# Inference Architectures @Xilinx

Graham Schelle, PhD
Principal Engineer
Xilinx Research Labs

**XILINX.**

# Xilinx Headlines



When Databases Meet FPGA – Achieving 1 Million TPS with X-DB Heterogeneous Computing
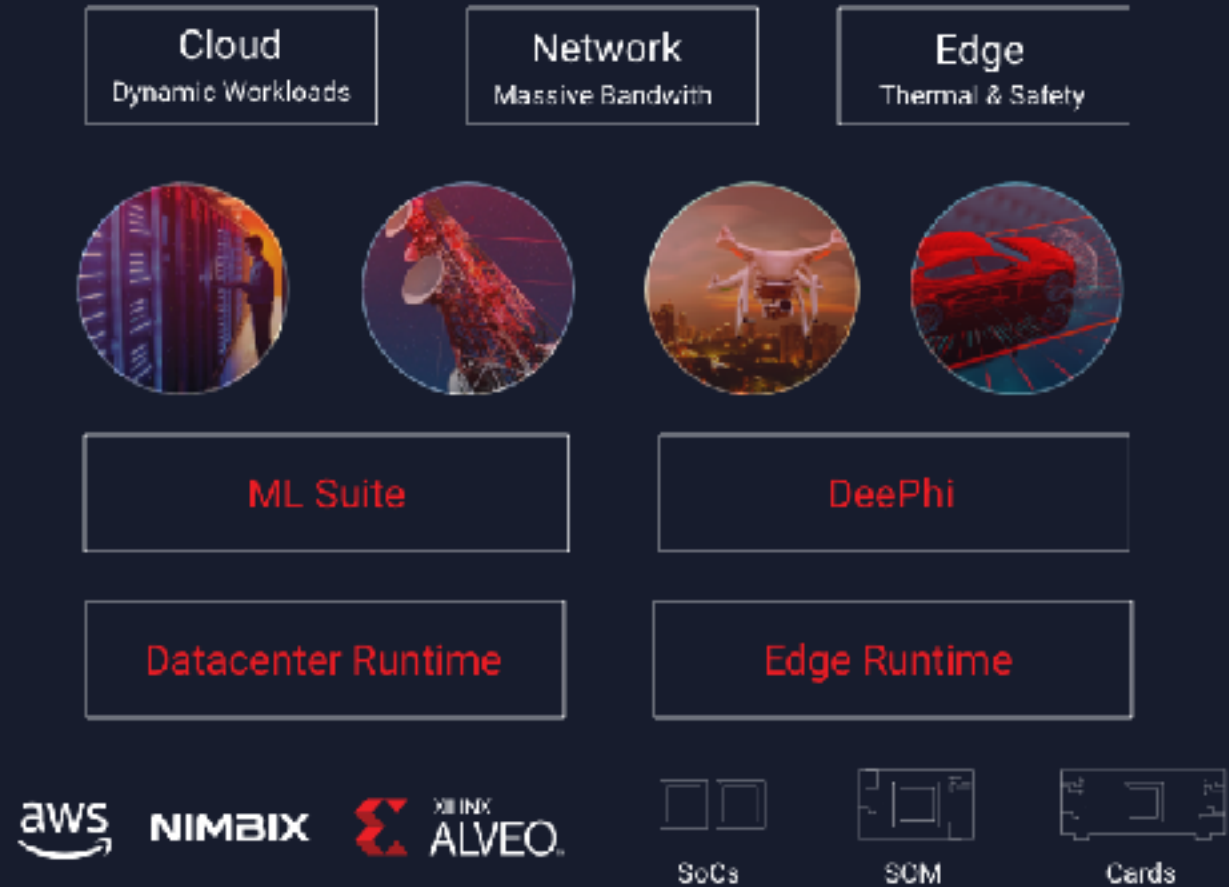alibabacloud.com



