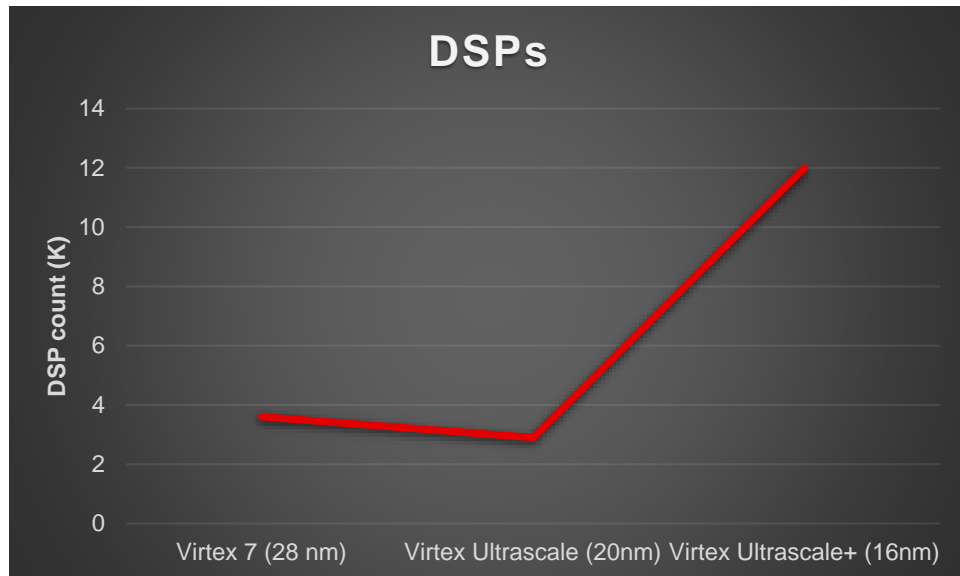
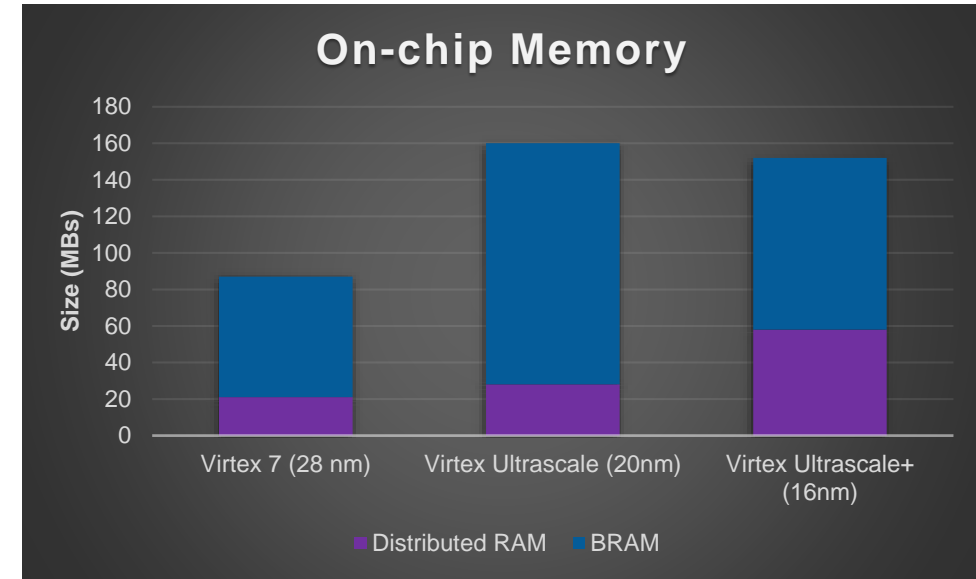
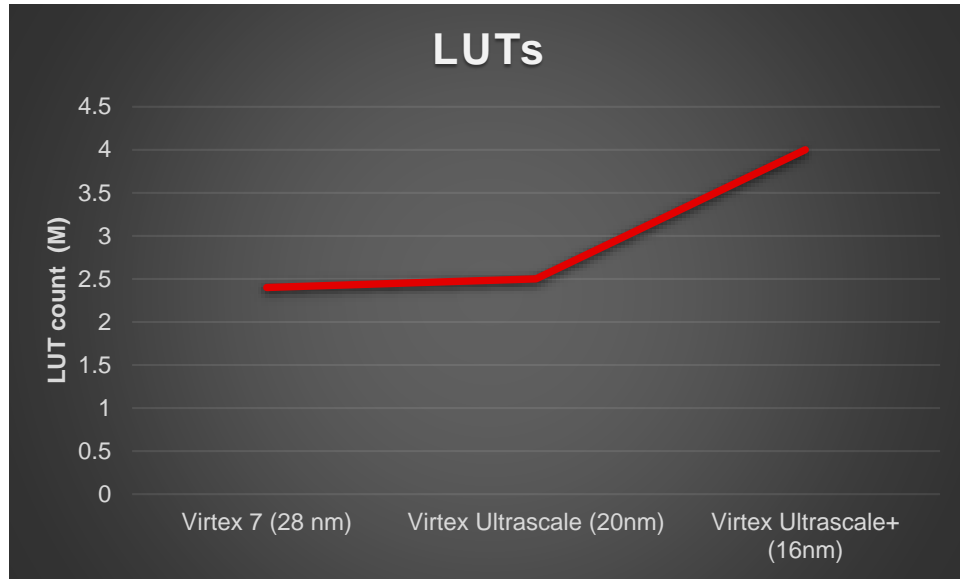
# The Technology Conundrum .. And the Need for a New Compute Paradigm



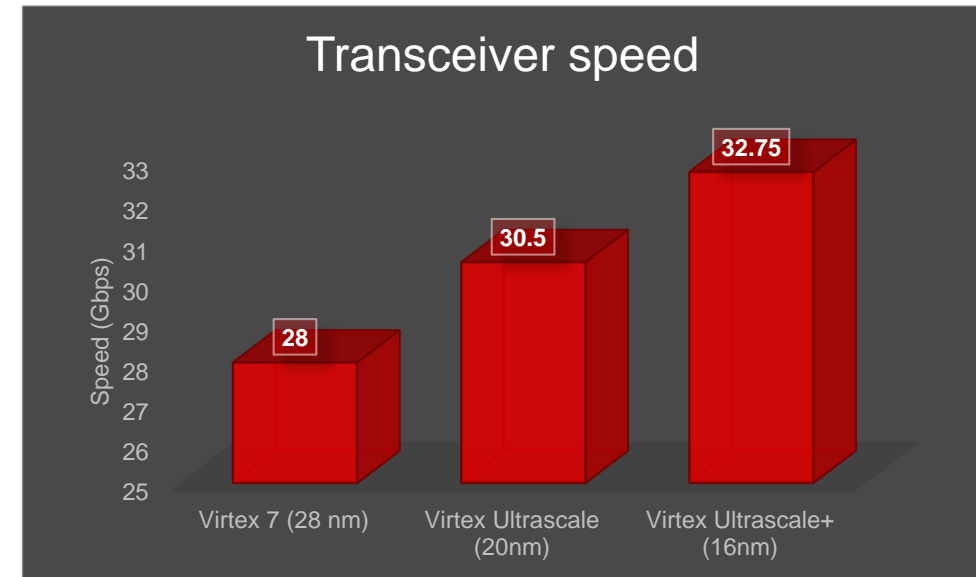
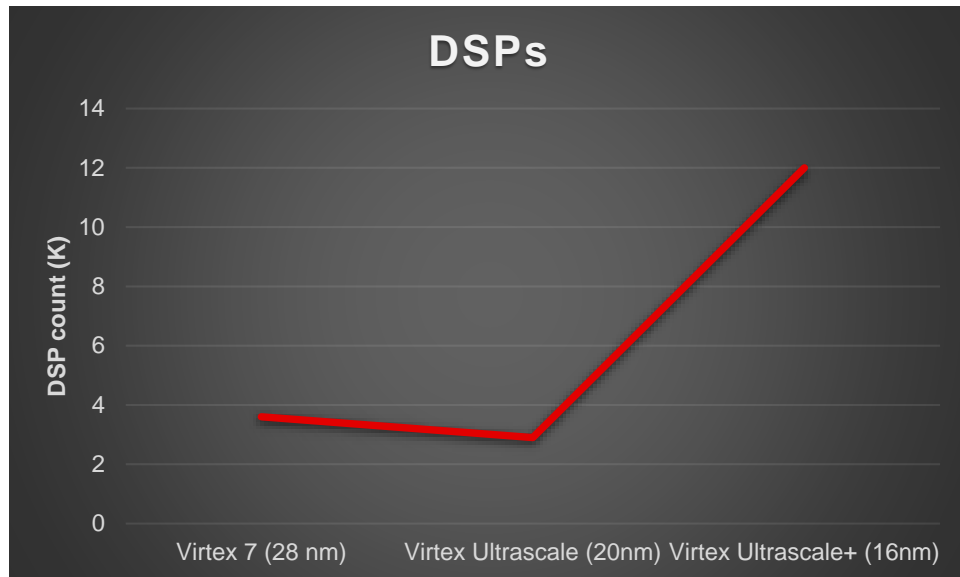
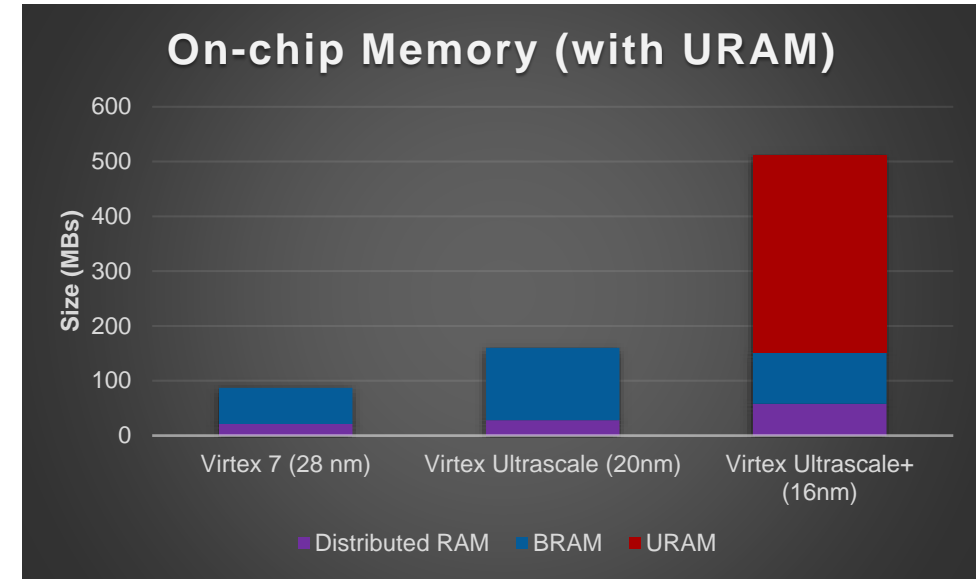
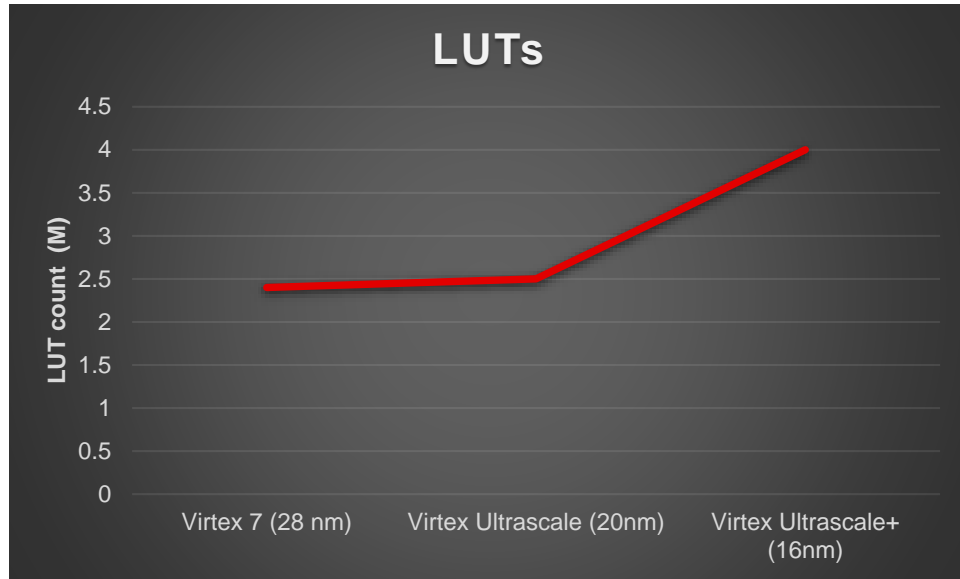
\*John Hennessy and David Patterson, Computer Architecture: A Quantitative Approach, 6/e. 2018

