

# Alveo Overview for Hyperscale and HPC Applications

**Viraj Paropkari**

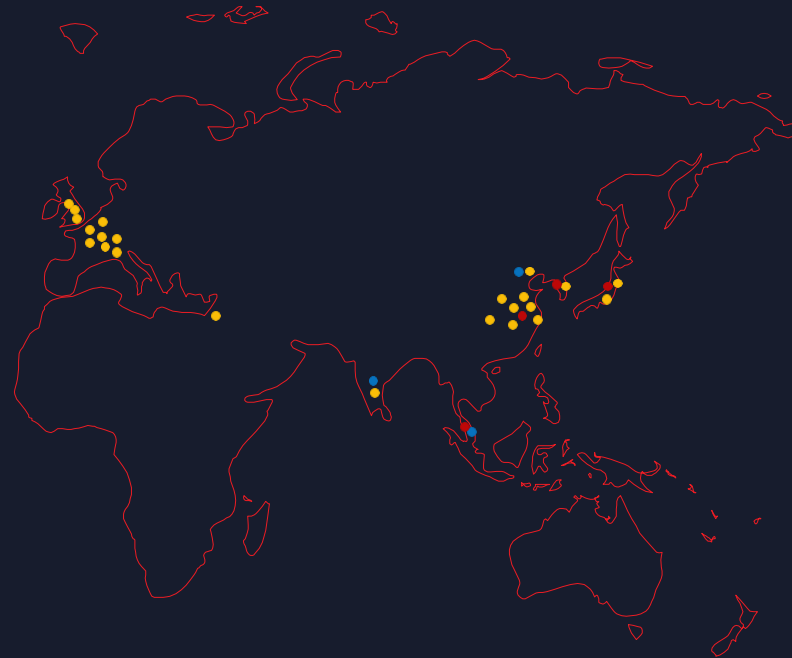
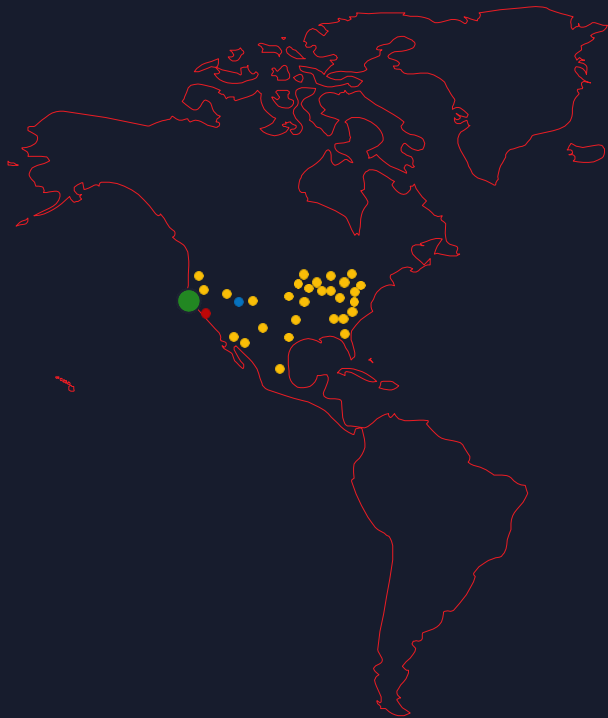
**Senior Manager, Global Data Center Marketing**

Arm HPC User Group (AHUG) 2019  
20<sup>th</sup> June 2019, Frankfurt, Germany



# Agenda

- > Data Center Focus & Strategy
- > Alveo Data Center Boards & Ecosystem  
Overview
- > Use Cases- Compute, Storage, Networking
- > Future of Computing
- > Getting Started



- Headquarters
- Research and Development
- Sales and Support



\$3.0B

Revenue



~4,400

Employees Worldwide



20K+

Customers



60+

Industry Firsts



4,000+

Patents

# Datacenter Focus & Strategy




# Data Center Opportunity



Hyperscale  
Public Cloud



Enterprise  
Private Cloud




Telco  
Cloud / Edge



High Performance  
Computing



Compute



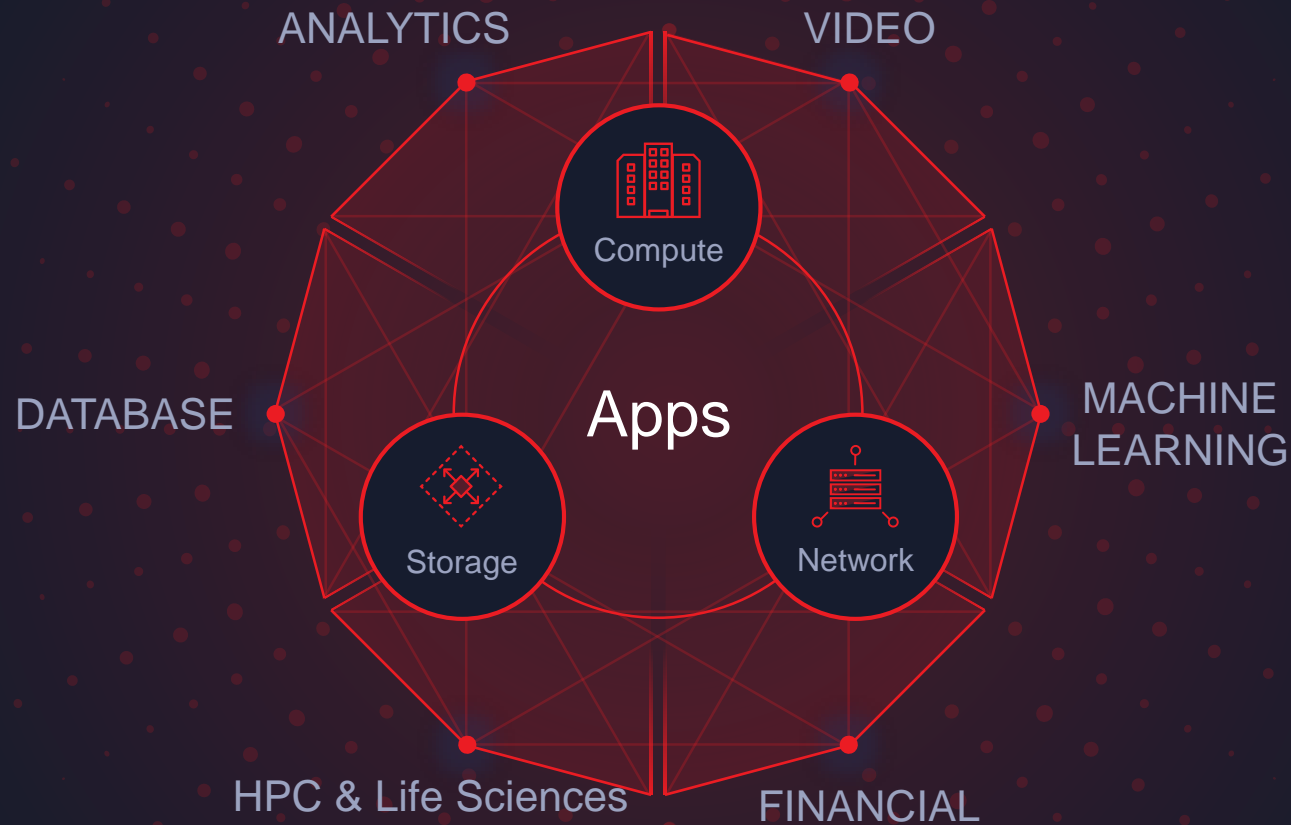
Storage



Network

- > Heterogeneous Computing post Moore's Law
- > Exponential Data Growth
- > Dawn of AI