Children's Hospital of Philadelphia And Edico Genome Achieve Fastest-Ever Analysis Of 1,000 Genomes
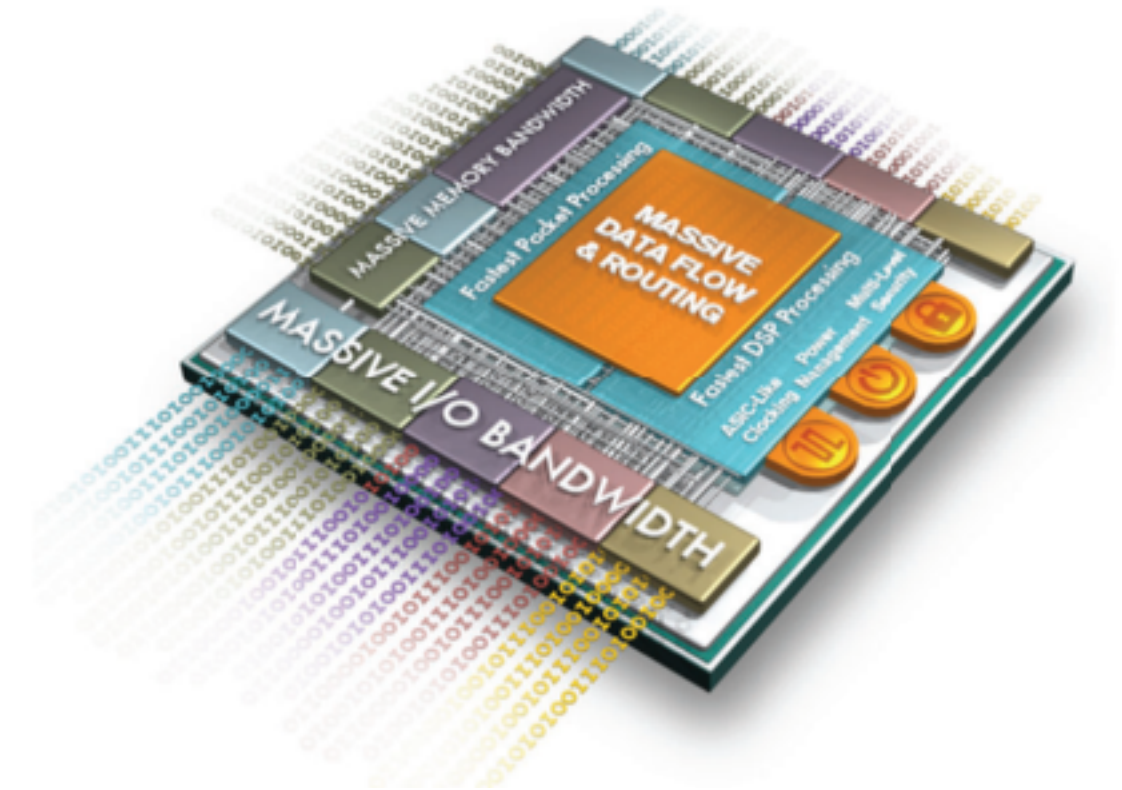
Twitch Chooses Xilinx to Enable its Broadcast-quality Livestream of eSports

