



XILINX.  
ALVEO™

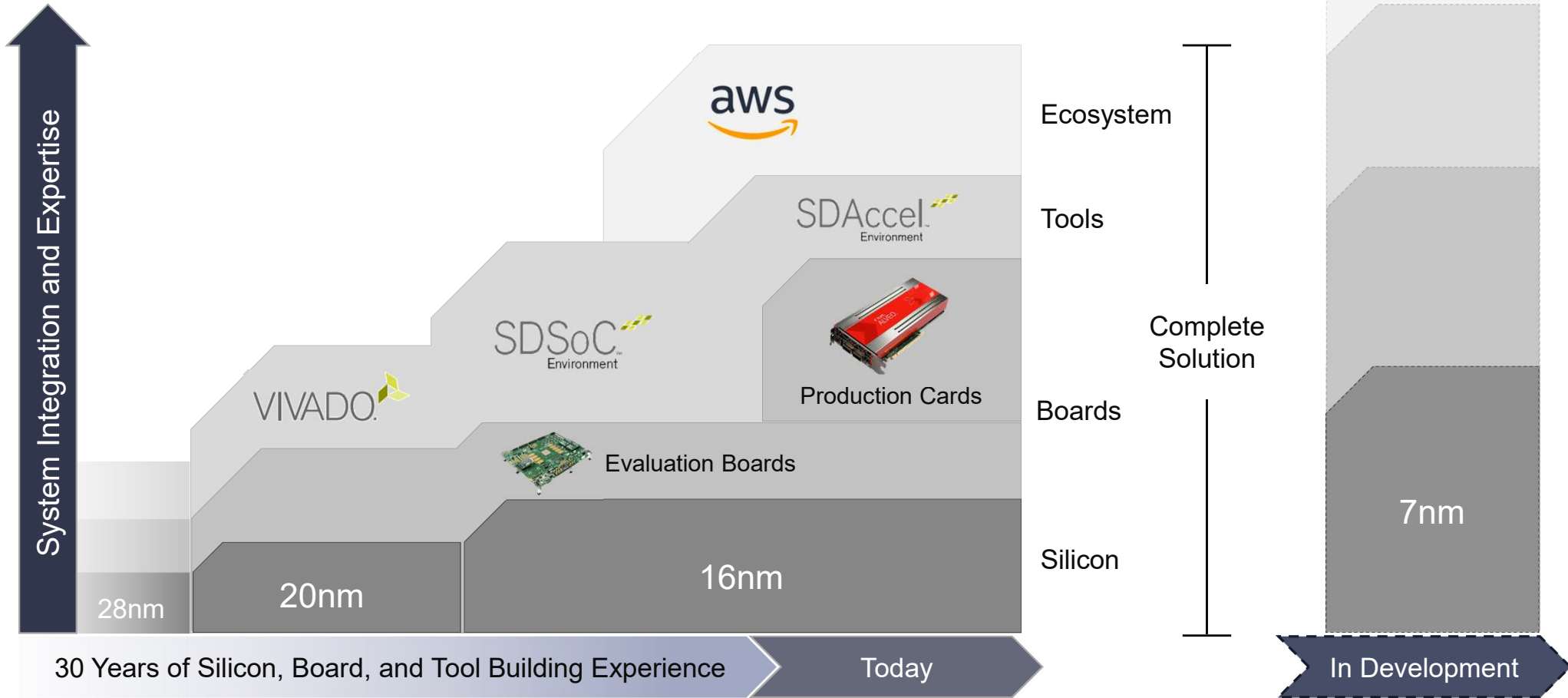
XILINX®

# Adaptable Machine Learning with Alveo Data Center Acceleration Cards

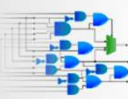
Yao Fu  
System Architect – Data Center Acceleration




# The Journey to a Complete Adaptable Solution




# Alveo – Breathe New Life into Your Data Center



**16nm  
UltraScale™ Architecture**




**Cloud Deployed**




**Off-Chip Memory Support**

- Max Capacity: 64GB
- Max Bandwidth: 77GB/s




**Cloud ↔ On-Premise Mobility**




**Internal SRAM**

- Max Capacity: 54MB
- Max Bandwidth: 38TB/s




**Ecosystem of Applications**

- Many available today
- More on the way



**PCIe Gen3x16**



**Server OEM Support**

- Major OEMs in Qualification



**Accelerate Any Application**

- IDE for compiling, debugging, profiling
- Supports C/C++, RTL, and OpenCL

# Alveo Accelerator Card Value Proposition



## Fast

*Highest Performance*

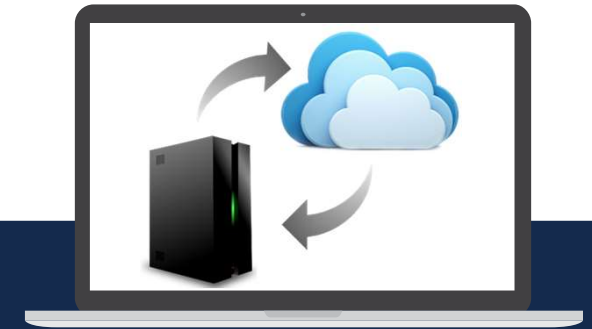
- Faster than CPUs & GPUs
- Latency advantage over GPUs



## Adaptable

*Accelerate Any Workload*

- Optimize for any workload
- Adapt to changing algorithms



## Accessible

*Cloud ↔ On-Premises Mobility*

- Deploy in the cloud or on-premises
- Applications available now

# Expanding Accelerator Card Portfolio



# Accelerator Cards That Fit Your Performance Needs



## Alveo U200

- 18.6 Peak INT8 TOPs
- 77GB/s DDR Memory Bandwidth
- 31TB/s Internal SRAM Bandwidth
- 892,000 LUTs

[Buy Now](#)

[Product Brief >](#)



## Alveo U250

- 33.3 Peak INT8 TOPs
- 77GB/s DDR Memory Bandwidth
- 38TB/s Internal SRAM Bandwidth
- 1,341,000 LUTs

[Buy Now](#)

[Product Brief >](#)



## Alveo U280

- 24.5 Peak INT8 TOPs
- 460GB/s HBM2 Memory Bandwidth
- 30TB/s Internal SRAM Bandwidth
- 1,079,000 LUTs

[Learn More](#)