# ➤ Era of Reconfigurable Accelerators



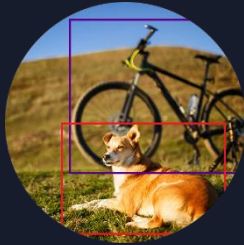
# ➤ Algorithm Diversity and Fast Evolution

## APPLICATIONS

Classification



Object Detection



Speech  
Recognition



Data  
Analytics



Recommendation  
Engine



Anomaly Detection



CNN

RNN, LSTM

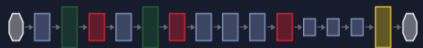
RF, LR

MLP

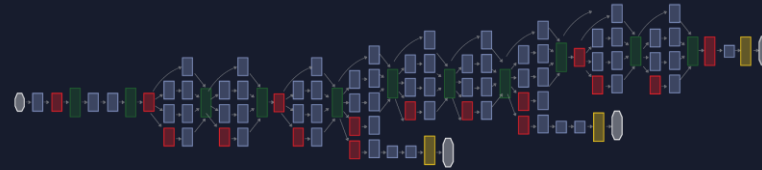
←..... Diverse AI models and Neural Networks (NN's) for a broad range of applications .....→

# ➤ New Algorithms Need New Architectures

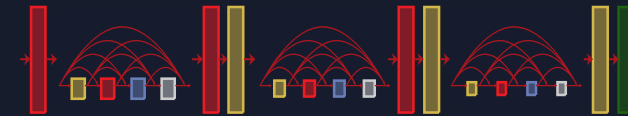
AlexNet



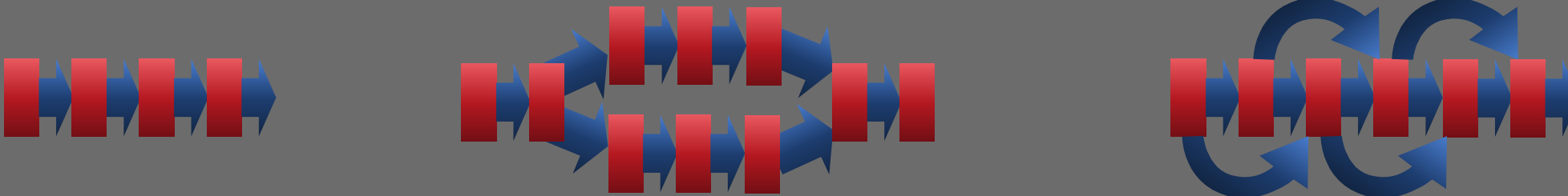
GoogLeNet



DenseNet



Highest throughput, latency, and efficiency requires different HW architecture

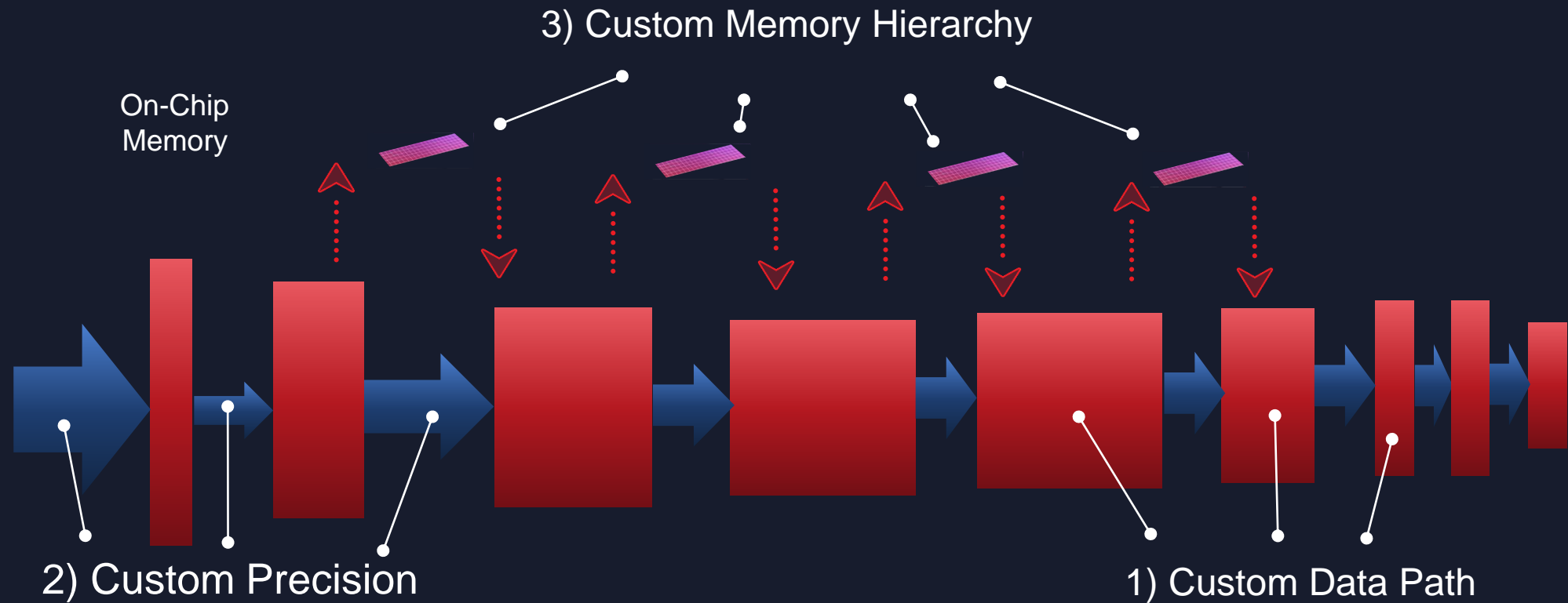




# ➤ Optimized Performance

Requires Custom Memory and Datapath

Off-Chip  
DDR



# ➤ Reconfigurable Data Center Accelerators

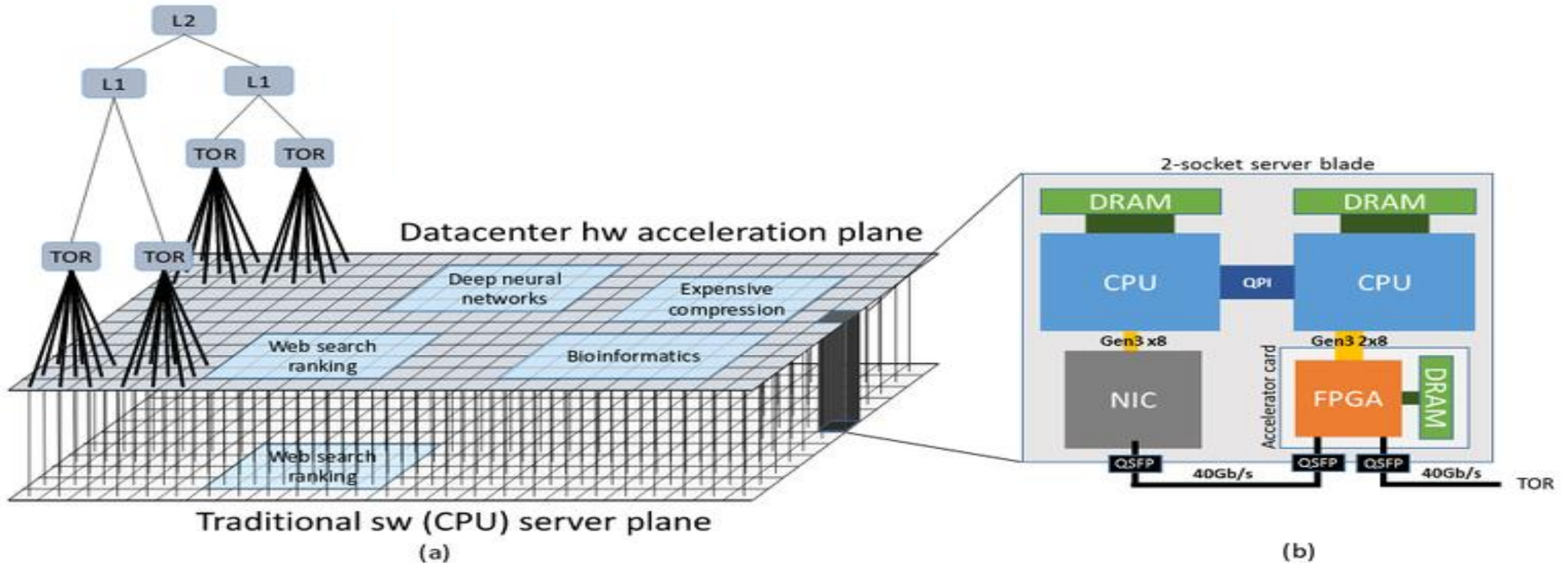


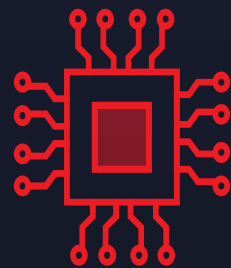
Fig. 1. (a) Decoupled Programmable Hardware Plane, (b) Server + FPGA schematic.

Tier1 HS adopting FPGA for general purpose acceleration across multiple in-house and external workloads