XILINX

# Agenda

> Xilinx Adaptive Architectures

> Inference Architectures
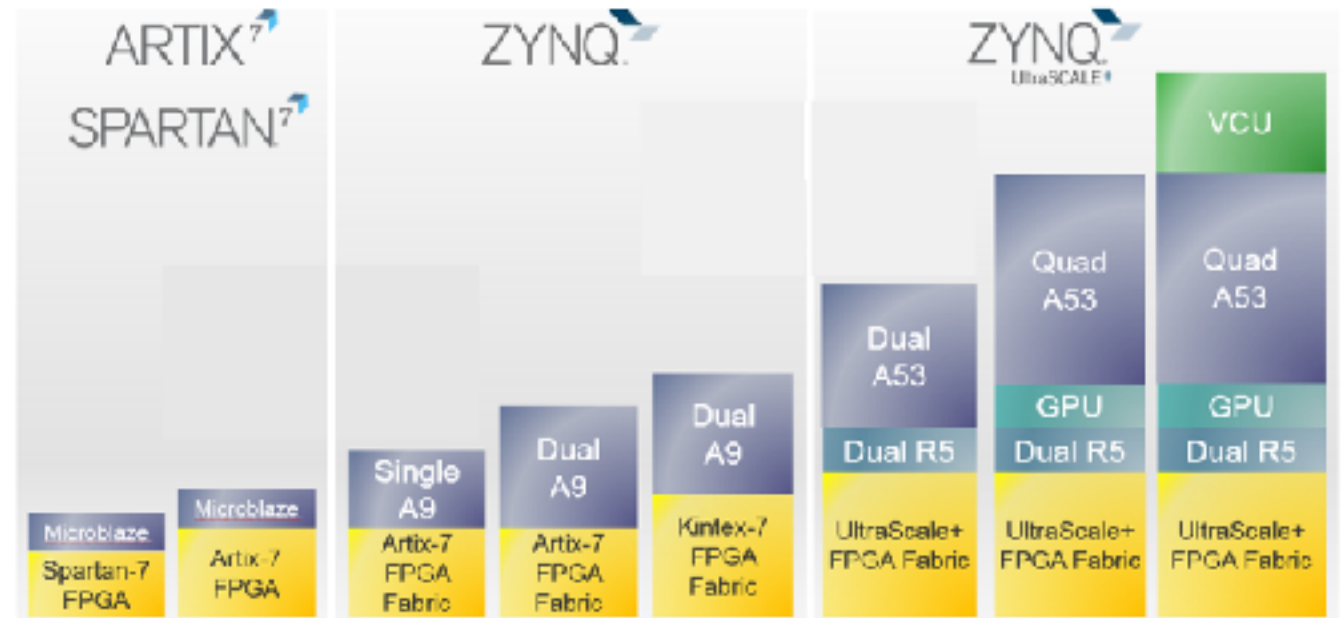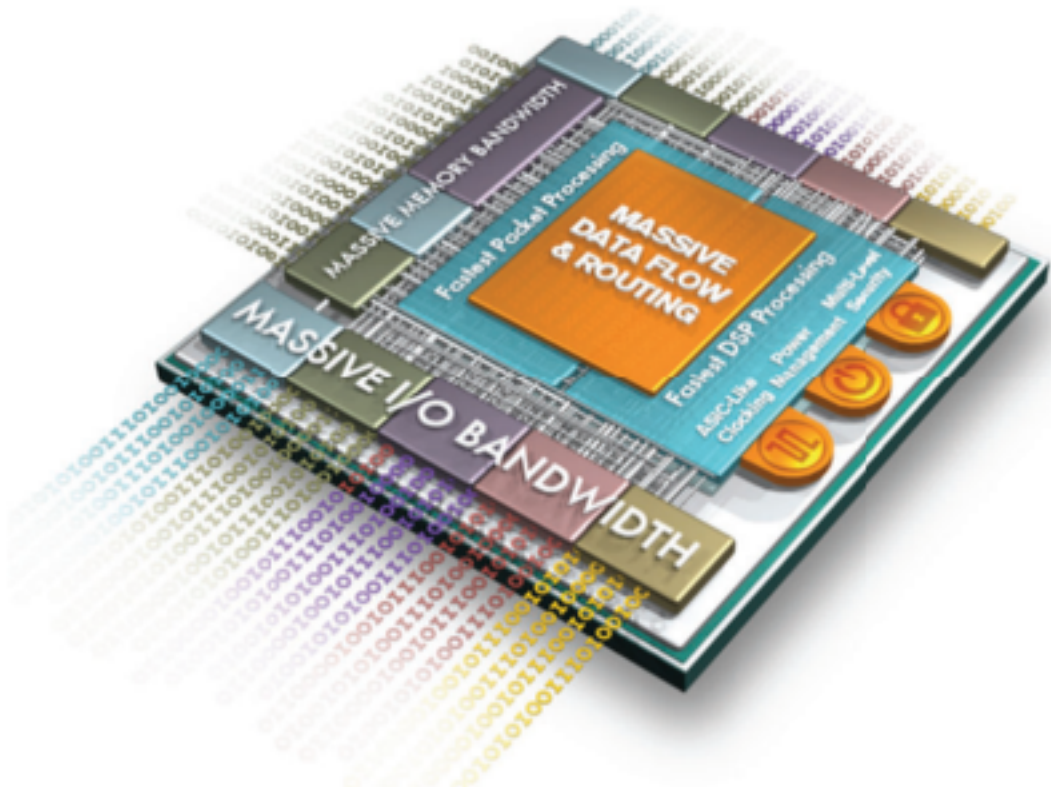
> Open Source