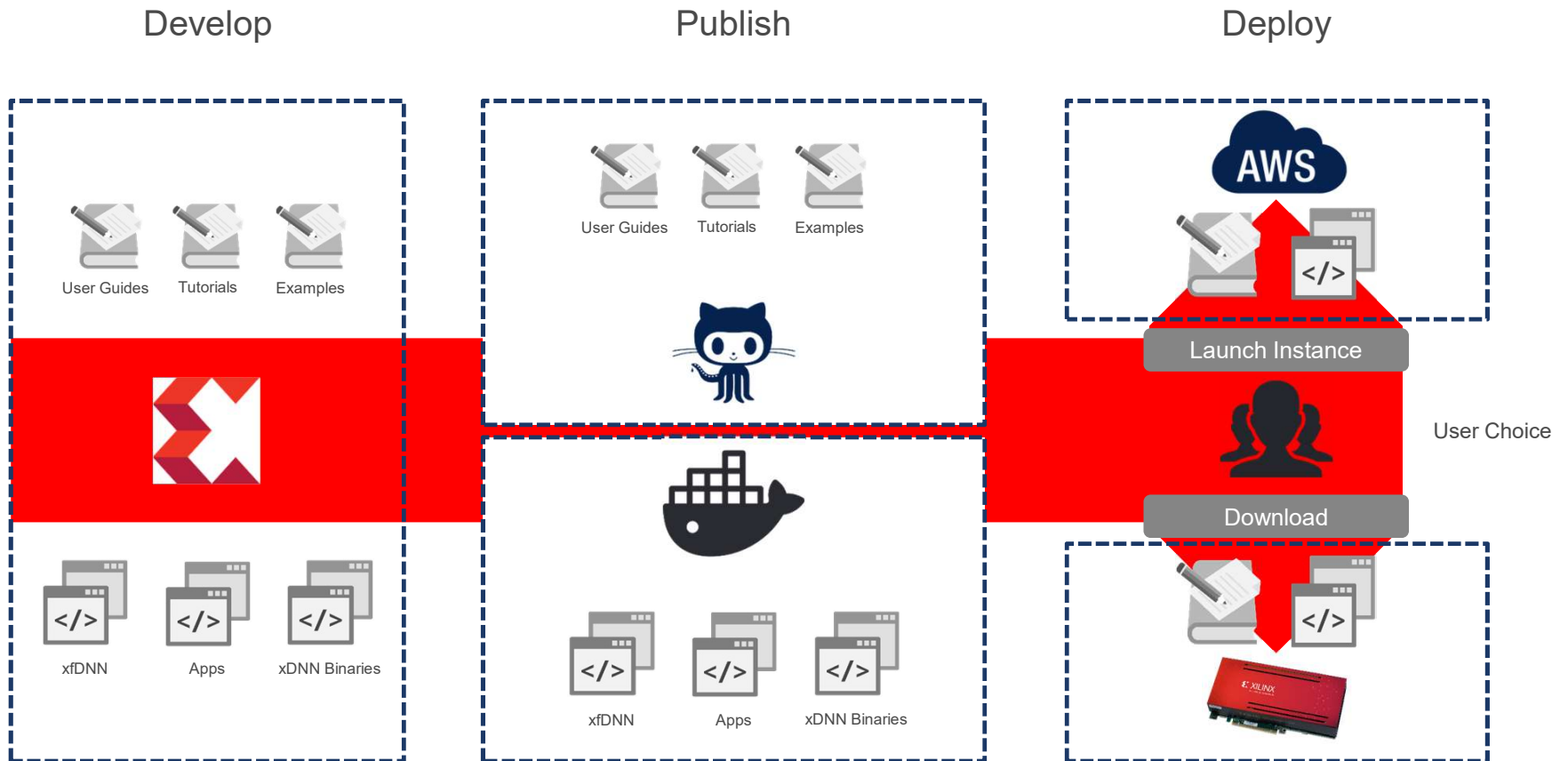
*Coming Soon*

# Adaptable Machine Learning with Alveo



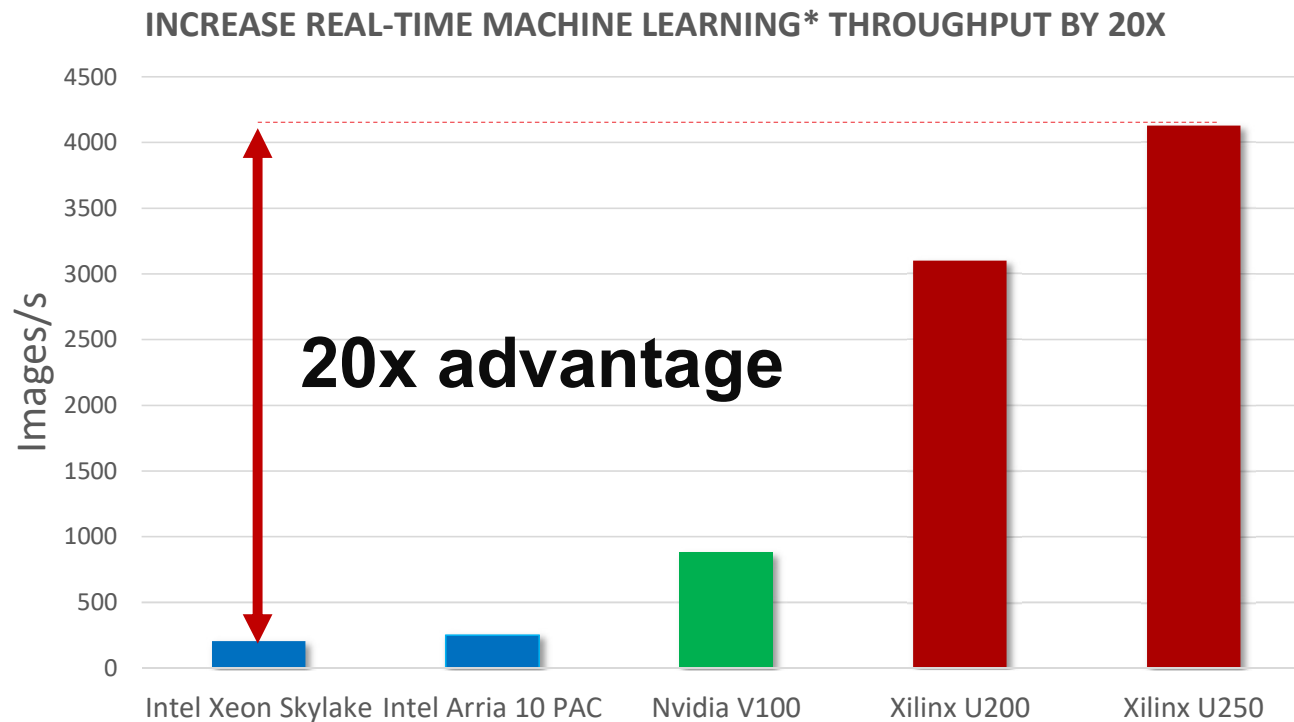


# Accessible Unified Simple User Experience from Cloud to Alveo



# Fast

## *Advantages in Machine Learning Inference*



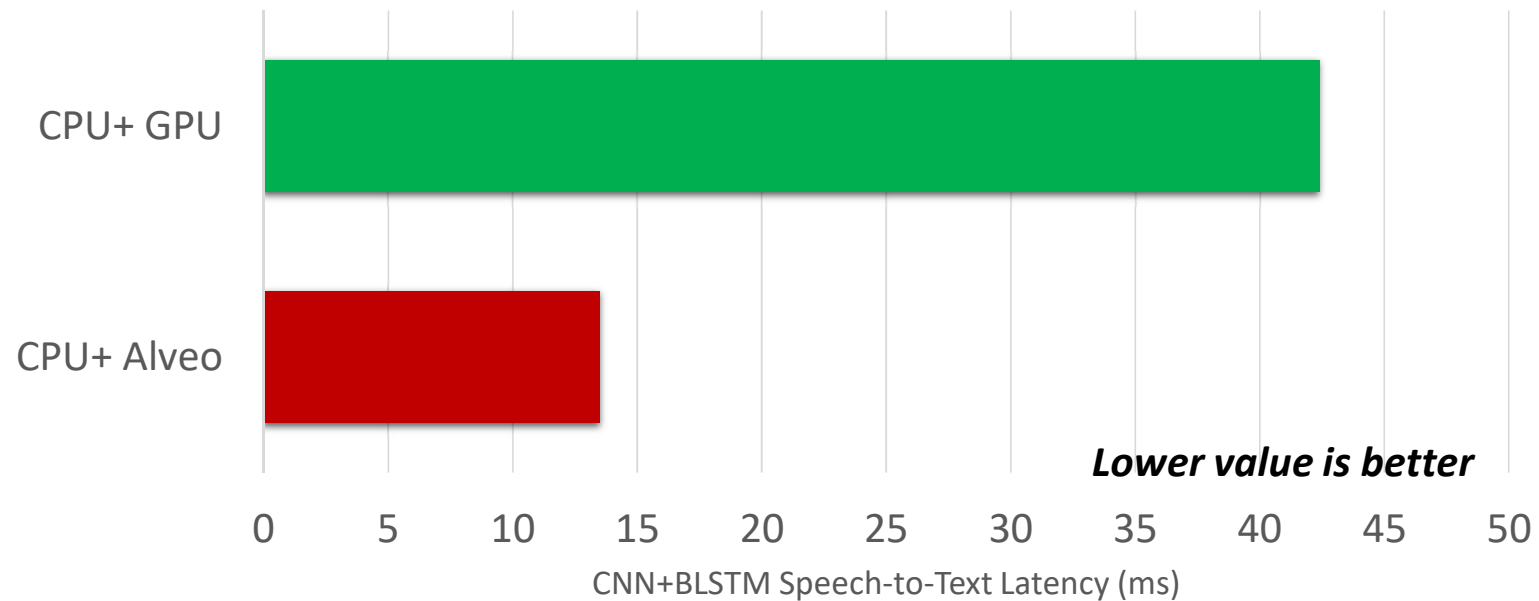
\* Source: [Accelerating DNNs with Xilinx Alveo Accelerator Cards White Paper](#)