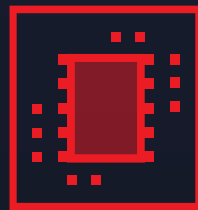
# Xilinx Transformation

From Devices to Platforms

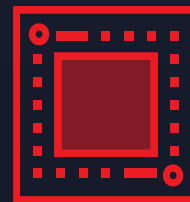
SW Programmability



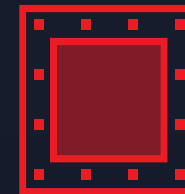
FPGA



SoC



MPSoC

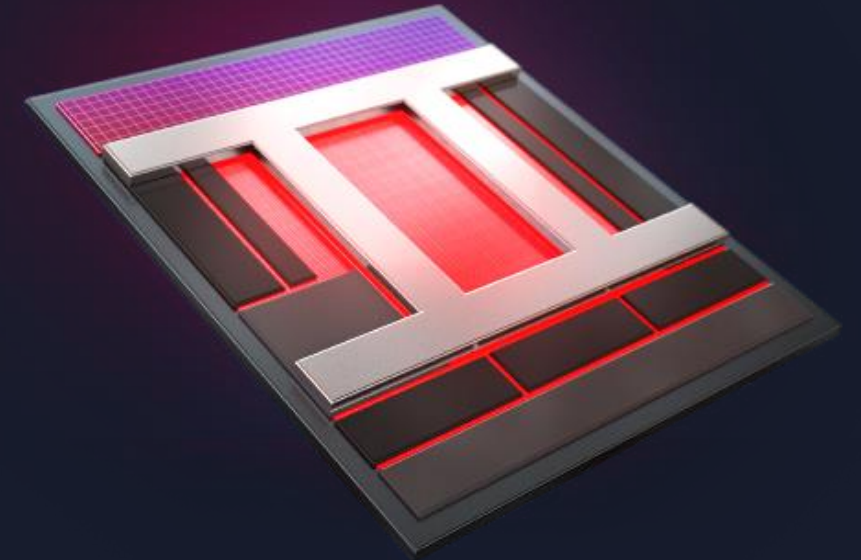


RFSoc



ACAP

Device Category



# Alveo Product & Ecosystem Overview



# 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





## U200

**892K**  
LUTs

**35MB**  
Internal SRAM  
Capacity

**31TB/s**  
Internal SRAM  
Bandwidth

**3100img/s**  
CNN Throughput\*

## U250

**1,341K**  
LUTs

**54MB**  
Internal SRAM  
Capacity

**38TB/s**  
Internal SRAM  
Bandwidth

**4100img/s**  
CNN Throughput\*

## U280

**1,079K**  
LUTs

**41MB**  
Internal SRAM  
Capacity

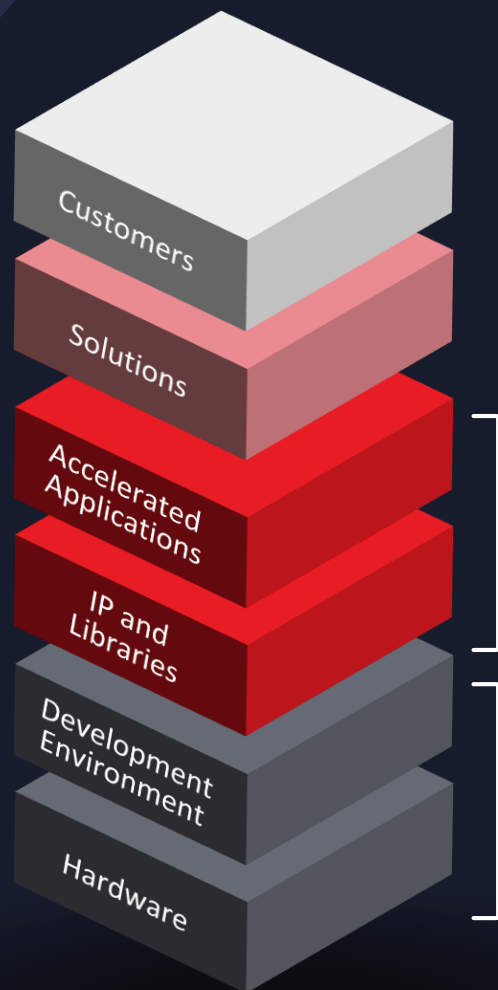
**30TB/s**  
Internal SRAM  
Bandwidth

**460GB/s**  
HBM2 Memory Bandwidth