© Copyright 2020 Xilinx

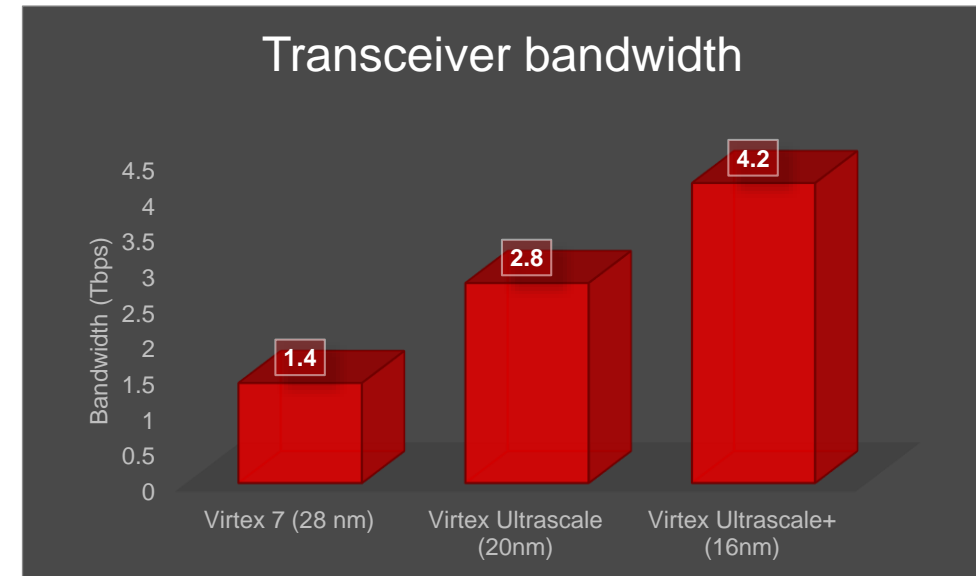
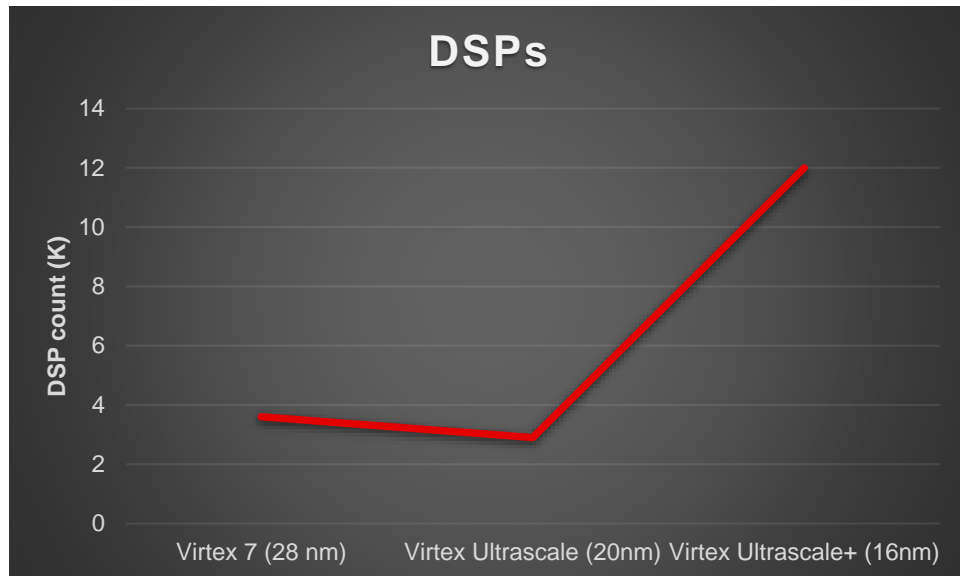
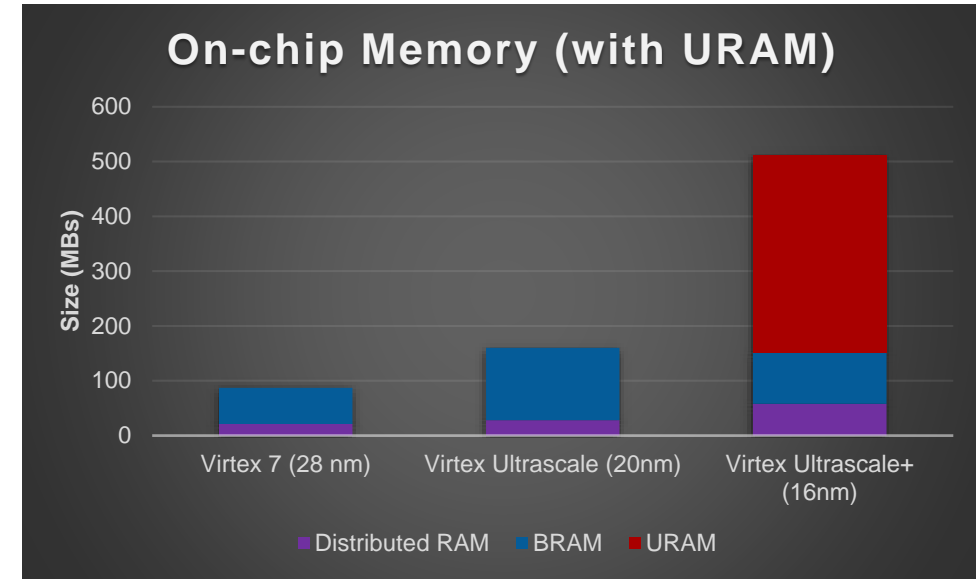
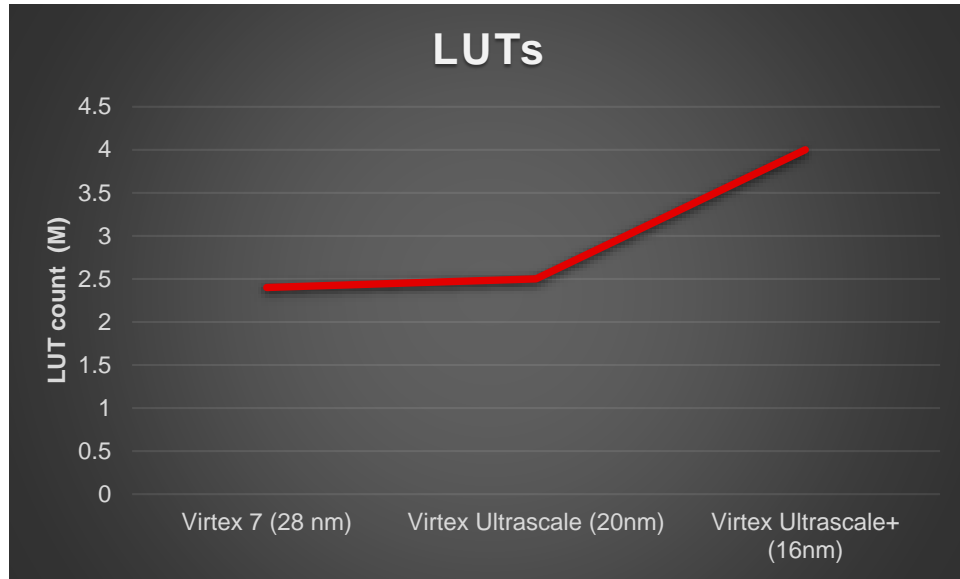
# FPGA scaling