# Fast

## *Advantages in Latency*

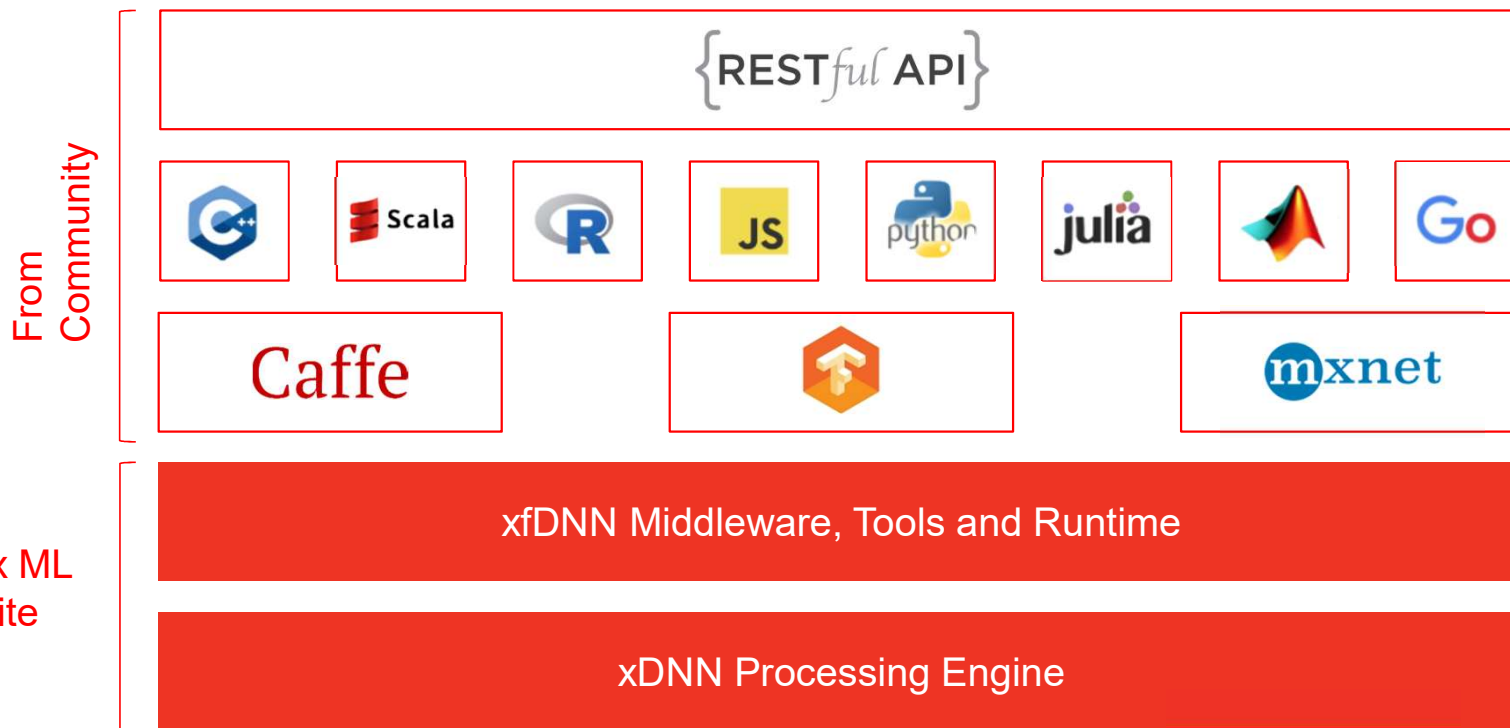


### REDUCE ML INFERENCE LATENCY BY 3X



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

# Accessible Xilinx Machine Learning Suite



- 4x higher real time inference over GPUs
- Highest Perf/Watt
- Easy to Use ML Frameworks and APIs