\*Low-latency GoogLeNet v1



# ➤ ALVEO Solution Stack



End  
Customers

Solution  
Providers

App & IP  
Developers

Channel  
Partners

Data  
Analytics

Video & Image  
Processing

Machine  
Learning

Life Science  
& HPC

Financial  
Computing

Titan  
IC

NGCODEC

edico genome  
an Illumina company

MAXELER  
Technologies  
HIGHER PERFORMANCE COMPUTING

DGEPhi  
深 金 科 技

PLUNIFY

CTACCEL

DEEPOLY

falcon  
COMPUTING

Mipsology

BLACKLYNX

byte  
LAKE

X E L E R A

skreens

swarm64

ALGO-  
LOGIC

MEGH  
COMPUTING

VITESSE DATA

boon

V-NOVA

Nextera  
Video

bigstream

1ENIAC

inaccel

BigZetta  
Systems

LegUp

mle  
measuring life decisions

aws

Cloud

Tencent Cloud

NIMBIX

SUPERMICR

inspur

DELL

Alibaba Cloud

Baidu 百度

HUAWEI

CLOUD



ON-PREMISE

# ➤ Growing Ecosystem

## Data Analytics



## Life Sciences & HPC



## Video Processing



## Machine Learning



## Financial Computing



## Image Processing





# ➤ Growing Acceleration Ecosystem

**>2x**

Applications

**27 ➤ 65**

May 2018

Today

**Published Apps**

**>6x**

Partners

**84 ➤ 555**

May 2018

Today

**Accelerator Program**

**>4x**

Developers Trained

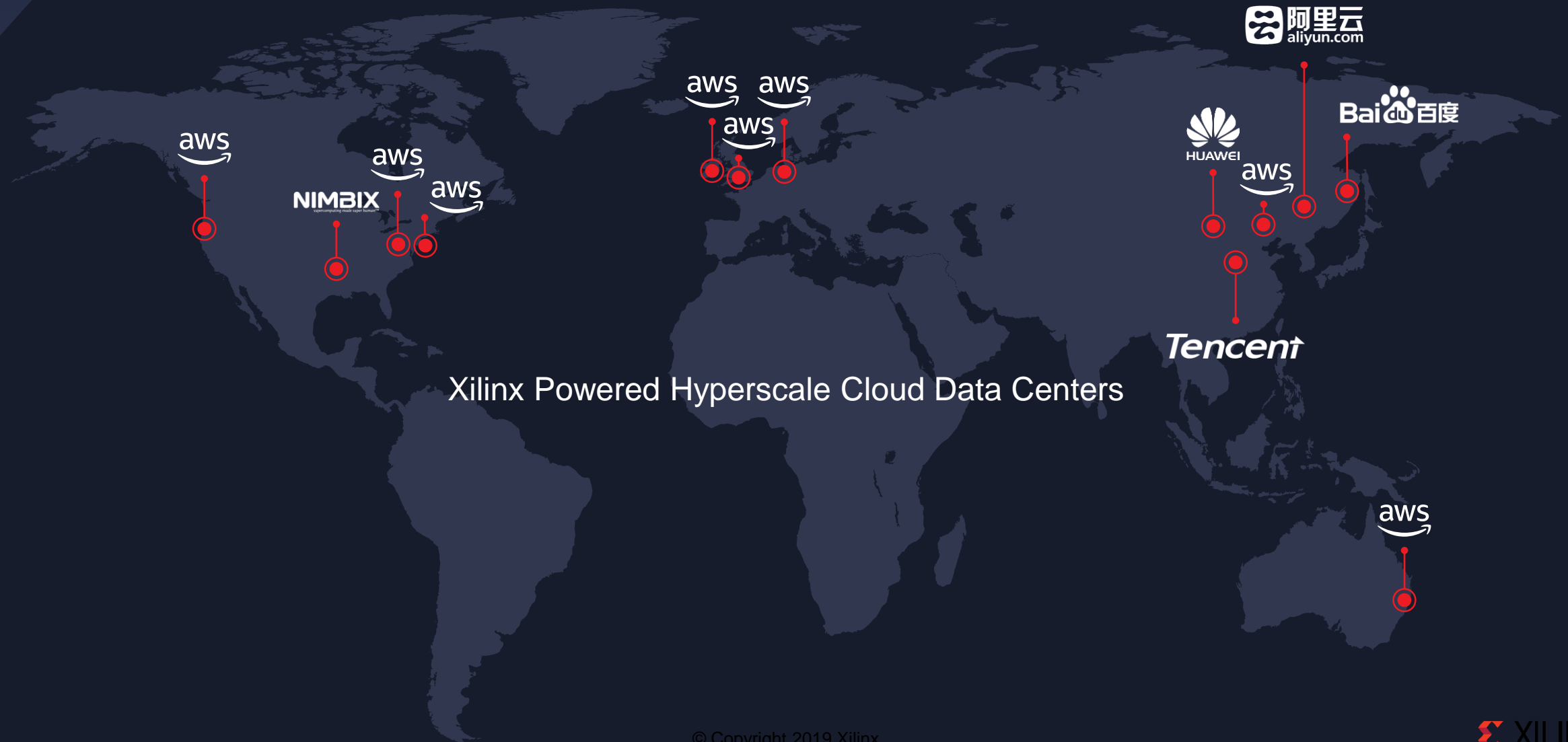
**716 ➤ 3089**

May 2018

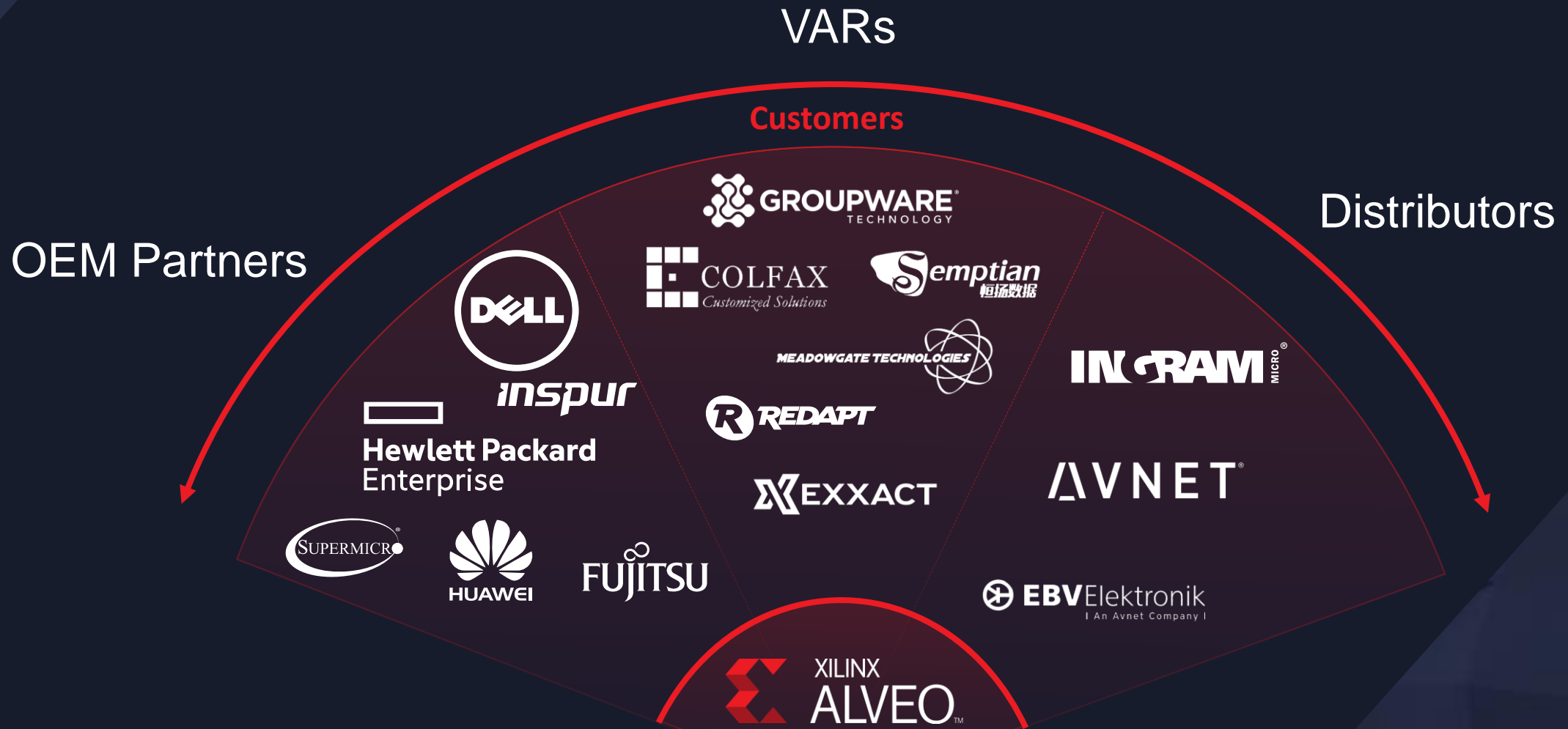
Today

**Companies + Academia**

# ➤ Growing Cloud Availability



# ➤ Growing Solution Partner Network



# Use Cases

# ➤ AI Accelerated Dark Matter Search (CERN)

Real-time ML Inference + Sensor pre-processing



Achieving 100ns Inference Latency on 150 Terabytes/Second Data Rates  
Unachievable by CPUs & GPUs

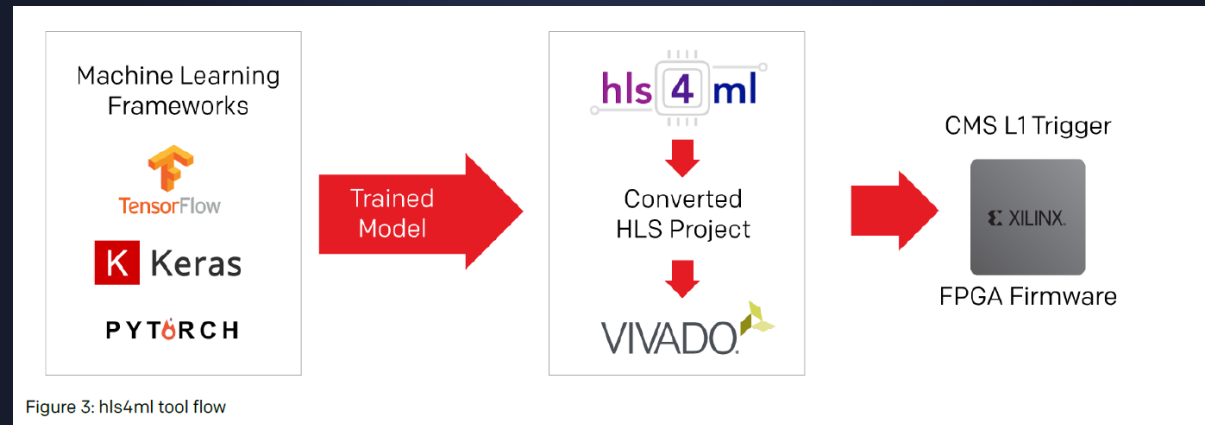
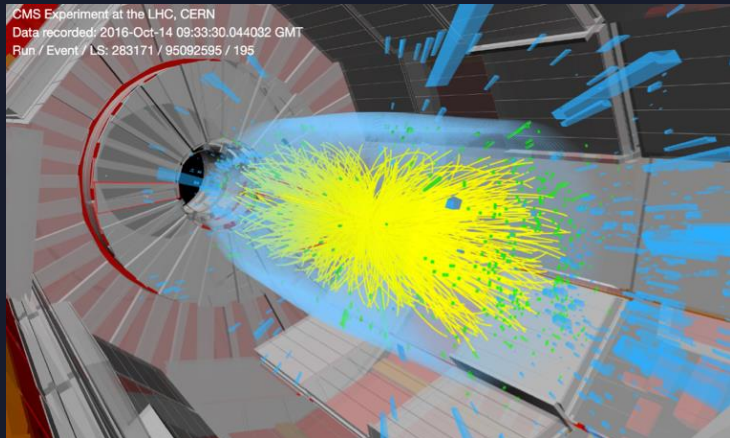