# FPGA scaling - URAM example



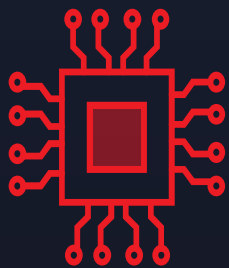
# FPGA scaling – transceivers example



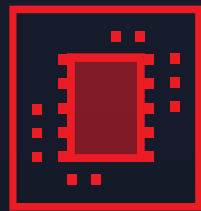
# Xilinx Transformation

From Devices to Platforms

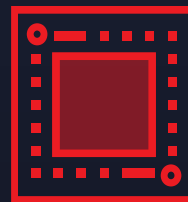
SW Programmability



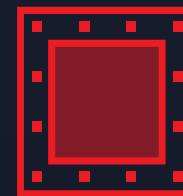
FPGA



SoC



MPSoC

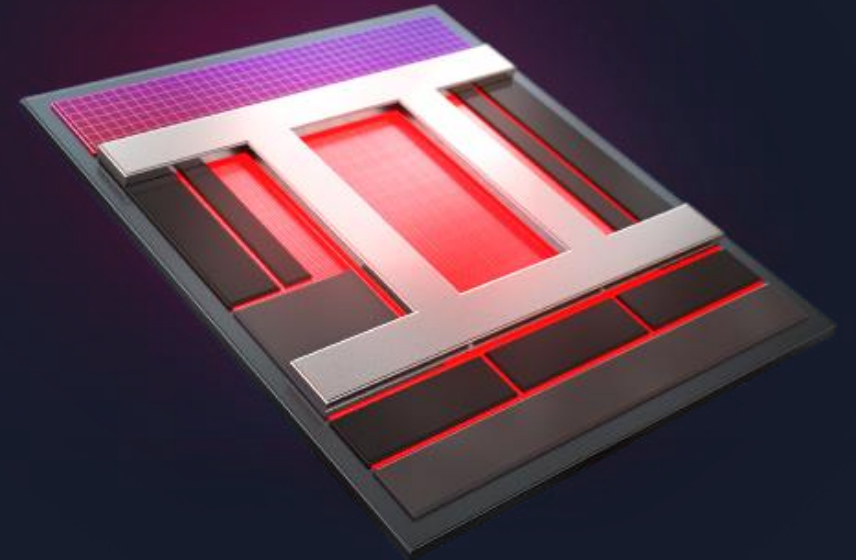


RFSoc



ACAP

Device Category





# Xilinx technologies



# ALVEO™



## FAST

Built for high throughput, ultra-low latency  
Accelerate compute, networking, storage



## ADAPTABLE

Deploy optimized domain-specific architectures  
Adapt to changing algorithms



## ACCESSIBLE

Deploy in the cloud or on-premises  
Rich set of accelerated Applications



# ➤ Data Center and AI Accelerator Cards



90x

Database Search  
& Analytics



89x

Financial  
Computing



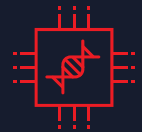
20x

Machine  
Learning



12x

Video  
Processing



10x

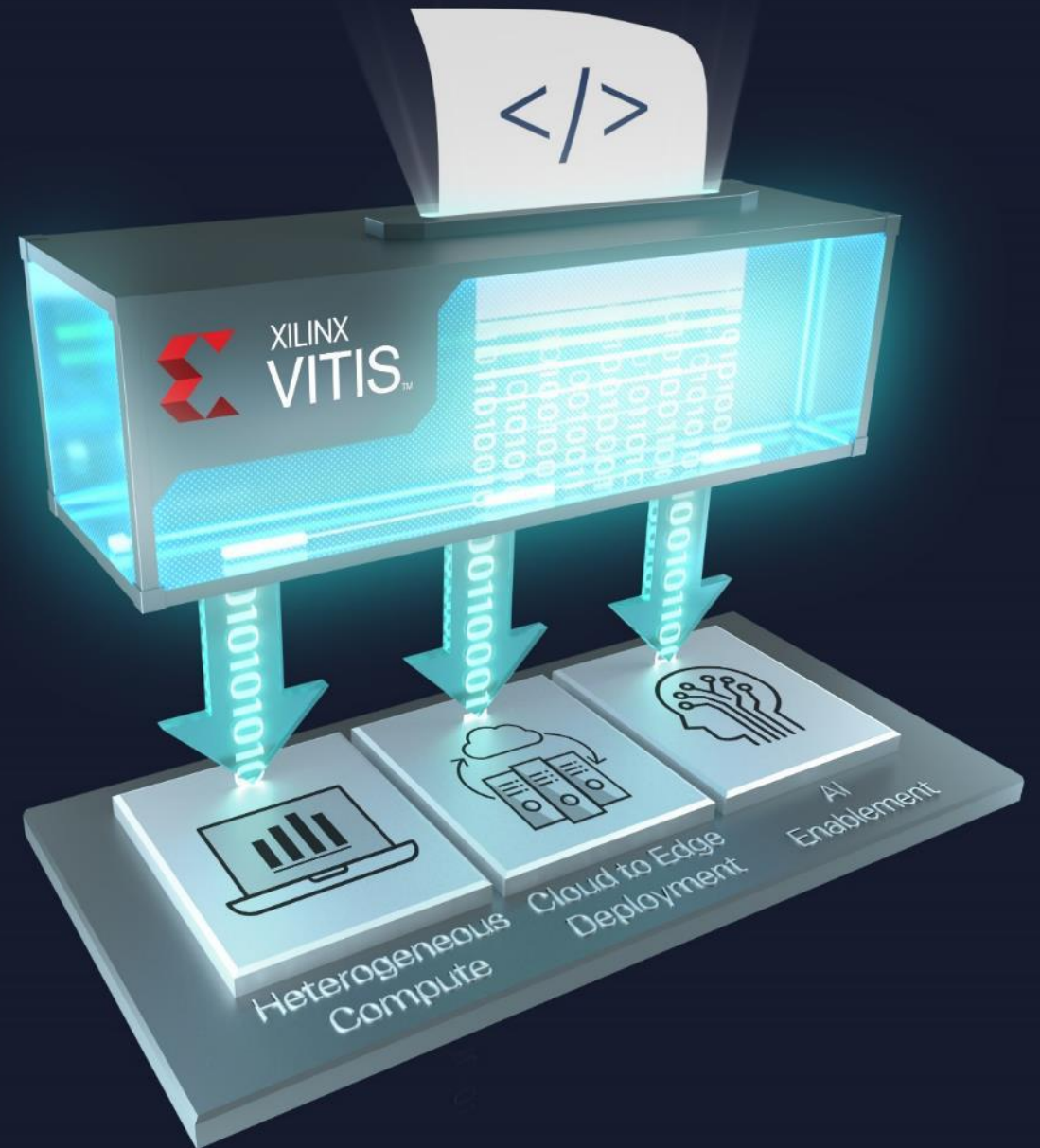
HPC &  
Life Sciences



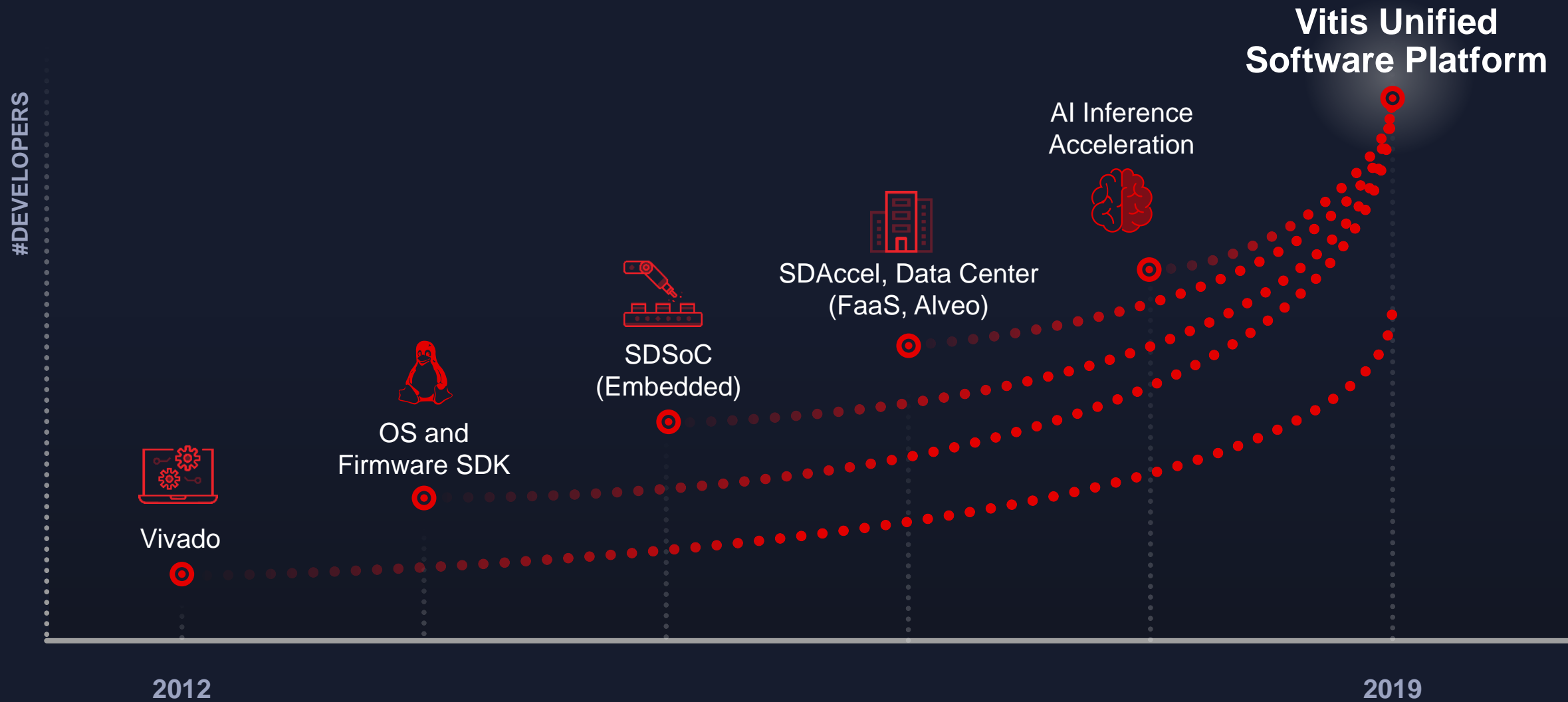
## Unified Software Platform

- Unified methodology edge to cloud
- Focus on platform and acceleration
- Available for free