# Xilinx Adaptive Architectures



Traditionally, FPGAs for massively data-parallel applications

XILINX

# Xilinx Adaptive Architectures



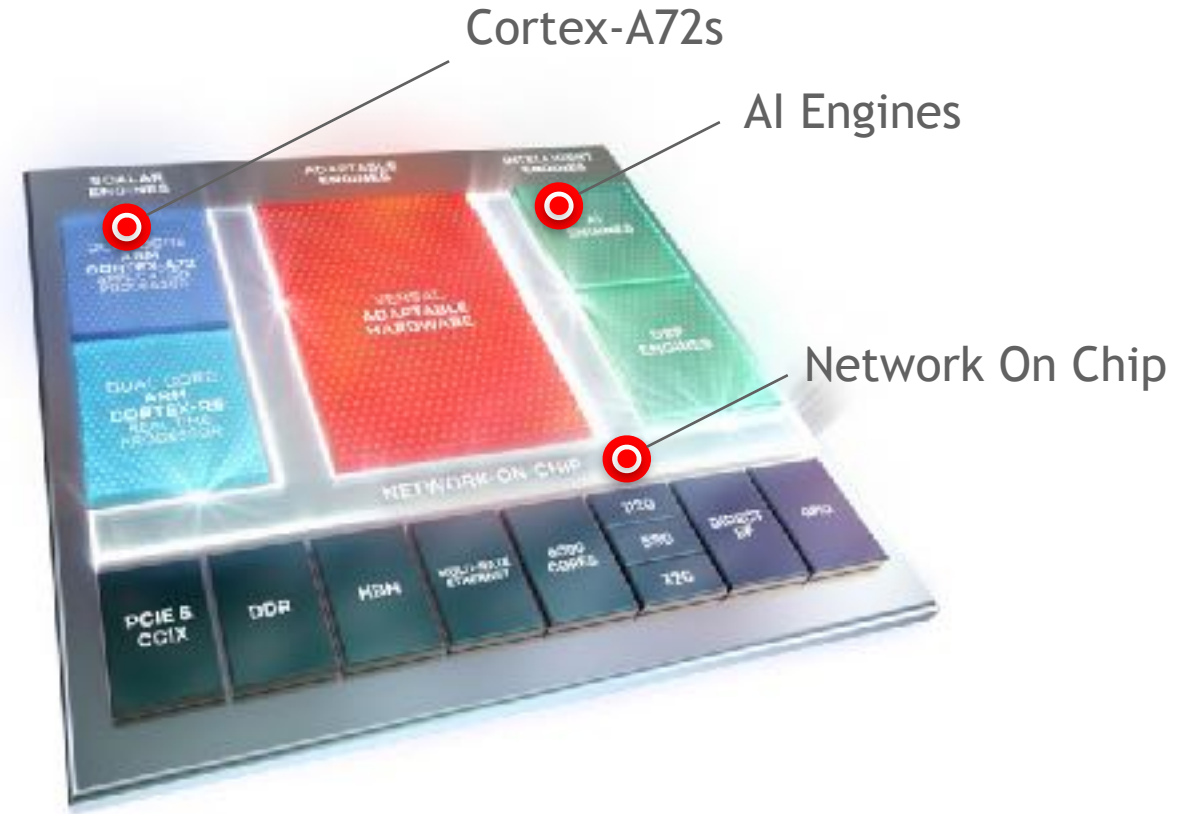Traditionally, FPGAs for massively data-parallel applications