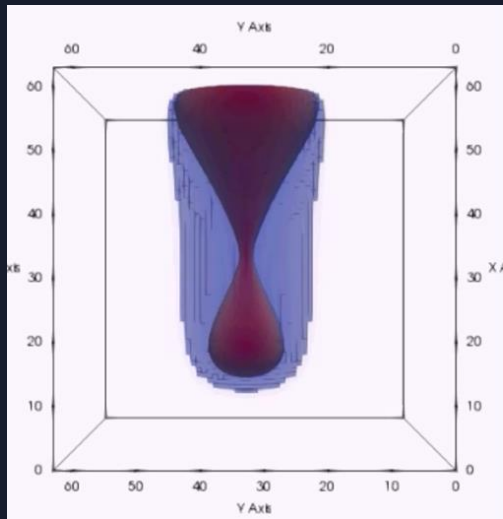
Figure 3: hls4ml tool flow

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

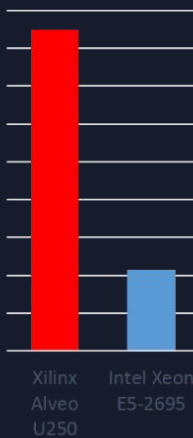
# ➤ Computational Fluid Dynamics

## ALVEO Accelerated CFD Kernels

byte  
LAKE



Performance



## Faster Time to insight, Fewer Nodes

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

# ➤ Precision Medicine

## Genomic Data Analytics

Accelerates Sequencing by 90x



illumina





# ➤ Live Video Streaming

## Video Transcoding for VP9 Live Stream

30x higher performance

25% cost reduction





# ➤ Video Processing + AI Inference

100x higher performance

Significant cost reduction

teral



# ➤ Computational Storage

**5x Speedup for  
Data-Intensive Workloads**

Compression

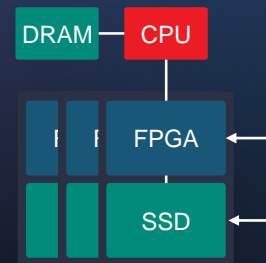
Encryption

Database query offloads

**Reduces TCO >2x**

Increased performance

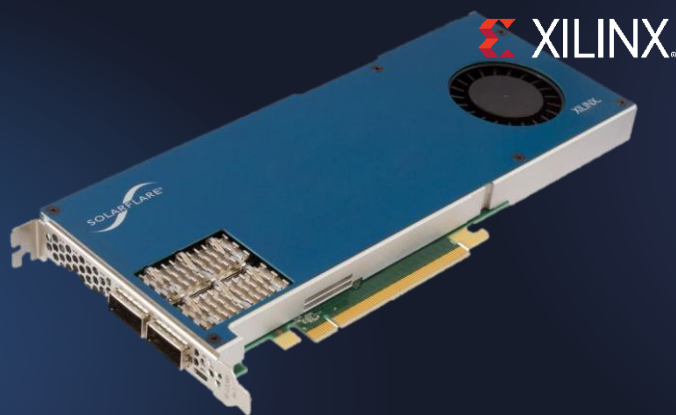
Frees up Host CPU cycles



**SAMSUNG**



# ➤ Smart Networking Acceleration



## **SolarFlare SmartNIC**

Processing >100 million packets/sec  
Dual 100G QSFP  
Under 75 watts

## **More Efficient Infrastructure**

Adaptable to support new protocols  
Extensible with programming in P4 & C/C++

# ➤ Data Center Platforms



# Future of Computing





# ➤ Expanding Xilinx Value with ACAP

Designed for SW programmability

Integrates acceleration functions

Higher performance

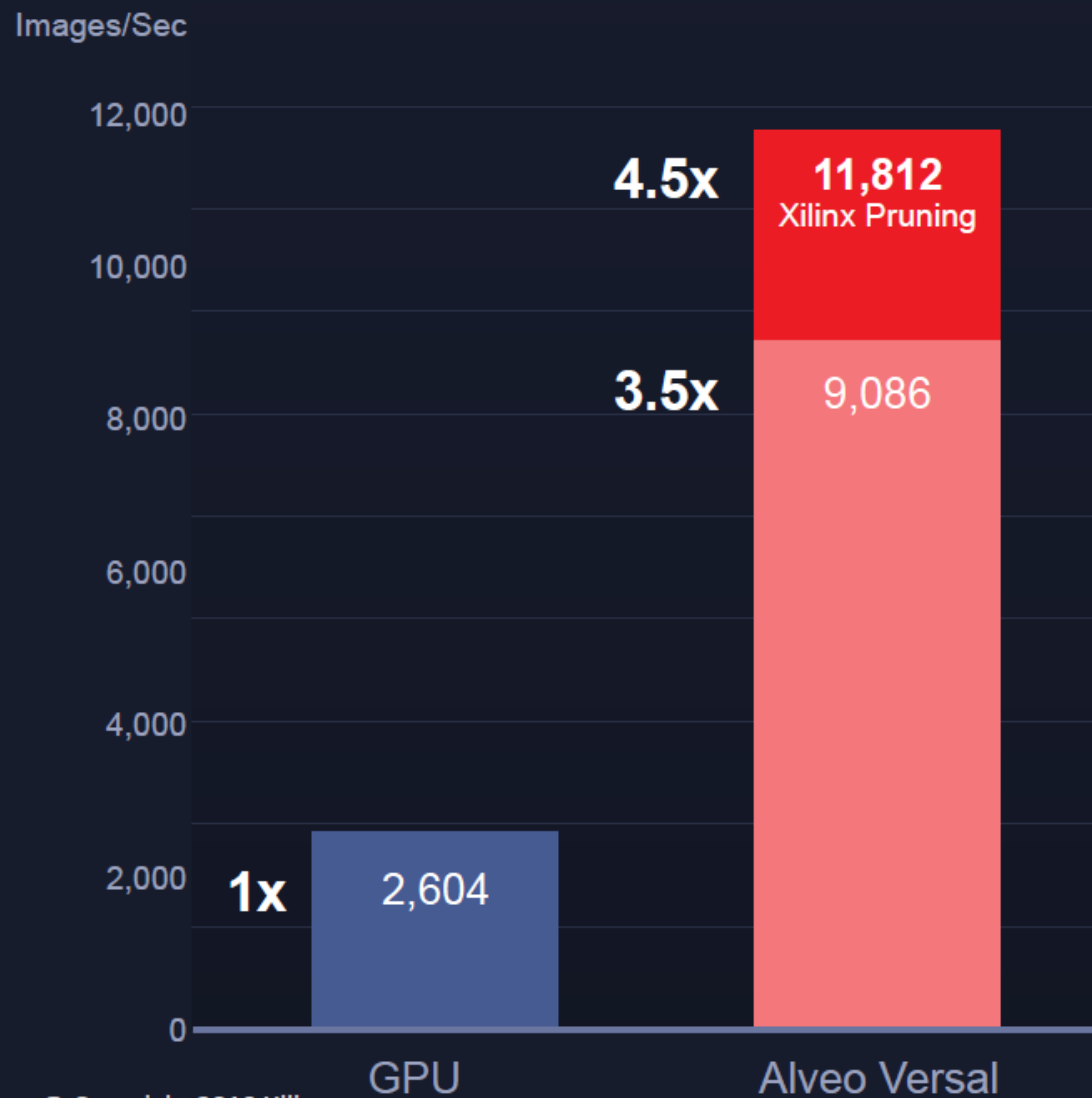
More power efficient

Now Shipping to multiple tier-1 customers.  
GA 2H 2019



# ➤ Versal AI Inference Performance

## Resnet50 Inference Performance



Sources:

