# Low Latency AI Inference Acceleration with Mipsology Zebra and Xilinx Alveo

**Andy Luo**

AI Product Maketing

**XILINX**®

# AI Applications Powered by Xilinx

AI Proliferation
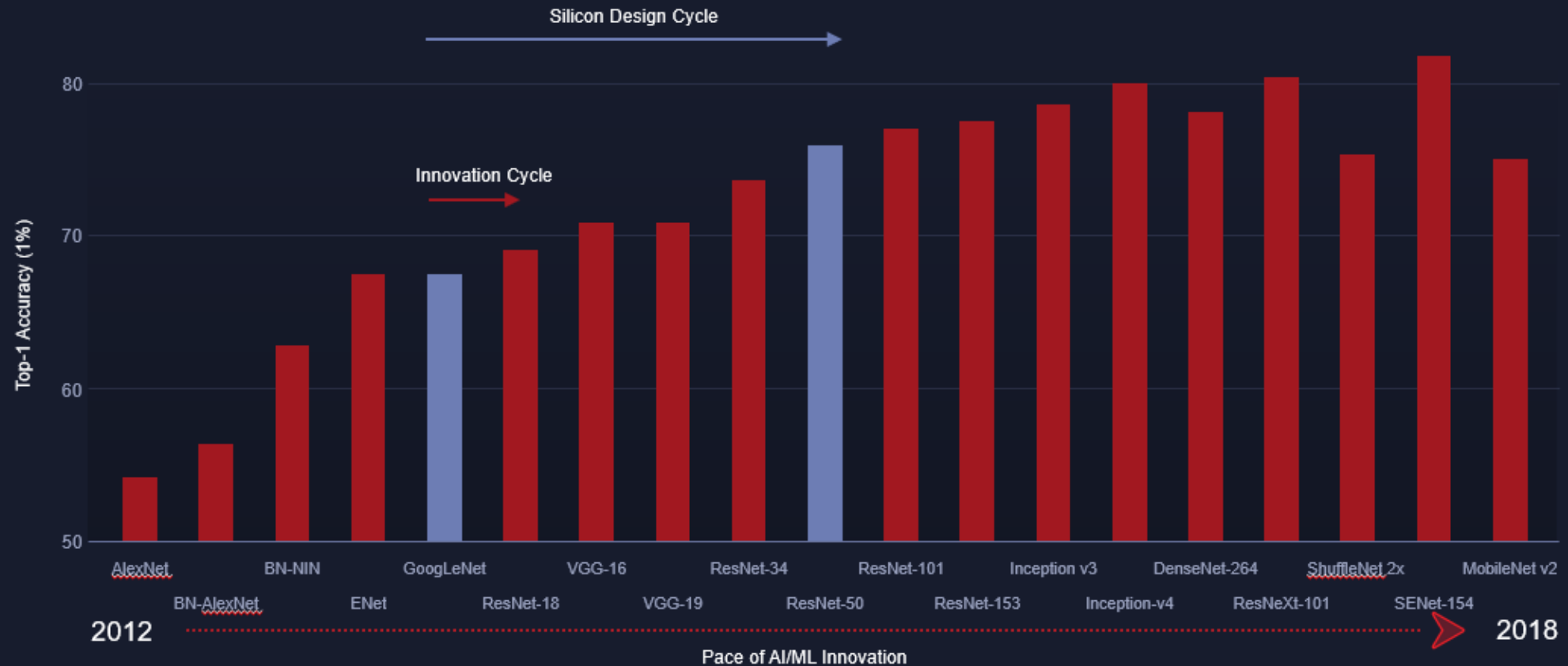
Data Center

5G

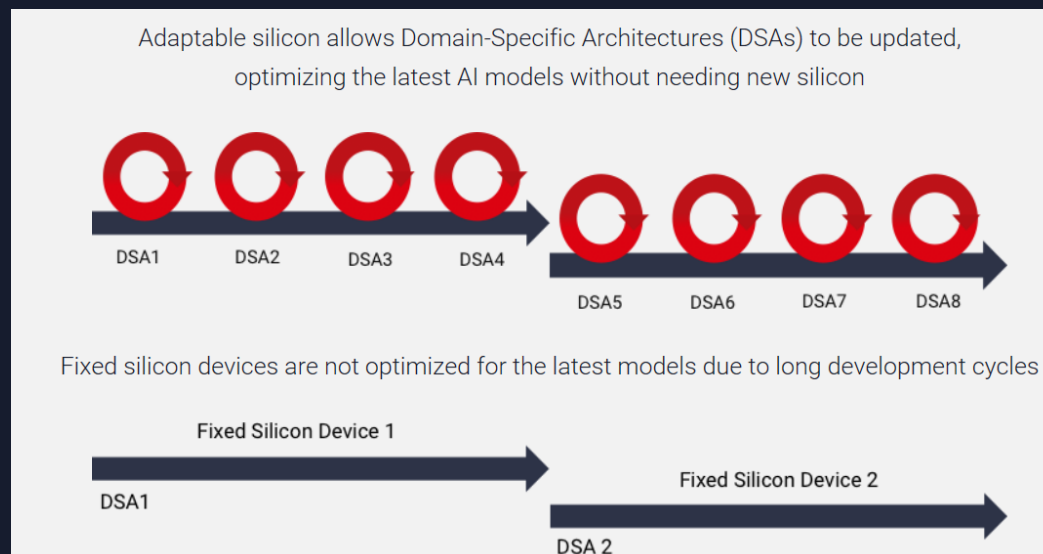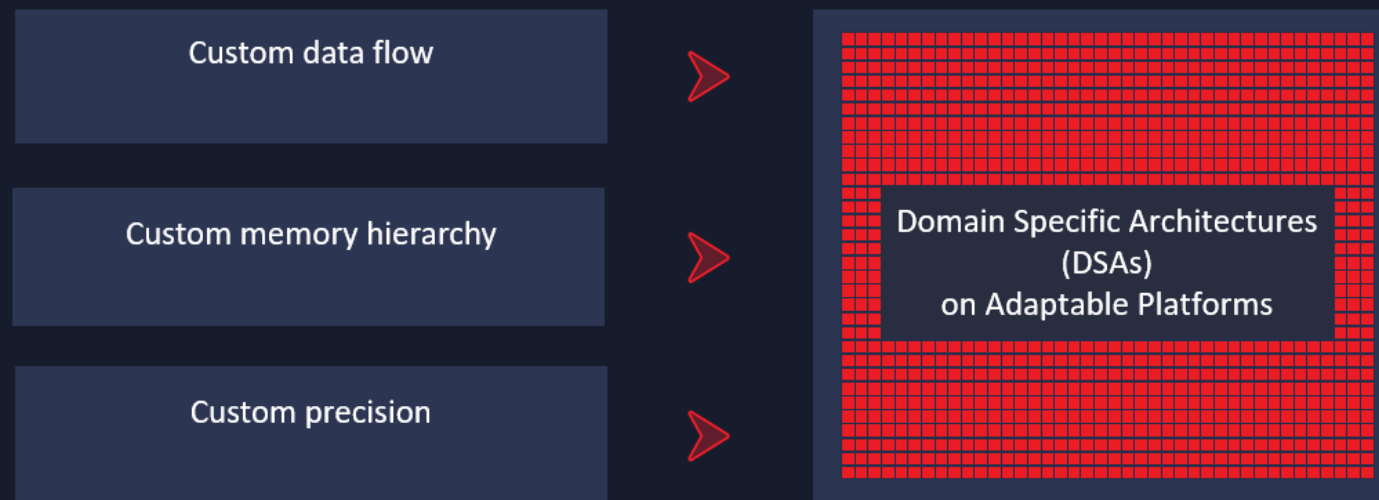Autonomous Driving