NIMBIX



Xilinx U200, U250

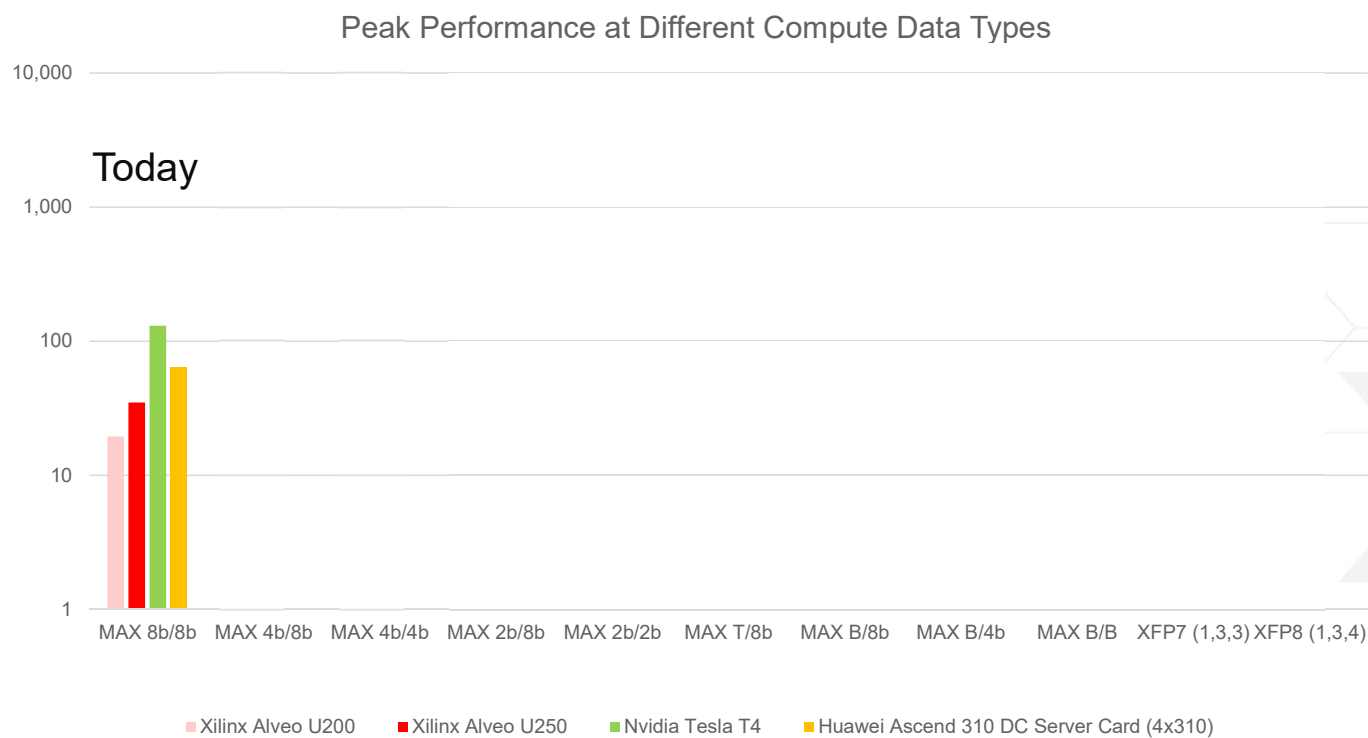
<https://www.xilinx.com/ml>

<https://github.com/xilinx/ml-suite>



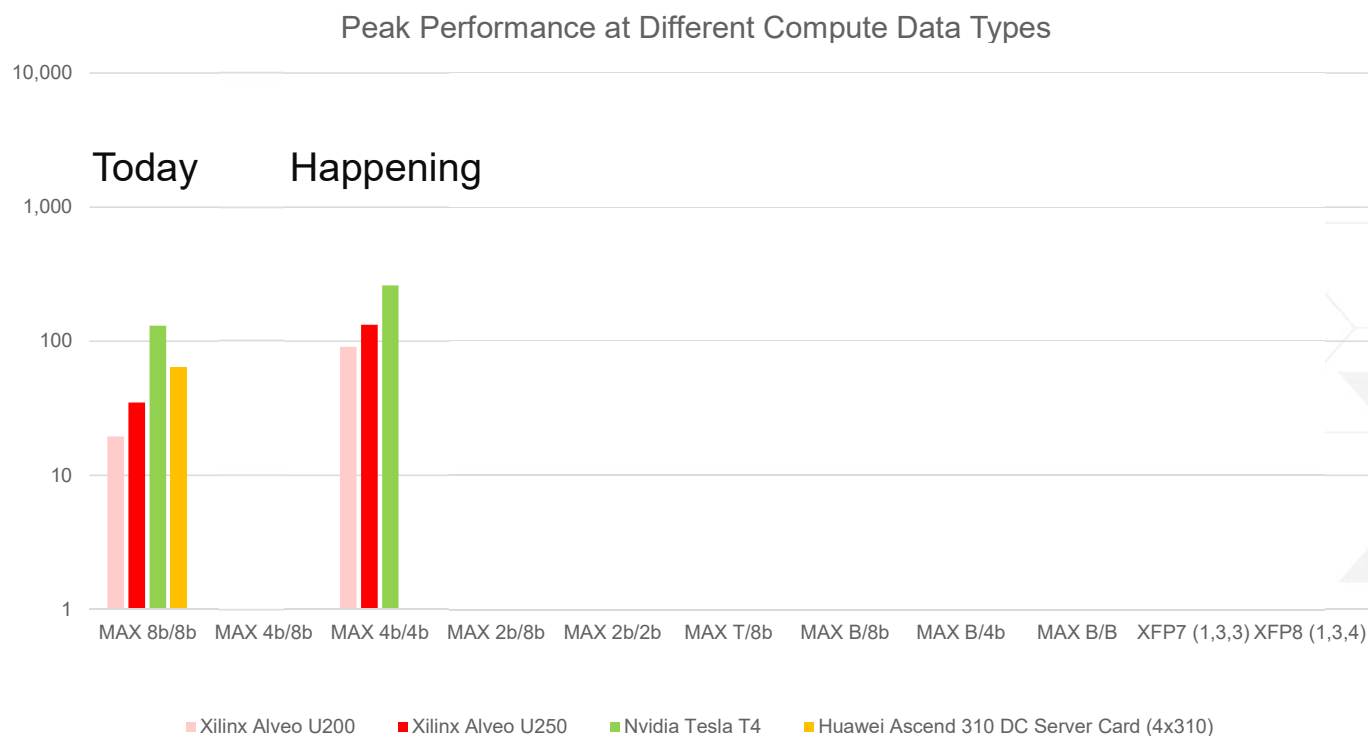
# Adaptable - Compute Capability

> INT8 is becoming standard for inference devices



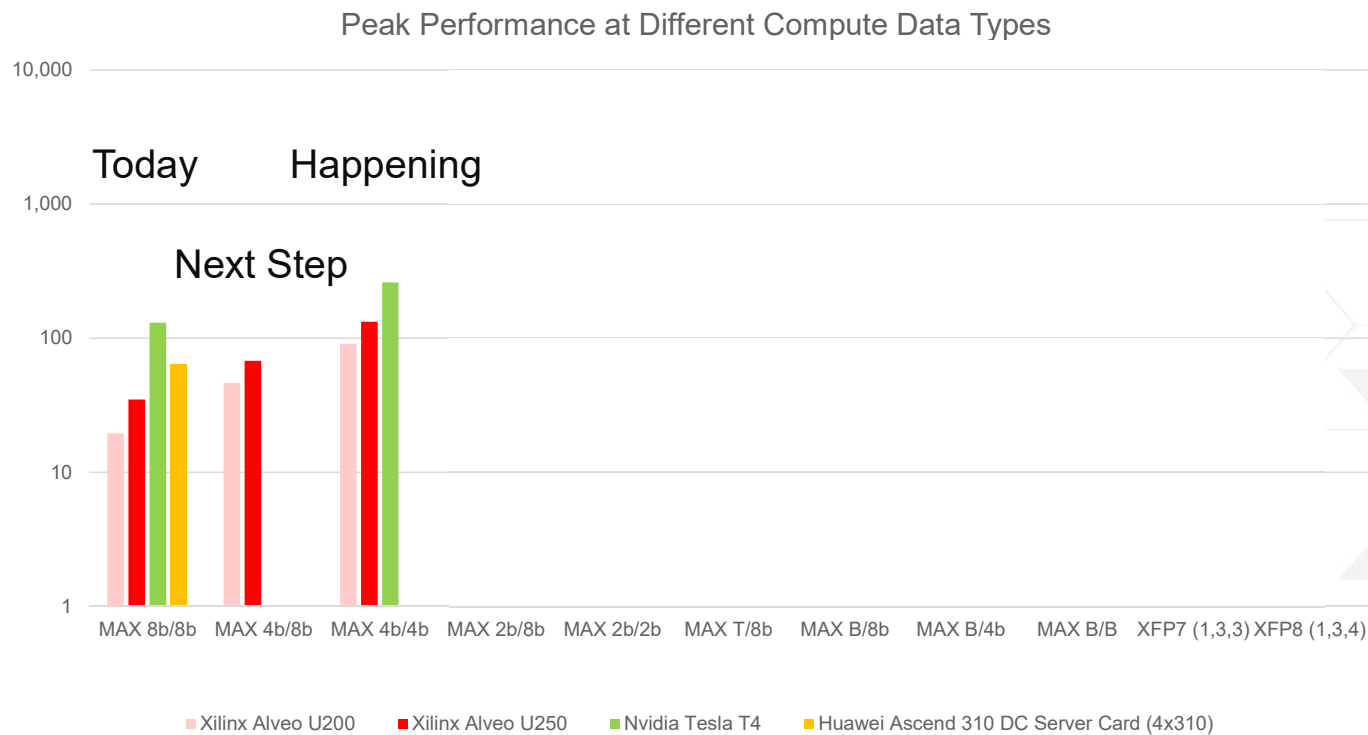
# Adaptable - Compute Capability

> Trends are moving toward lower precision, e.g. INT4 is emerging