GPU: Nvidia T4 TensorRT 5, Published March 2019

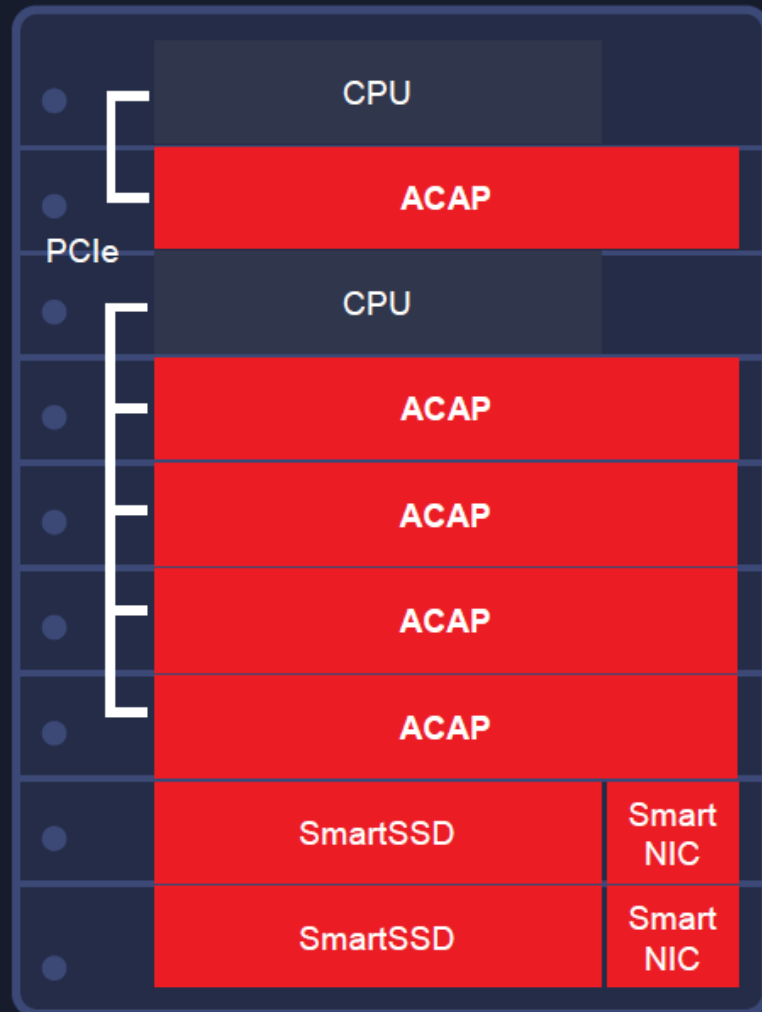
(INT8, Batch=4, 1.5ms Latency)

Alveo Versal Card, Projected (INT8, Batch=8, 1.5ms Latency)

@ Copyright 2019 Xilinx

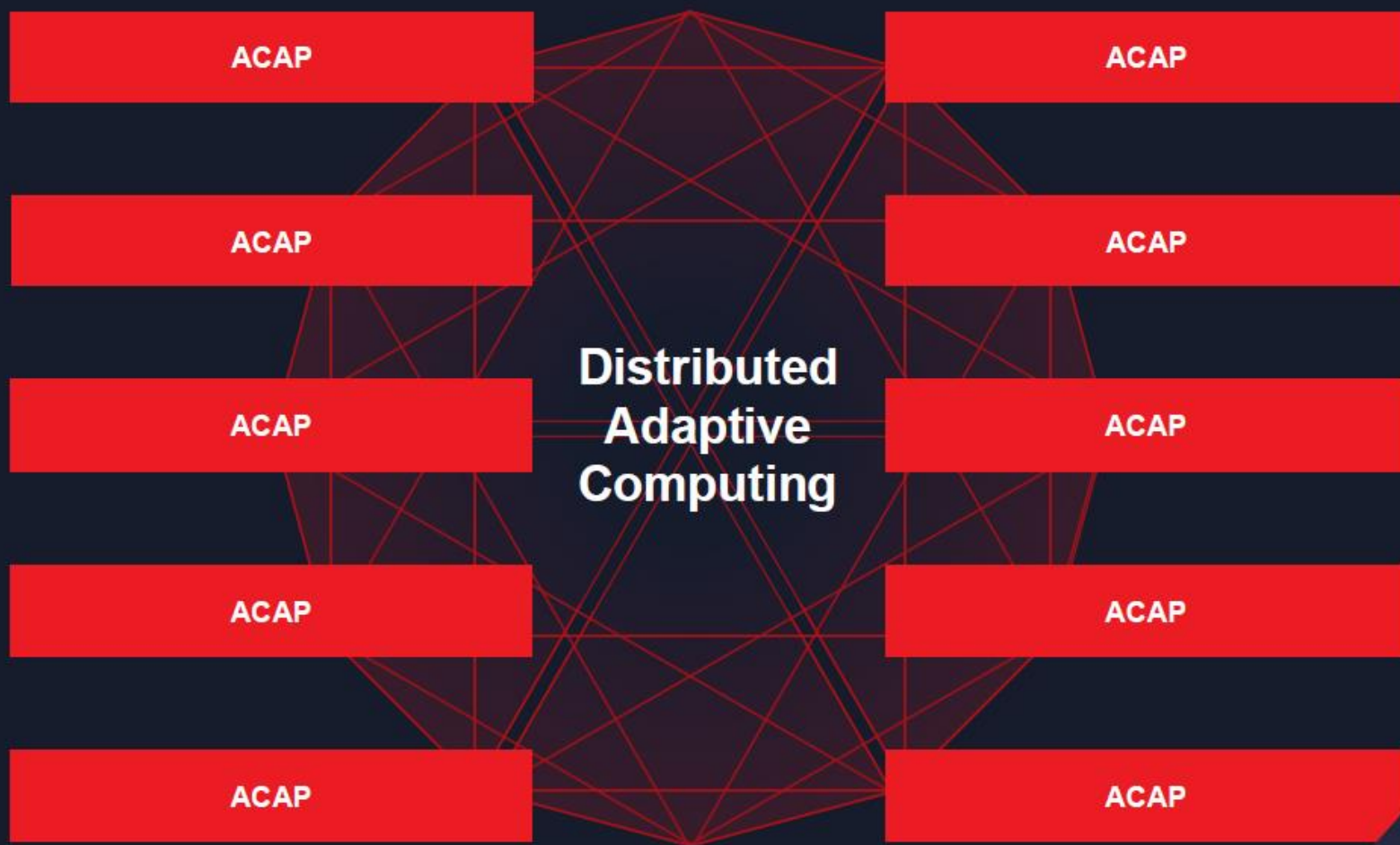
# ➤ Vision: ACAP for Distributed Adaptive Computing