In 2011, Zynq introduced (ZU+ in 2015) ARM CPUs added for embedded applications

XILINX

# Xilinx Adaptive Architectures – Alveo & Versal
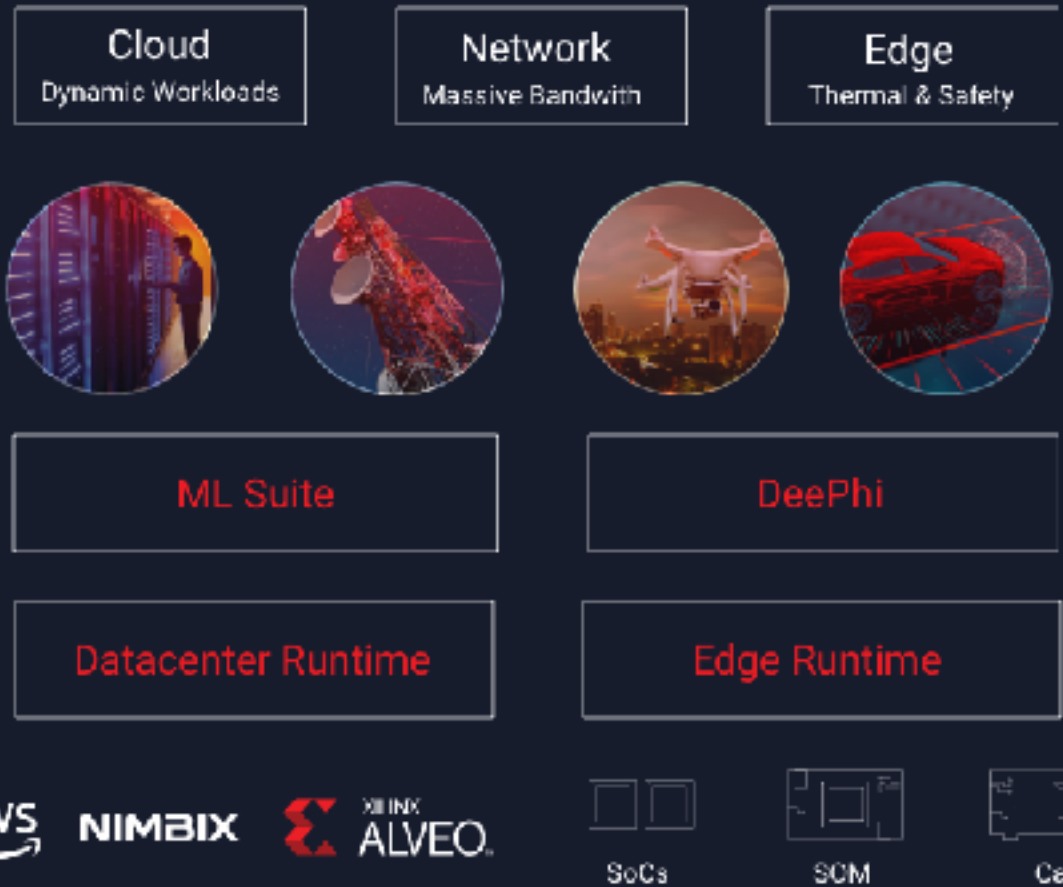


Cortex-A72s

AI Engines

Network On Chip

In 2018, Alveo introduced
Accelerator cards for data center workloads

Coming in 2019, Versal Platform
Adaptive compute acceleration platform (ACAP)

# Inference Architectures

| Cloud | Network | Edge |
|-------|---------|------|
| Dynamic Workloads | Massive Bandwith | Thermal & Safety |

| ML Suite | DeePhi |
|----------|--------|

| Datacenter Runtime | Edge Runtime |
|--------------------|--------------|

aws    NIMBIX    XILINX ALVEO.