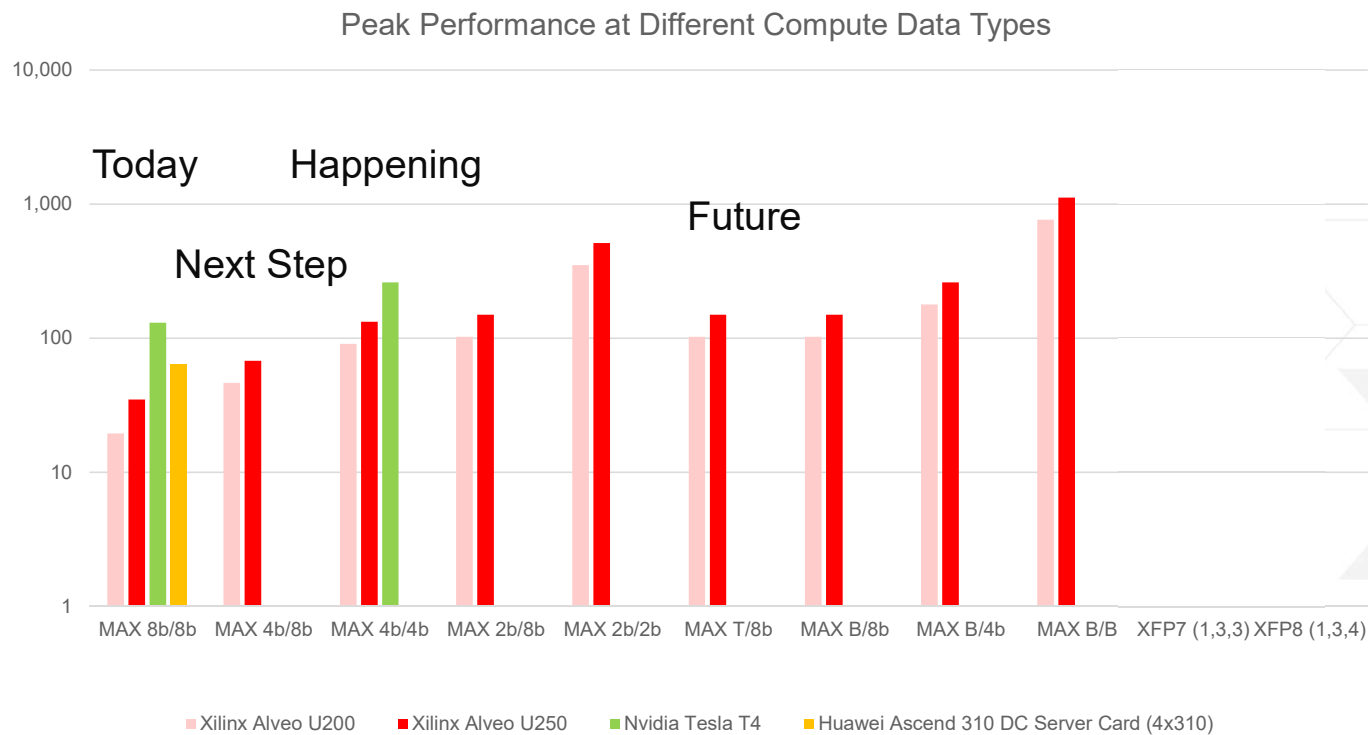
# Adaptable - Compute Capability

- > Mixed weight and activation precisions provide additional optimization for inference while keeping high accuracy



# Adaptable - Compute Capability

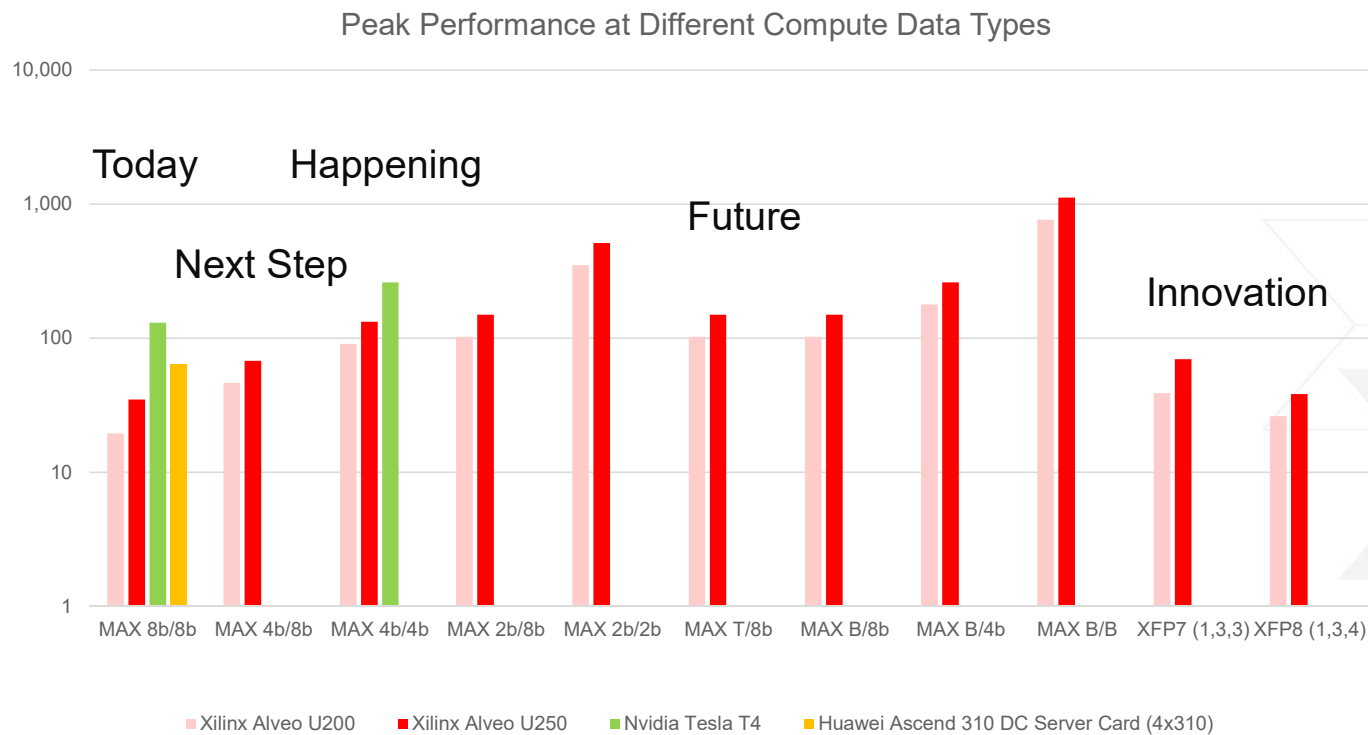
- > The future we see: each model will have its own mix of datatype for optimal efficiency and accuracy