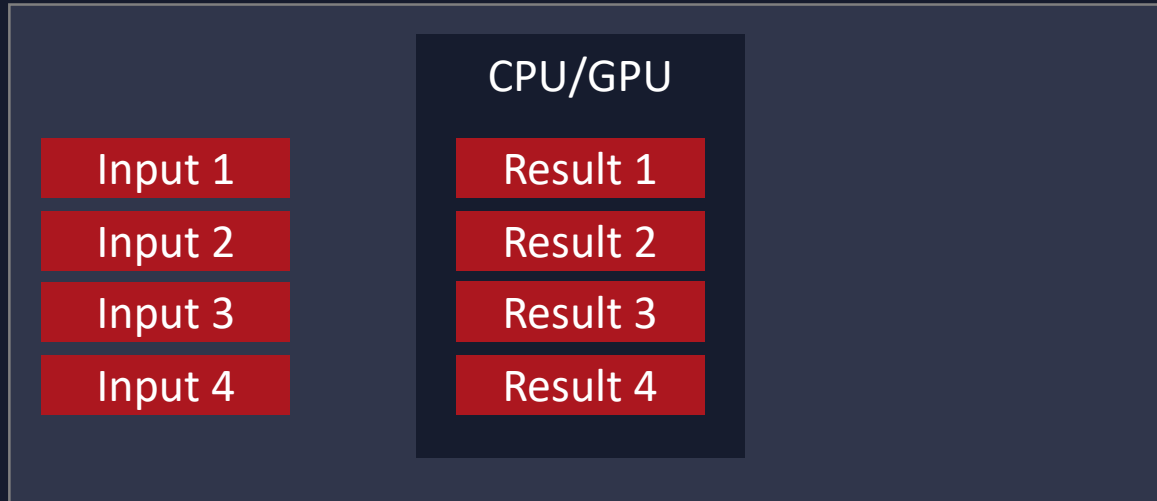Security

Genomics

Video Analytics

Healthcare

Finance

# AI is Evolving Rapidly



Silicon Design Cycle

Innovation Cycle

Top-1 Accuracy (1%)

2012

2018

Pace of AI/ML Innovation

AlexNet
BN-AlexNet
BN-NIN
ENet
GoogLeNet
ResNet-18
VGG-16
VGG-19
ResNet-34
ResNet-50
ResNet-101
ResNet-153
Inception v3
Inception-v4
DenseNet-264
ResNeXt-101
ShuffleNet 2x
SENet-154
MobileNet v2

XILINX

# Adapt to Your AI Workloads

Custom data flow

Custom memory hierarchy

Custom precision

Domain Specific Architectures
(DSAs)
on Adaptable Platforms

Adaptable silicon allows Domain-Specific Architectures (DSAs) to be updated,
optimizing the latest AI models without needing new silicon

DSA1    DSA2    DSA3    DSA4

DSA5    DSA6    DSA7    DSA8

Fixed silicon devices are not optimized for the latest models due to long development cycles

Fixed Silicon Device 1

DSA1

Fixed Silicon Device 2

DSA 2

XILINX.

# Low Latency AI Inference

**CPU/GPU**

Input 1    Result 1

Input 2    Result 2

Input 3    Result 3

Input 4    Result 4

**FPGA/ACAP**

Input 1    Result 1

Input 2    Result 2

Input 3    Result 3

Input 4    Result 4

50ms latency
response

3ms latency
response

# Adapt to Your AI Application

Performance
Critical
Functions

AI
Inference

Performance
Critical
Functions

## Xilinx – Matched Throughput

Performance
Critical Functions

AI Inference

Performance
Critical Functions

## GPU & CPU – Mismatched Throughput

Performance
Critical Functions

AI Inference

Performance
Critical Functions

XILINX

# How to Deploy AI Inference to Xilinx Platform

**AI DSA**

**Direct Framework Compilation**

**Minutes of Compile Times**

MIN          HRS

AI Model

TensorFlow

PyTorch

AI Toolchain

Platform

ALVEO

XILINX.