\*Open source Xilinx Runtime library (XRT), Accelerated libraries, AI Models



# Platform Transformation



# ➤ Vitis: Unified Software Platform

Domain-specific  
development  
environment

Vitis accelerated  
libraries

Vitis core  
development kit

OpenCV  
Library

BLAS  
Library

Finance  
Library

TensorFlow

Vitis AI

FFmpeg

Vitis Video

**Partners**  
Genomics,  
Data Analytics,  
And more

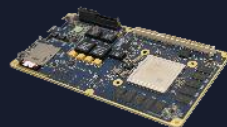
Compilers

Analyzers

Debuggers

Xilinx runtime libraries (XRT)

Vitis target platform

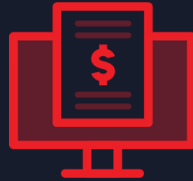


# ➤ Build: Extensive, Open Source Libraries

## Domain-Specific Libraries



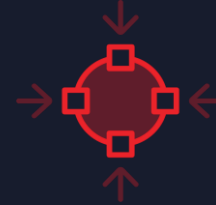
Vision &  
Image



Finance



Data Analytics &  
Database

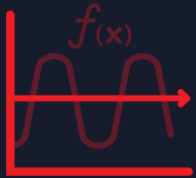


Data Compression



Data Security

## Common Libraries



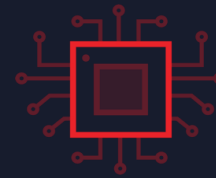
Math



Linear Algebra



Statistics



DSP



Data Management

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



# ➤ Vitis AI: Deep Learning Acceleration

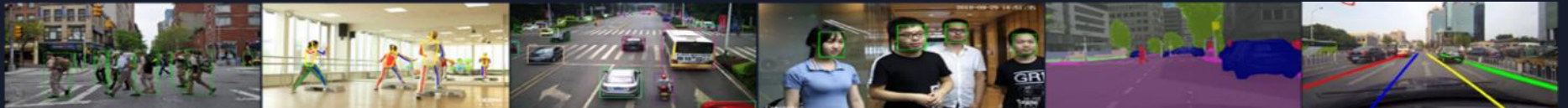
Frameworks

TensorFlow

Caffe

PyTorch

Vitis AI  
models



Vitis AI  
development kit

AI Optimizer

AI Quantizer

AI Compiler

AI Profiler

AI Library

Xilinx runtime library (XRT)

Deep Learning  
Processing Unit

CNN DPU

LSTM DPU

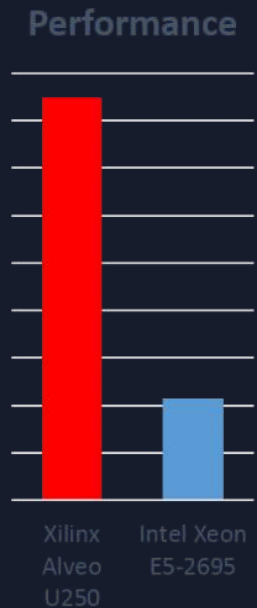
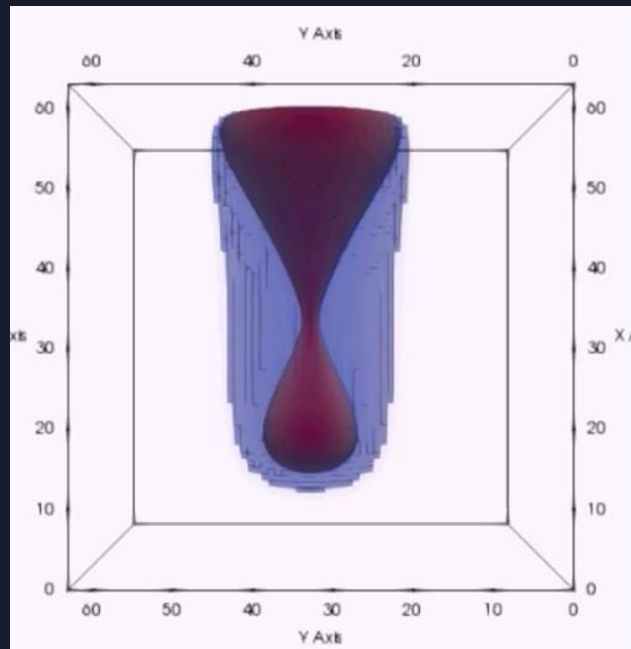
MLP DPU



# Example performance

# ➤ Computational Fluid Dynamics

## ALVEO Accelerated CFD Kernels



byte  
LAKE

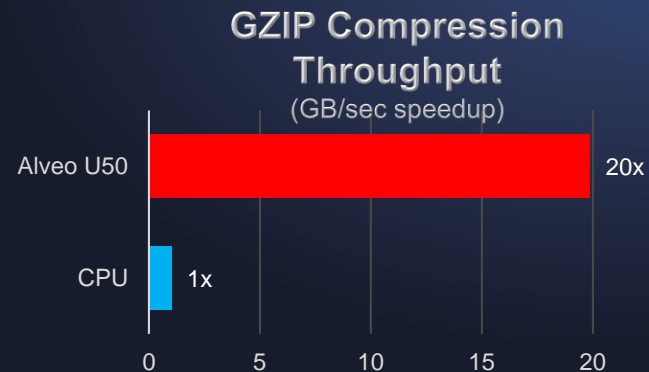
## Faster Time to insight, Fewer Nodes

- 4x Faster simulation time
- 80% lower energy consumption
- 6x better performance per Watt

# Computational Storage

## Line-rate Data Compression Acceleration

Compression, decompression, erasure coding, encryption all accelerated on one platform



Intel Skylake-SP 6152 @2.10GHz 22-core CPU (Ubuntu 16.04), GB/s compression per CPU core = .0229. Alveo U50 = 10GB/s