# Adaptable - Compute Capability

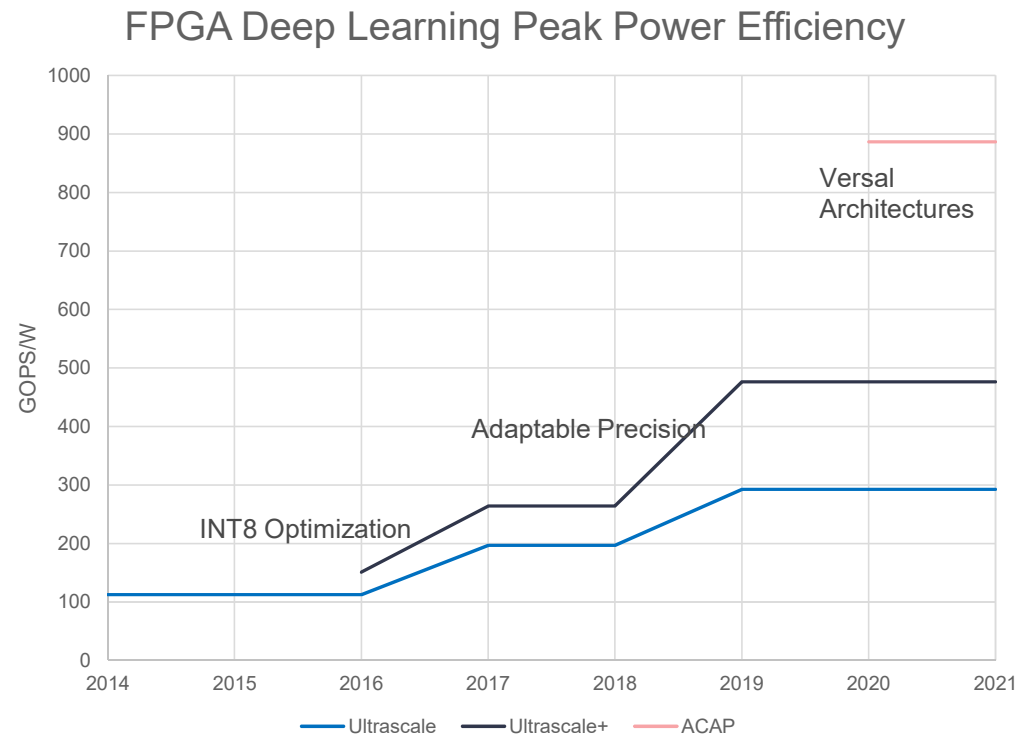
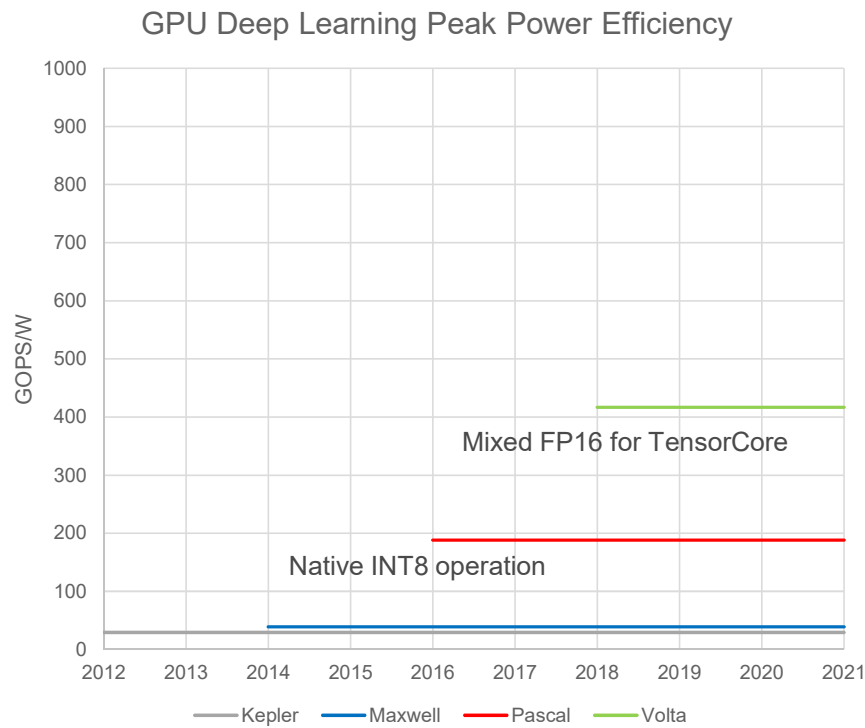
- > Alveo: Most future proof architecture
- > Scalable performance for any data type



# Adaptable - Break Through on Peak Performance

➤ GPU: Introduce new architectures and silicon

➤ Xilinx: Adapt the break through of emerging domain knowledge

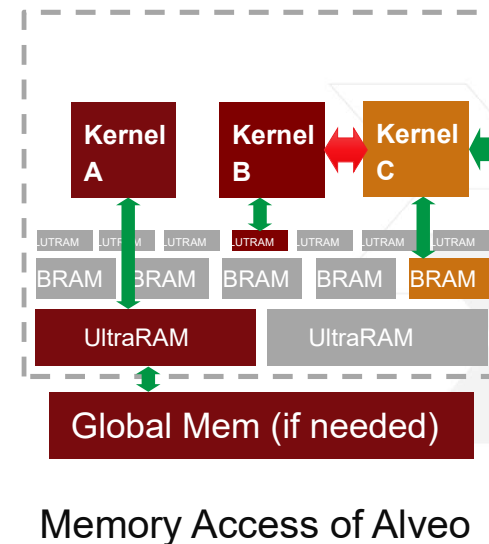
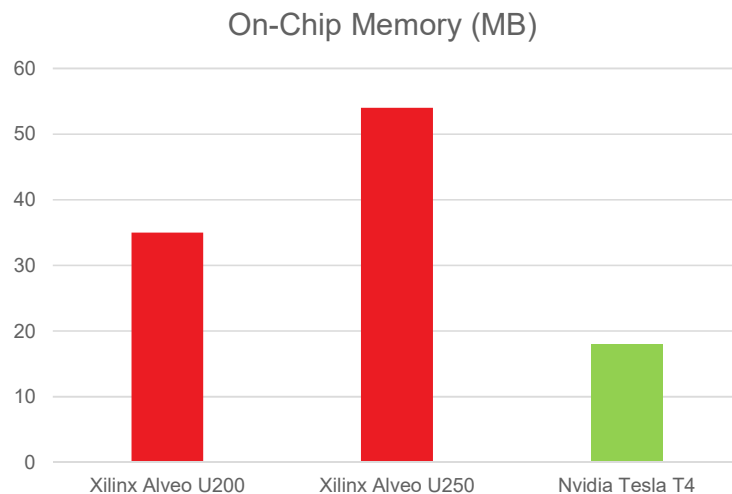


# Adaptable - On-chip memory

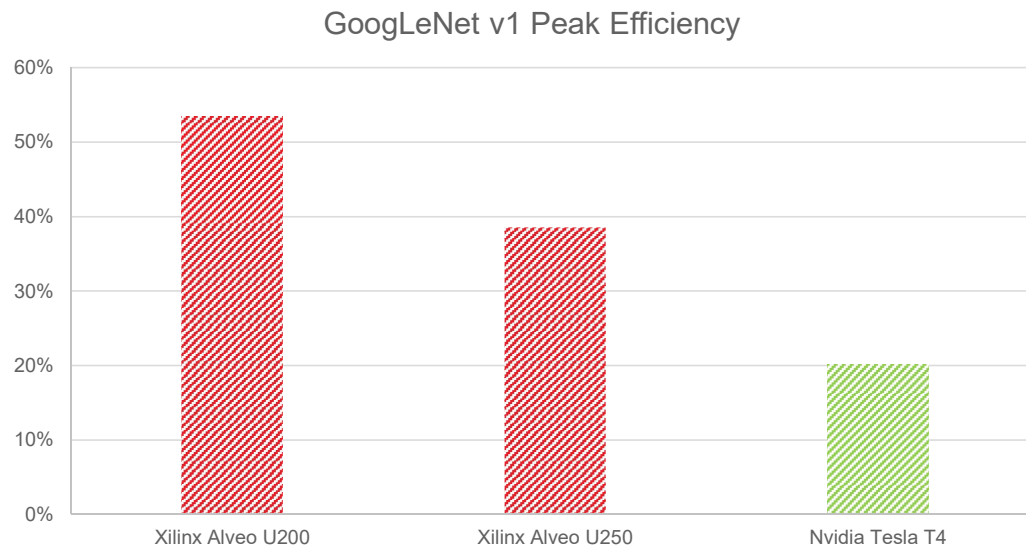
The critical asset on-chip to assist and feed the compute