# Adaptable.
# Intelligent.

**XILINX**

# Mipsology

# Zebra™ by Mipsology

Accelerating computation for machine learning

Presented by: Ludovic (Ludo) Larzul, Founder and CEO, Mipsology

Date: February 25, 2020

*We Focus on Acceleration, You Focus on Your Application!*

# Welcome

Low Latency AI Inference Acceleration with **Zebra** by Mipsology and Xilinx Alveo

Ludovic Larzul, Mipsology CEO

Robert Lara, Mipsology Sr. Director

Mario Trentini, Mipsology Sr. Director



For more information email us at zebra@mipsology.com

# What is Zebra by Mipsology?

Zebra™ software accelerates inference computation faster and easier for machine-learning AI based systems

Zebra works on most Xilinx Alveo Cards

Mipsology

# What Makes Zebra Unique

**Mipsology**

## Ease of Use
With a single command and zero change, any engineer can replace CPU or GPU by Zebra without any knowledge of FPGA.

## Large Support of Neural Network
With a large support of neural network, its start from the same training, results with similar accuracy, Zebra does not require scarce AI resources to do any change.

## Best Performance
Zebra delivers immediately the best performance at low latency from small to large FPGA, matching AI needs from embedded to data centers.

## Lower Cost of Ownership
Quick transition for low NRE and large choice of FPGA and cards, with long life span to lower the TCO.

# How does Zebra work?

✓ Keep your GPU for training

✓ Keep your same neural network for GPU/CPU

✓ Keep your same application software for GPU/CPU

✓ Keep your same framework for GPU/CPU

✓ **Type a Single Linux Command**

✓ Zebra will do automatic proprietary quantization

✓ Experience best performance with low latency
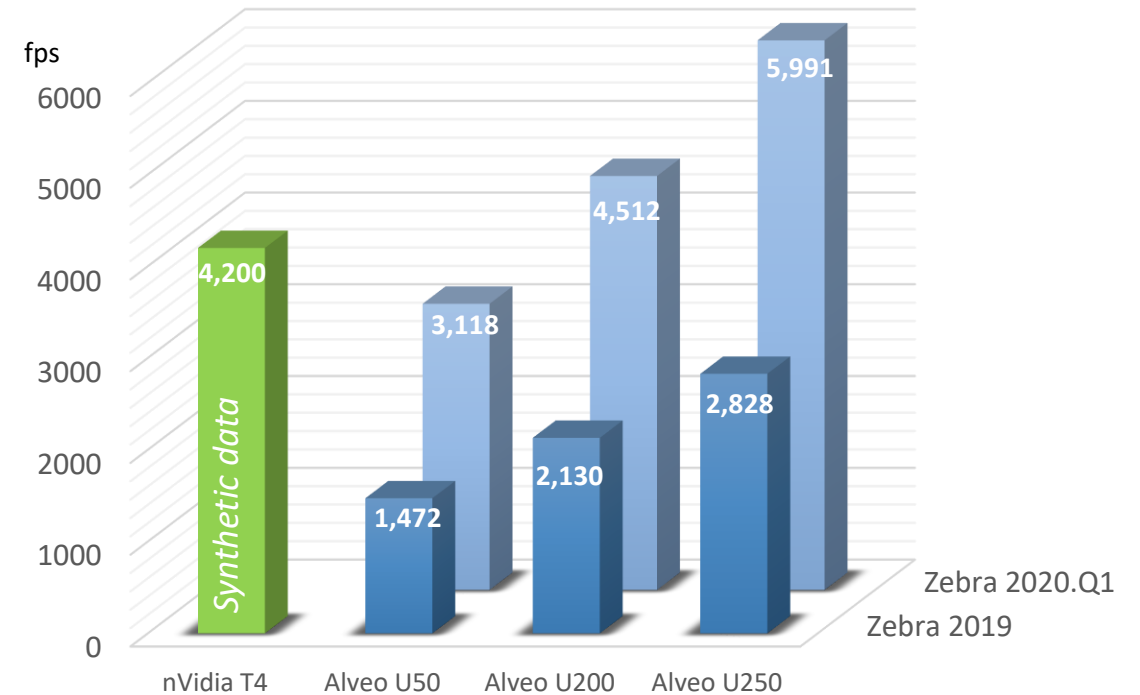
✓ Of course, Zebra works on Xilinx Alveo Boards