# ➤ Computational Storage

Line-rate Data Compression Acceleration

20x Throughput Per Node  
2x Less Nodes  
40% Lower Total Cost



Alveo U50  
Acceleration

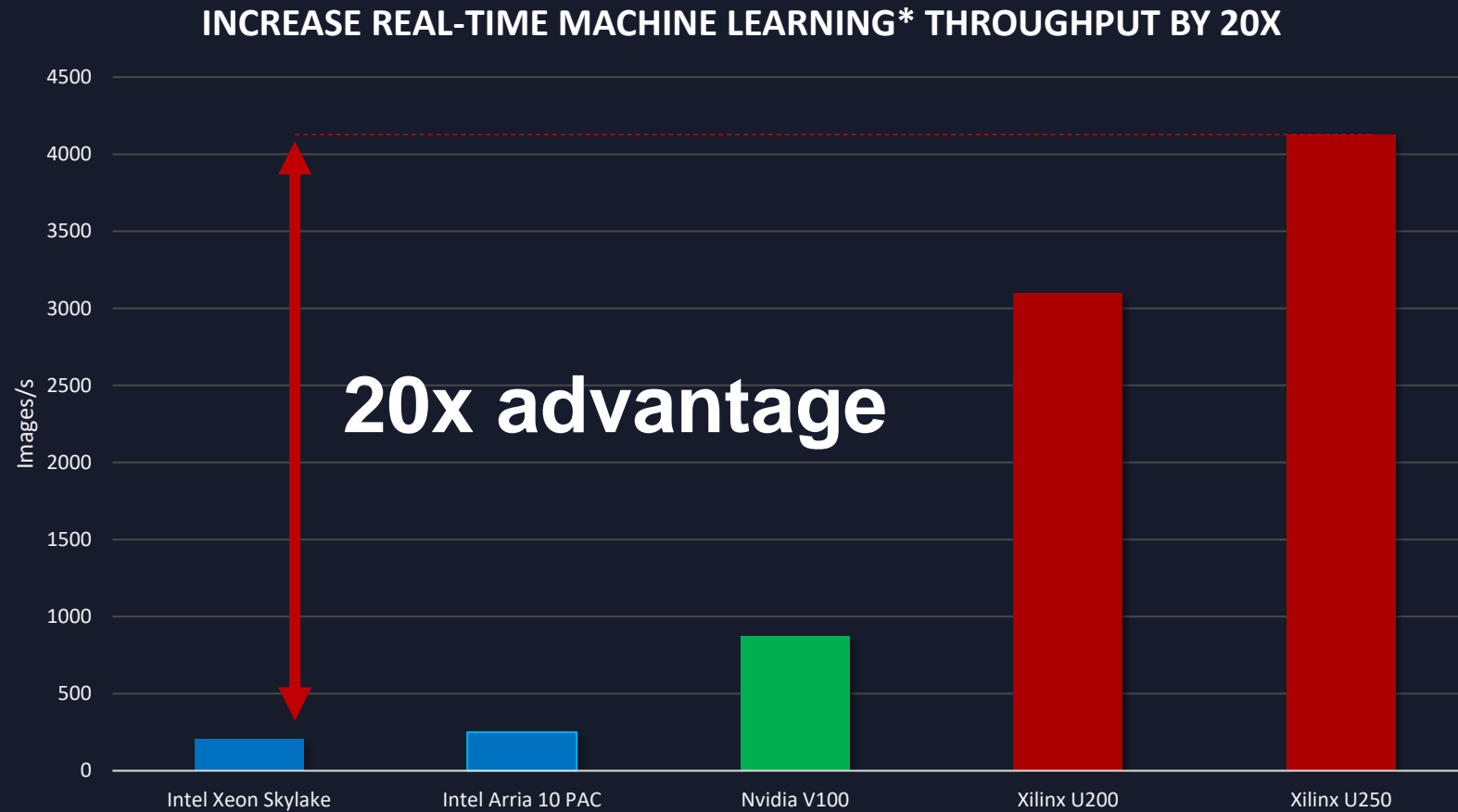


**2x Dual CPU Servers**  
192TB SSDs, 1GB/sec Per Node  
Compression Throughput

**Alveo Server with 2x Alveo U50**  
96TB SSDs (192TB effective),  
20GB/sec Per Node Compression  
Throughput

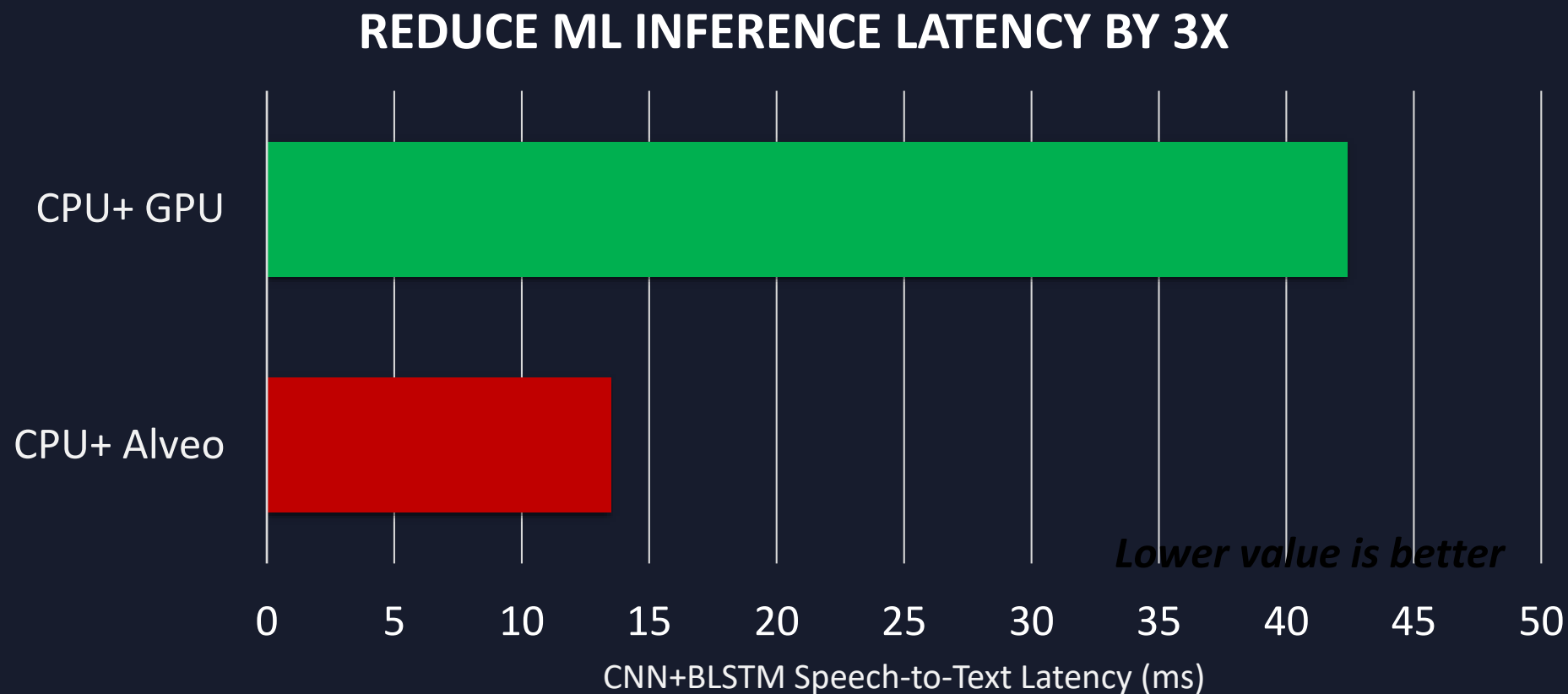
Intel Skylake-SP 6152 @2.10GHz CPU (Ubuntu 16.04), GB/s compression per CPU core = .0229. Alveo U50 = 10GB/s, Assume 2:1 compression

# ➤ Advantages in Machine Learning Inference



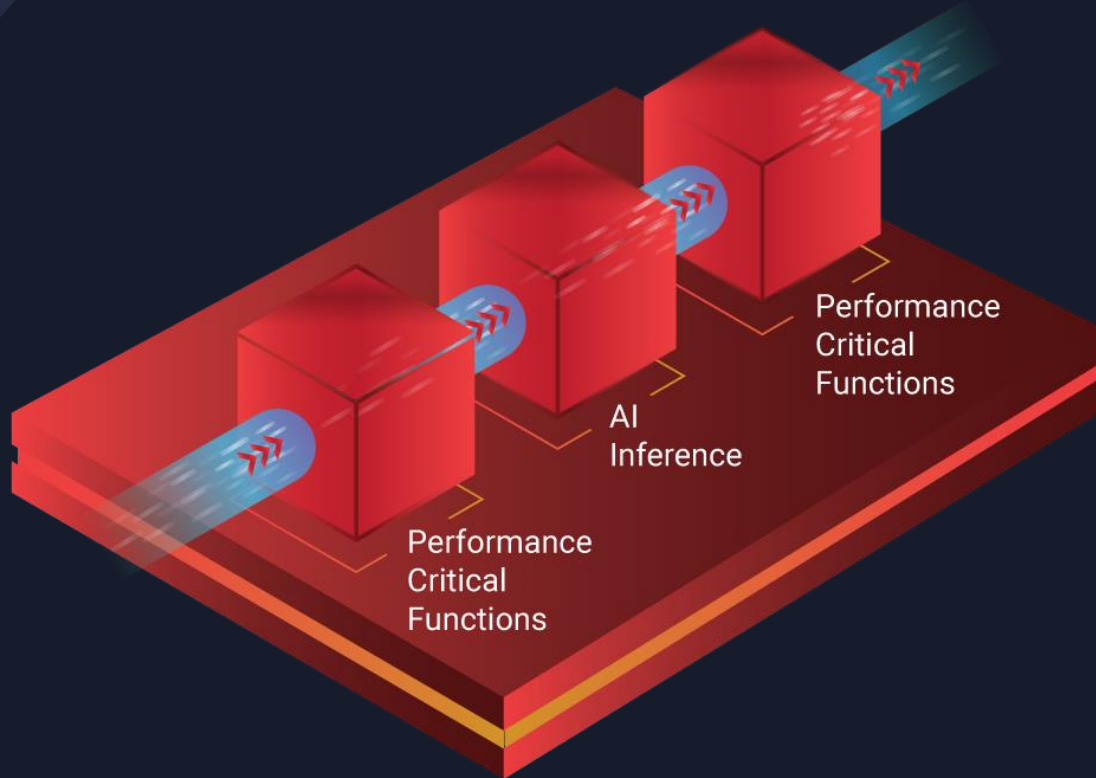
\* Source: [Accelerating DNNs with Xilinx Alveo Accelerator Cards White Paper](#)