Store parameters, buffer intermediate activations, and move data around

- > Alveo: Highest on-chip memory capacity
- > Alveo: Most adaptable memory architecture

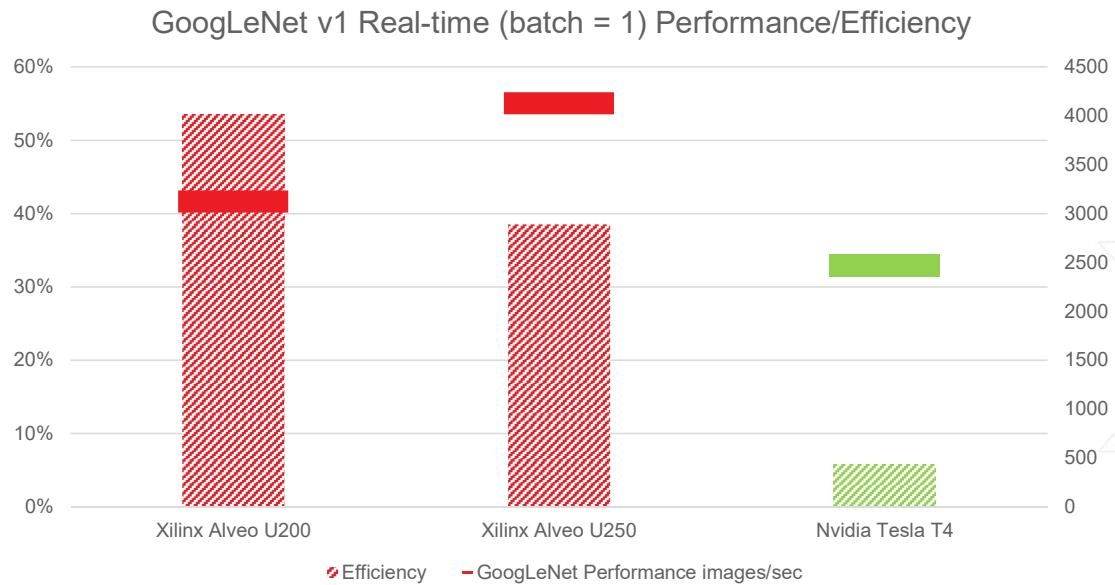


# Adaptable - Neural Network Inference Efficiency



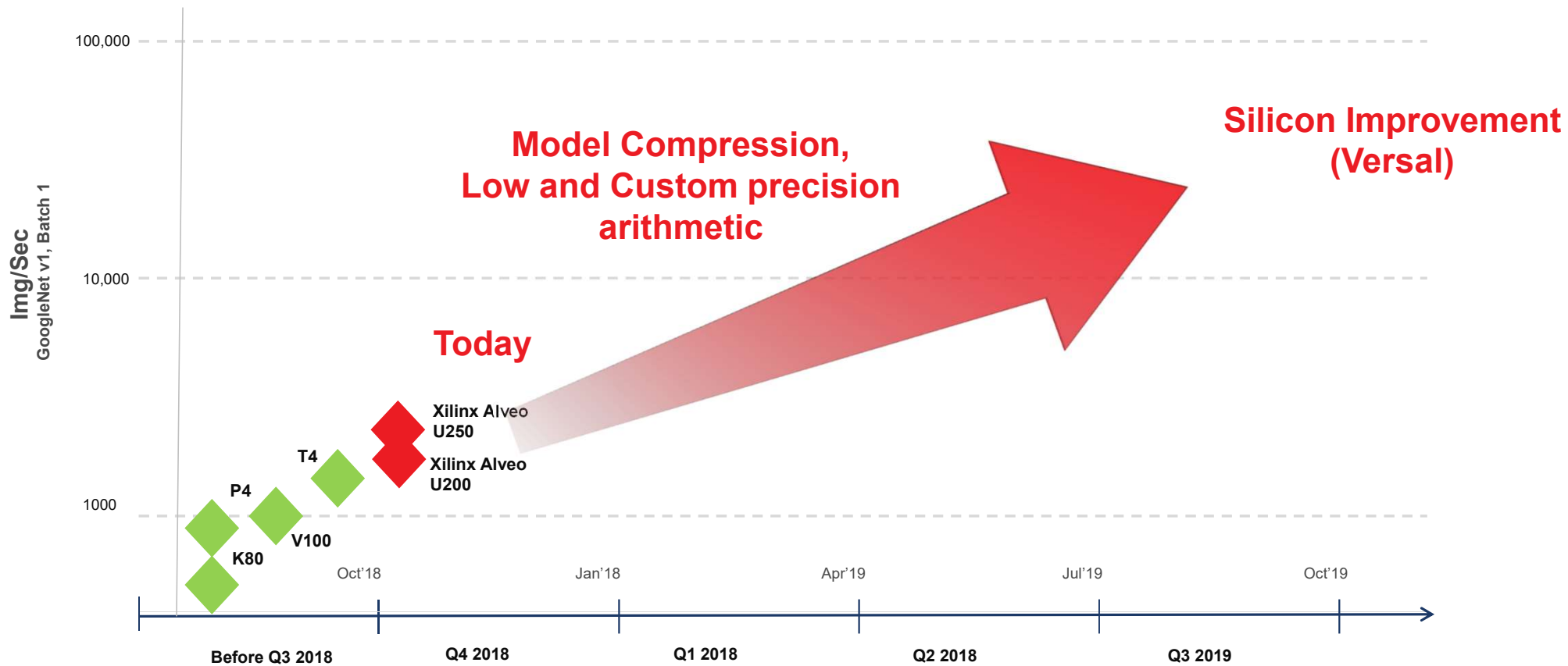
\*T4 efficiency assumes 2x power efficiency improvement vs. P4 as claimed in Nvidia whitepaper: <https://www.nvidia.com/content/dam/en-zz/Solutions/design-visualization/technologies/turing-architecture/NVIDIA-Turing-Architecture-Whitepaper.pdf>

# Fast Real-time Performance



\*T4 Performance and efficiency assumes 2x power efficiency improvement vs. P4 as claimed in Nvidia whitepaper:  
<https://www.nvidia.com/content/dam/en-zz/Solutions/design-visualization/technologies/turing-architecture/NVIDIA-Turing-Architecture-Whitepaper.pdf>