Mipsology

# Zebra Metrics – Best Performance

**ResNet50 Performance on actual data, single-image batch**

fps

- nVidia T4 — Synthetic data: **1,077**
- Alveo U50 — Zebra 2019: **1,104**, Zebra 2020.Q1: **1,656**
- Alveo U200 — Zebra 2019: **1,427**, Zebra 2020.Q1: **2,069**
- Alveo U250 — Zebra 2019: **2,087**, Zebra 2020.Q1: **2,922**

**ResNet50 Performance on actual data, multi-image batch**

fps

- nVidia T4 — Synthetic data: **4,200**
- Alveo U50 — Zebra 2019: **1,472**, Zebra 2020.Q1: **3,118**
- Alveo U200 — Zebra 2019: **2,130**, Zebra 2020.Q1: **4,512**
- Alveo U250 — Zebra 2019: **2,828**, Zebra 2020.Q1: **5,991**

Zebra 2019: batch below 64. Production version. All performance measured.
Zebra 2020.Q1: batch below 64. Alpha version, production end of Q3. Performance measured on Alveo U250, scaled for U200 and U50.
nVidia T4: performance for same latency as Alveo. ResNet50 performance number from nVidia website measured with synthetic data. nVidia does not publish Yolo performance.

Mipsology

# Zebra Metrics –
## Keeps Accuracy with Proprietary quantization



Top1 classification:

Mipsology

# Zebra enables real-time processing of cameras at full speed

At 25 frame per second, a reactive application has 40ms to process an image from a camera:



Long NN processing means late actions:



➤ A slow robot moving at 12ft/s → 1 inch per 10ms

➤ A car moving at 65mph → 8 inches per 10ms

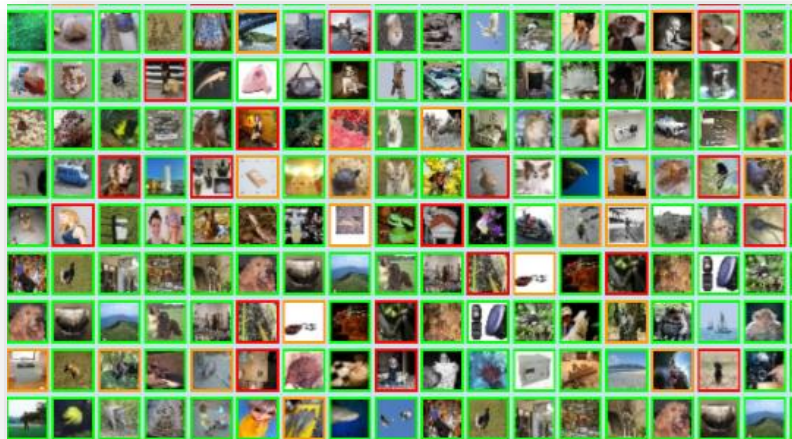Zebra on Alveo U50 reacts **3x faster than best performing GPU**.



**Latency: LOWER IS BETTER**

(computing time for an image, based on customer Yolo-class CNN)

# Zebra for All Applications

*Tested networks: AlexNet, CaffeNet, GoogLeNet, inceptionV3, inceptionV4, ResNet50, ResNet152, NiN, VGG16, VGG19, YoloV1, YoloV2, YoloV3, VDSR, SR_ResNet, MobileNet, SSD, and many custom ones*

**Classification**



**Segmentation**



**Super Resolution**

VDSR    Prior Art

VS



**Body Positioning**



**Mipsology**

# Zebra Demo on Alveo U50

Mipsology © 2015 - 2020

Mipsology

# Mipsology

# REQUEST OUR ZEBRA SOLUTION OVERVIEW

Zebra Solution Overview
https://mipsology.com/product/#download

Visit us at www.mipsology.com    Or email us at zebra@mipsology.com