# ➤ *Advantages in Latency*

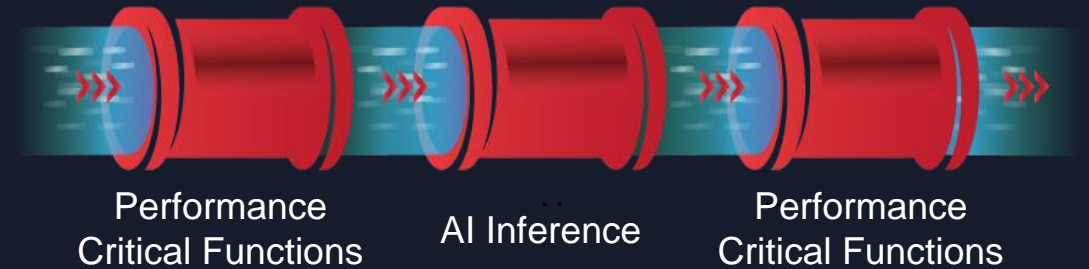


**Alveo Provides Massive Parallel Compute with Lowest Latency vs GPUs**

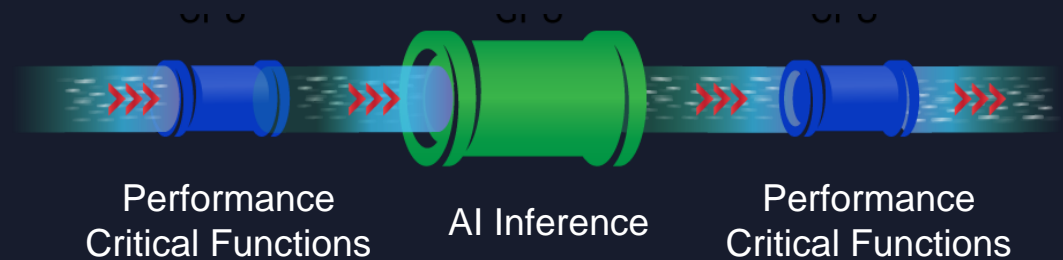
# ➤ Whole Application Acceleration



Xilinx – Matched Throughput



Other solutions – Mismatched Throughput





# ➤ AI Accelerated Dark Matter Search (CERN)

Real-time ML Inference + Sensor pre-processing



100ns Inference Latency on 150 Terabytes/Second Data Rates

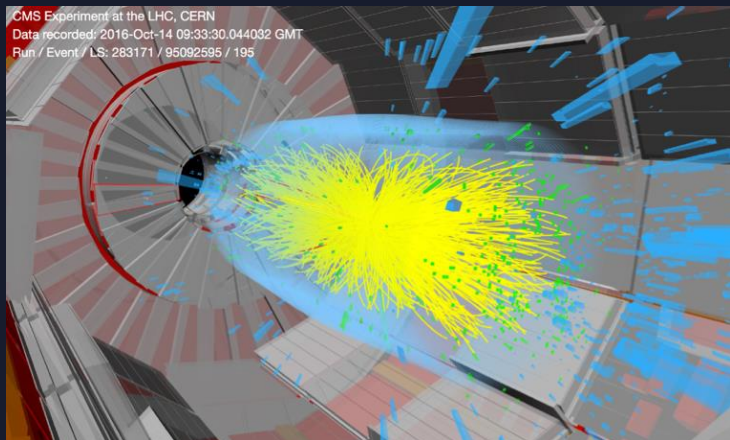


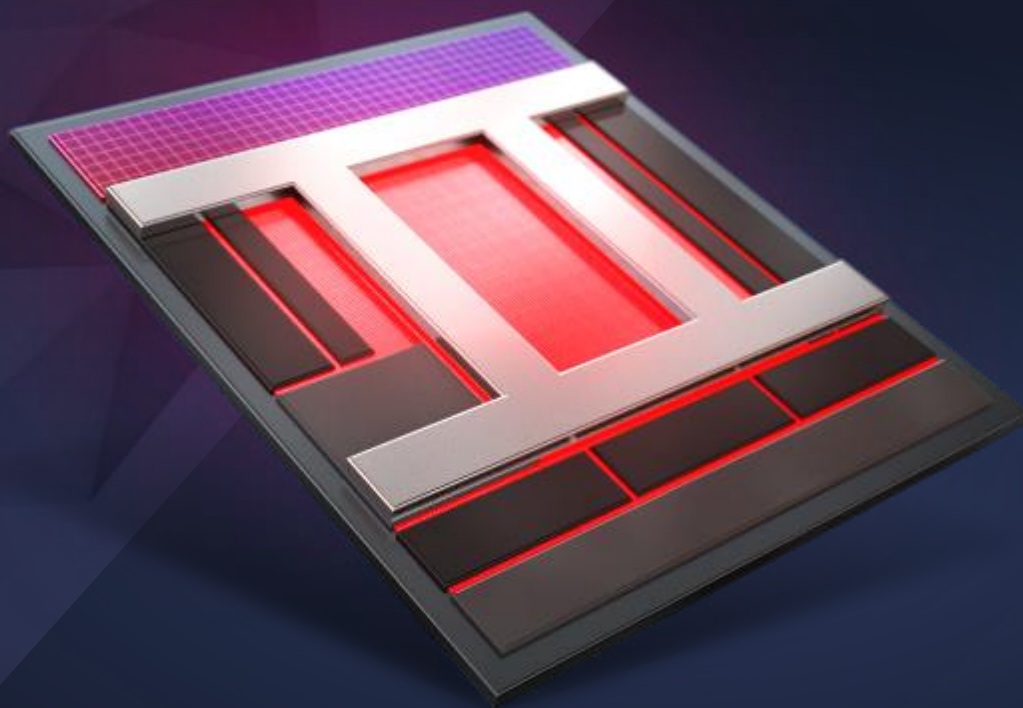
Figure 3: hls4ml tool flow

<https://www.xilinx.com/content/dam/xilinx/publications/powered-by-xilinx/cern-case-study-final.pdf>