# Xilinx Performance Roadmap



# Beyond Machine Learning



# Fast

## *Advantages Across Many Workload Types*



Database



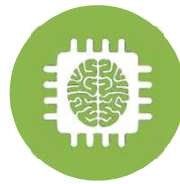
90x

Financial



89x

Machine Learning



20x

Video



12x

HPC & Life Sciences



10x



# Application Adaptability

Strong Acceleration  
Some Acceleration  
No Acceleration

Machine Learning

Video

Database

HPC & Life Sciences

Financial



FPGA



GPU



ASIC

# It's All About the Applications



## Database



## Machine Learning



## Video



## Financial



## HPE & Life Sciences



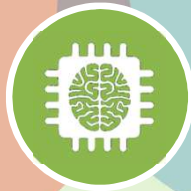
Application Ecosystem Continues to Grow

# Infuse Machine Learning with other accelerations

Database



Machine Learning



HPC & Life Sciences



Video



Financial



© Copyright 2018 Xilinx

 XILINX

# Solution Stack



Accelerated Solutions

**DEVELOPERS**

**100%**

Growth of Published Applications

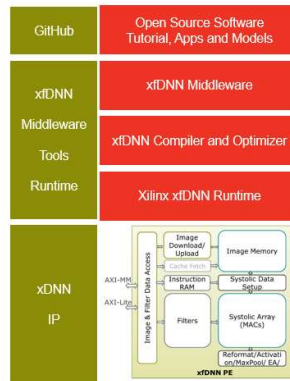
**Hundreds**

of Developers Trained

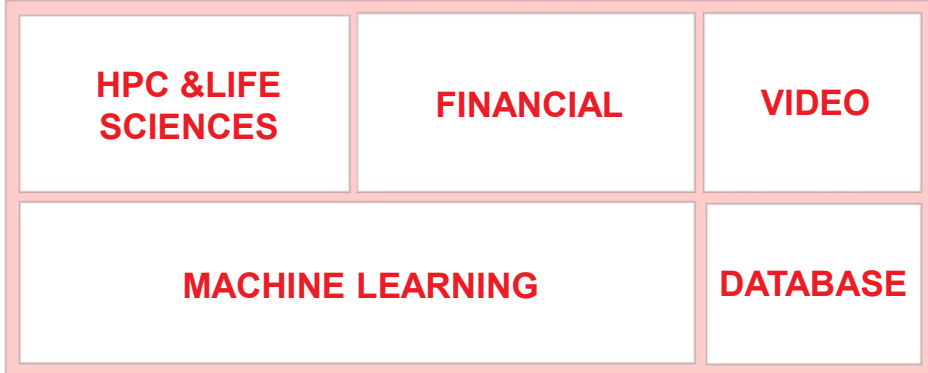
RTL, C, C++, OpenCL



Xilinx ML Suite



Framework, API, Python/Java/C++ Programmability



Solutions  
Xilinx  
ISVs

Developer Package

Platforms



Cloud

FPGA as a Service (FaaS)



On-premise

Platform

# Xilinx is Qualifying with Major Server OEMs

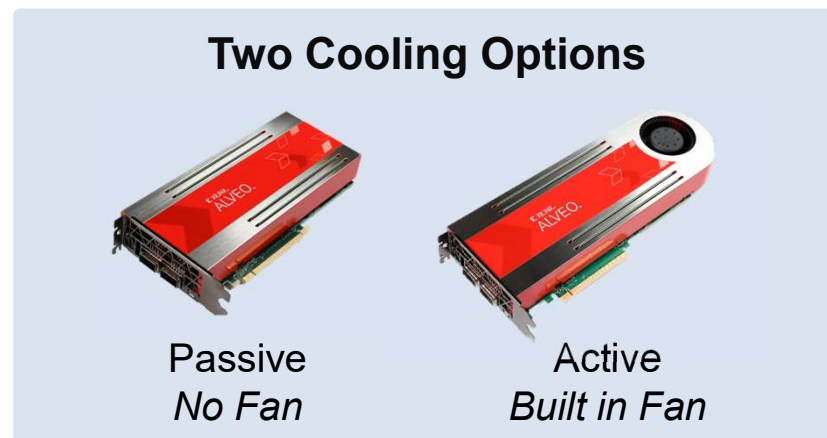


## Server Support & Qualification Strategy



# Ordering Information

<b>A</b>	<b>U###</b>	<b>P</b>	<b>64G</b>	<b>PQ</b>	<b>G</b>
<b>Alveo</b>	<b>Kit Name</b> U200 U250	<b>Cooling</b> P: Passive A: Active	<b>Memory</b> 64G: 64GB	<b>Solution Qualification</b> ESx: Engineering Sample PQ: Production Qualified	<b>RoHS Indicator</b> G: RoHS 6/6

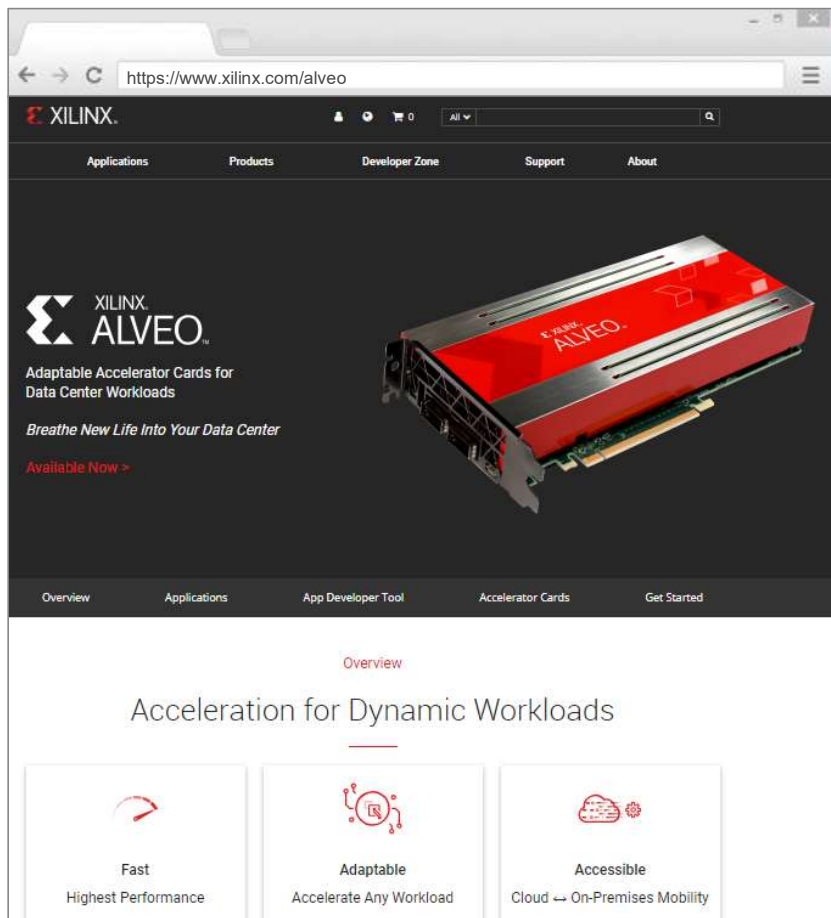


Alveo Accelerator Card	Part Number	SRP 1pc	Developer Discount 1-2pc limit <sup>1</sup>
<b>U200</b>	A-U200-A64G-PQ-G A-U200-P64G-PQ-G	\$8,995	\$3,995
<b>U250</b>	A-U250-A64G-PQ-G A-U250-P64G-PQ-G	\$12,995	\$6,495

1 – Requires qualification through [Accelerator Program](#)

- > Buy via standard quote/PO process from Xilinx or Avnet
- > Buy via eCOM (at SRP) from Xilinx, Avnet, DigiKey, EBV, or Premier Farnell

# More Information Available on Xilinx.com



## Xilinx.com

[Product Brief](#)

[Product Selection Guide](#)

[Getting Started Guide](#)

[Data Sheet](#)

[ML Solution Brief](#)

[SDAccel Solution Brief](#)

[ABR Transcoding Solution Brief](#)

[Accelerating DNNs with Alveo White Paper](#)

[Applications Directory](#)

**Adaptable.**  
**Intelligent.**