## ACAP Integrated Accelerator

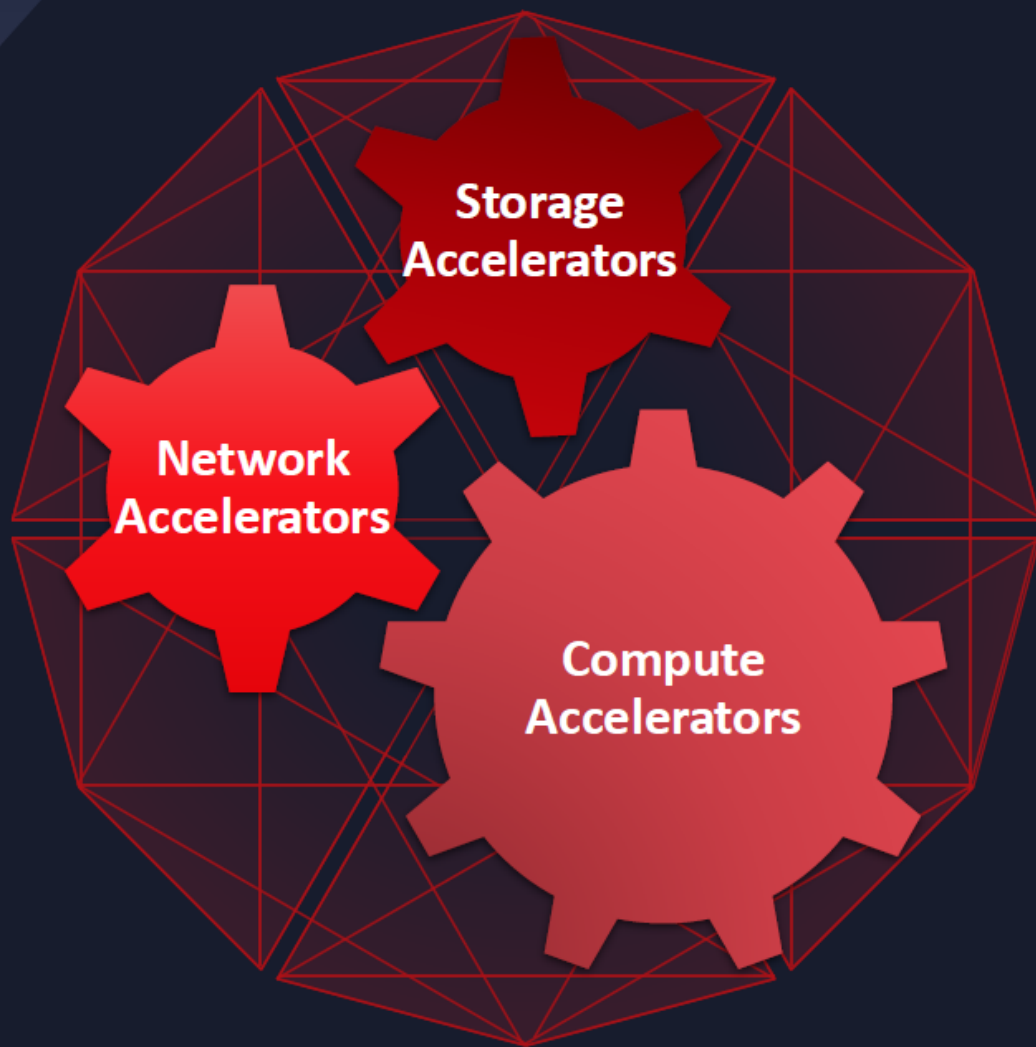




# ➤ Harnessing the Power and Flexibility of Multiple Accelerators



# ➤ Value of Distributed Adaptive Computing



Whole application acceleration across compute, networking and storage

Virtualization and independent scale-out of resources

Best compute optimization and utilization within the lowest power envelope

Significant TCO savings

# Getting Started



# Accelerator Program

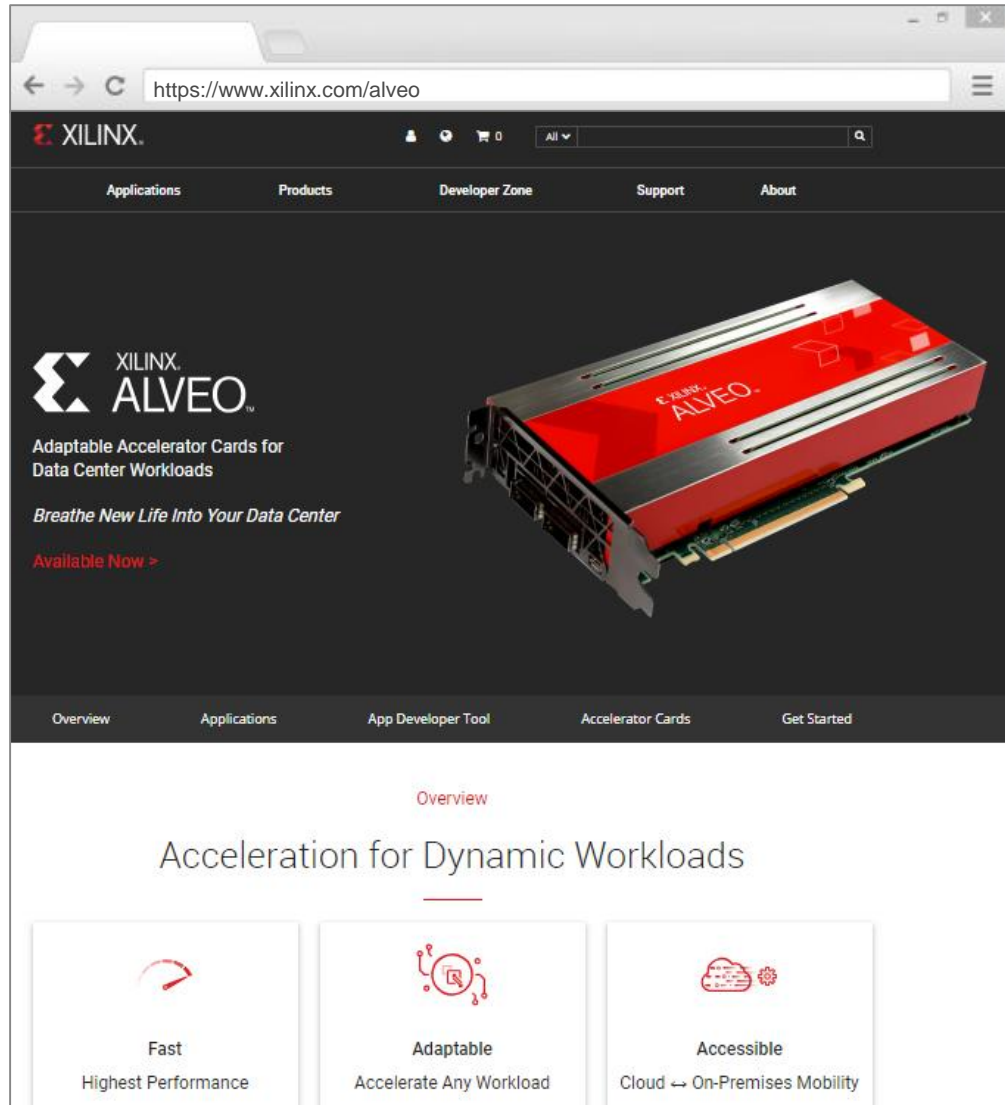
A quick start program to enable companies to accelerate products and services with Alveo™ Data Center accelerator cards and FPGAs in the cloud. Target workloads include data analytics, genomics, video processing, machine learning, financial technology, security, and storage.

Apply Now ➞



[WWW.XILINX.COM/ACCELERATOR-PROGRAM](http://WWW.XILINX.COM/ACCELERATOR-PROGRAM)

# 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](#)

**Thank you!**  
**Contact email : [virajp@xilinx.com](mailto:virajp@xilinx.com)**



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