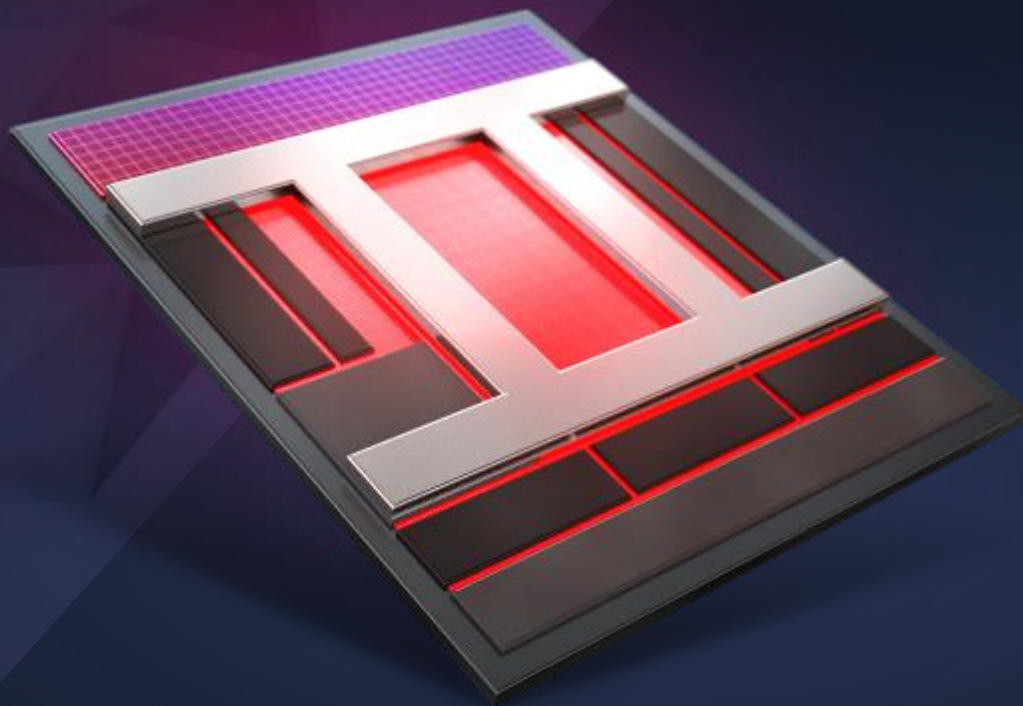




# Introducing the Versal ACAP

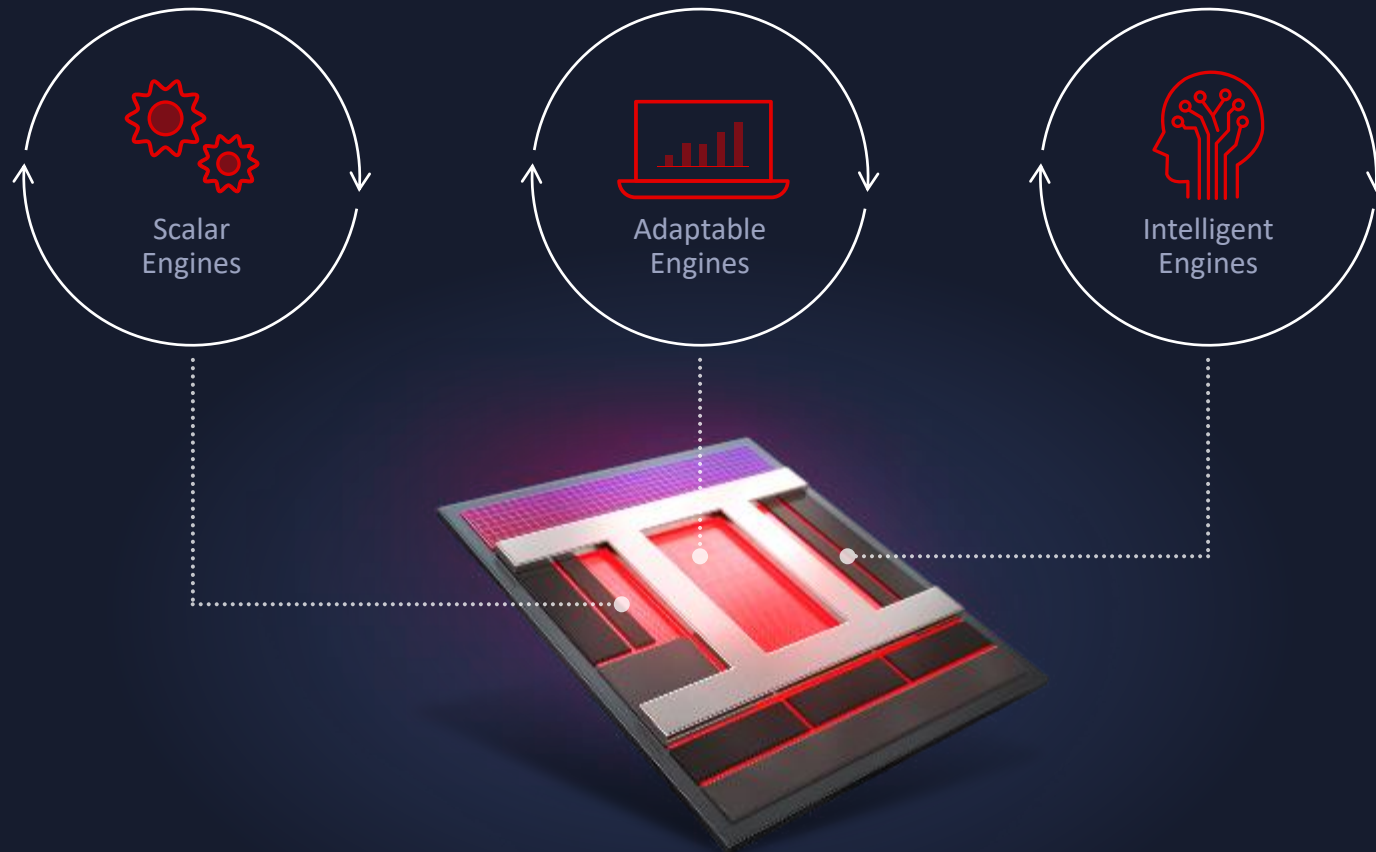


# ACAP



# Adaptive Compute Acceleration Platform

# Compute Acceleration





XILINX®  
**VERSAL™**

**The Industry's First ACAP**

Heterogeneous Acceleration

For Any Application

For Any Developer



7nm  
FinFET





# Versal ACAP Technology Tour



Scalar Processing Engines



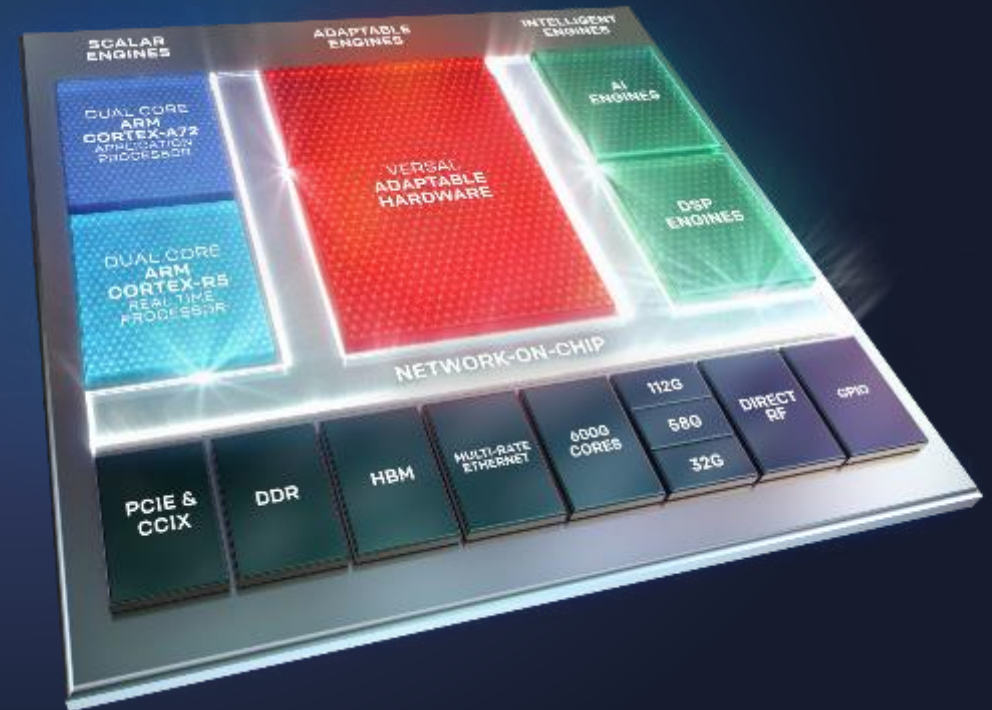
Adaptable Hardware Engines



Intelligent Engines  
SW Programmable, HW Adaptable



Breakout Integration of Advanced  
Protocol Engines

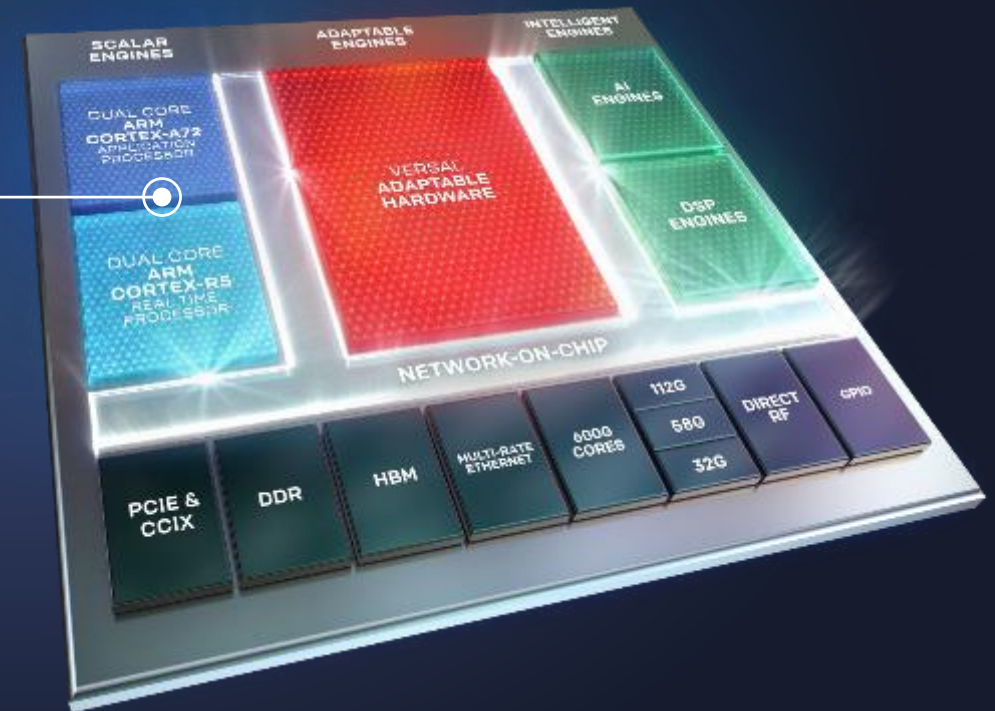


# Scalar Processing Engines

Arm Cortex-A72  
Application Processor

Arm Cortex-R5  
Real-Time Processor

Platform Management Controller

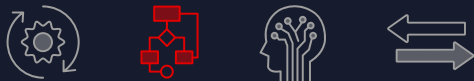
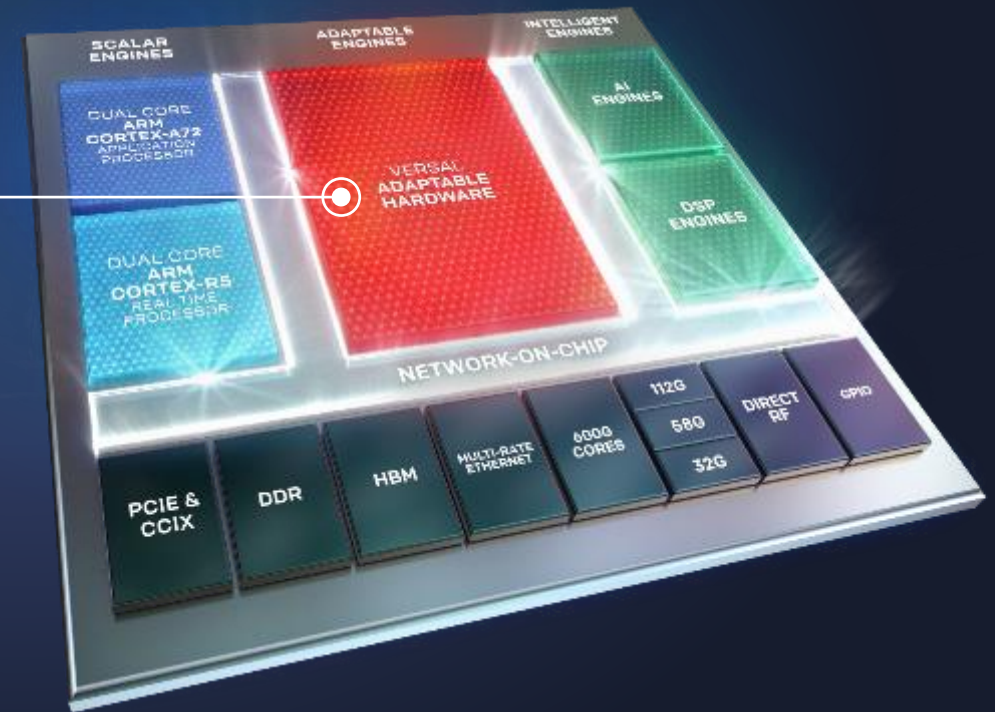


# Adaptable Hardware Engines

Re-architected foundational HW fabric for greater compute density

Enables custom memory hierarchy

8X Faster Dynamic Reconfiguration (“on-the-fly”)





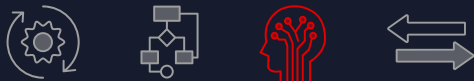
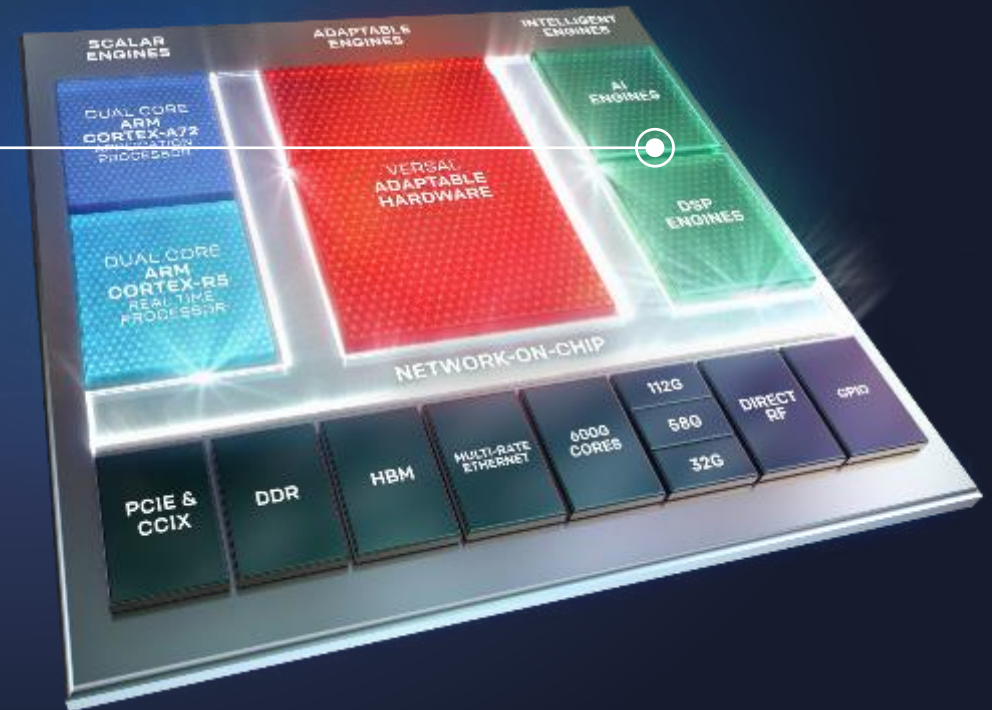
# Intelligent Engines

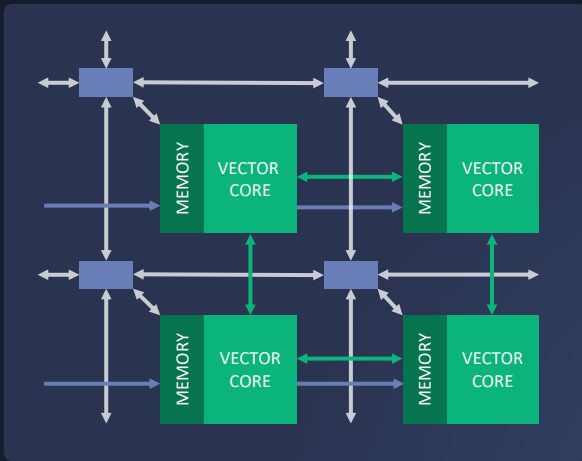
## DSP Engines

High-precision floating point & low latency  
Granular control for customized datapaths

## AI Engines

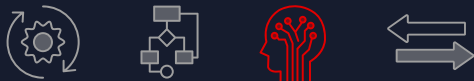
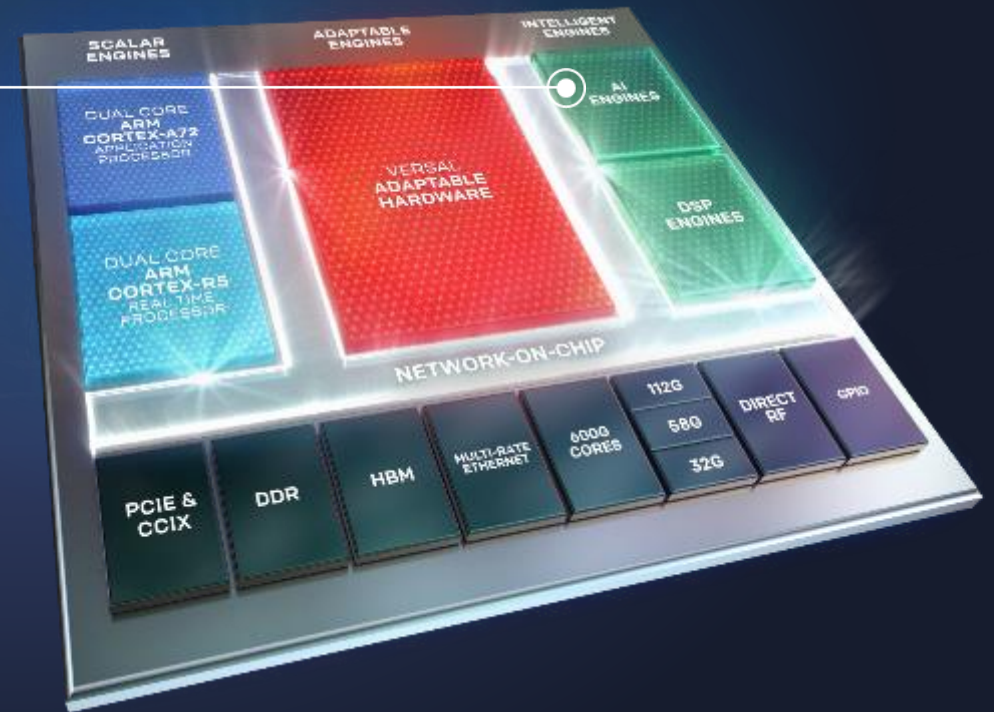
High throughput, low latency, and power efficient  
Ideal for AI inference and advanced signal processing

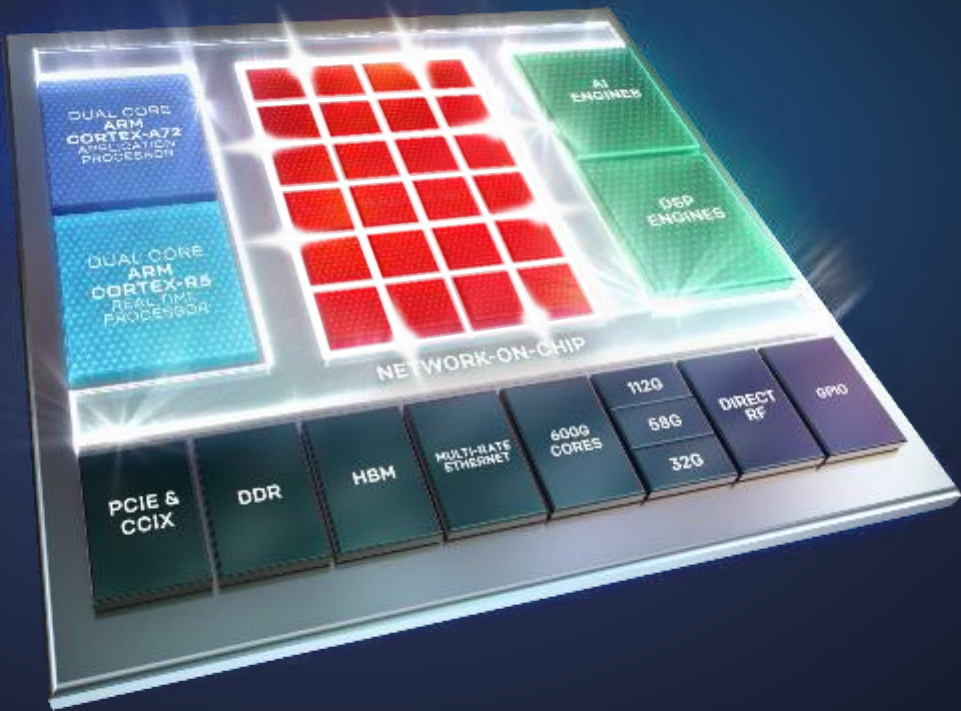




# AI Engines

Optimized for AI Inference and  
Advanced Signal Processing Workloads





# Network-on-Chip (NoC)

## Ease of Use

Inherently software programmable  
Available at boot, no place-and-route required

## High Bandwidth and Low Latency

Multi-terabit/sec throughput  
Guaranteed QoS

## Power Efficiency

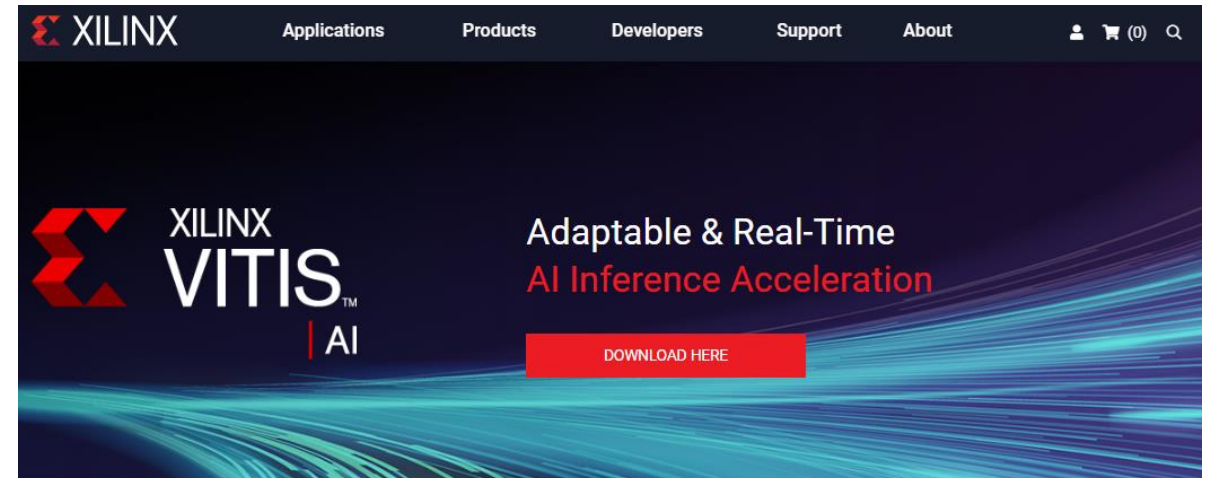
8X power efficiency vs. soft implementations  
Arbitration across heterogeneous engines



# Xilinx University Program

- ▶ Donation program
- ▶ Training materials
- ▶ Request tutorials

cathal.mccabe@xilinx.com



#### XILINX POWERS BAIDU'S ACU

Xilinx Powers Baidu's Production-Ready ACU-Advanced Platform for Automated Valet Parking

[Read Now >](#) [View All News >](#)



#### VITIS AI AVAILABLE FOR DOWNLOAD

Unleash the full potential of AI inference acceleration on Xilinx edge devices and Alveo accelerator cards.

[Read Now >](#) [View All News >](#)



#### POWERING ADVANCED ADAS & AD

Automotive-qualified adaptive portfolio now scales from edge sensors to complex domain controllers.

[Read Now >](#) [View All News >](#)

[www.Xilinx.com/university](http://www.Xilinx.com/university)



# Next generation compute efficiency with Xilinx FPGAs and the new Versal ACAP

- ▶ Try the new Vitis software for platform design free
- ▶ Test drive Alveo – production ready accelerator cards
- ▶ Next generation Versal ACAP



Performance



Data rates



Power



Cost



Compute density



Machine Learning



Adaptability



Cloud scalability



---

# Thank You



# Xilinx Mission

---

**Building the Adaptable,  
Intelligent World**