SoCs    SCM    Cards

XILINX.

# Inference Architectures – Evolving Frameworks



**Andrej Karpathy on Twitter**

> **Increasing, Evolving Workloads**
>> New acceleration needs & algorithms
>> ML "infused" in many applications
>> Adaptable HW a key benefit

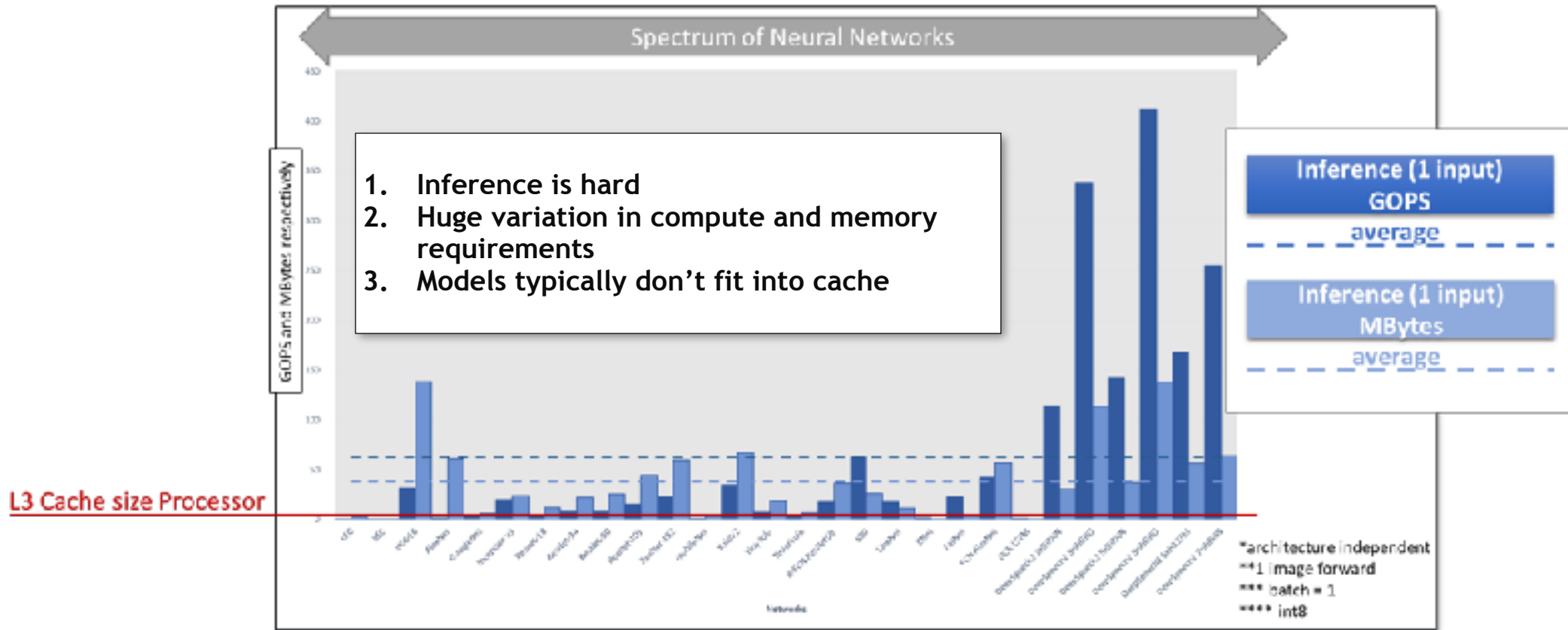**XILINX**

# Inference Architectures – Evolving Frameworks



**Andrej Karpathy on Twitter**

> **Increasing, Evolving Workloads**
> >> New acceleration needs & algorithms
> >> ML "infused" in many applications
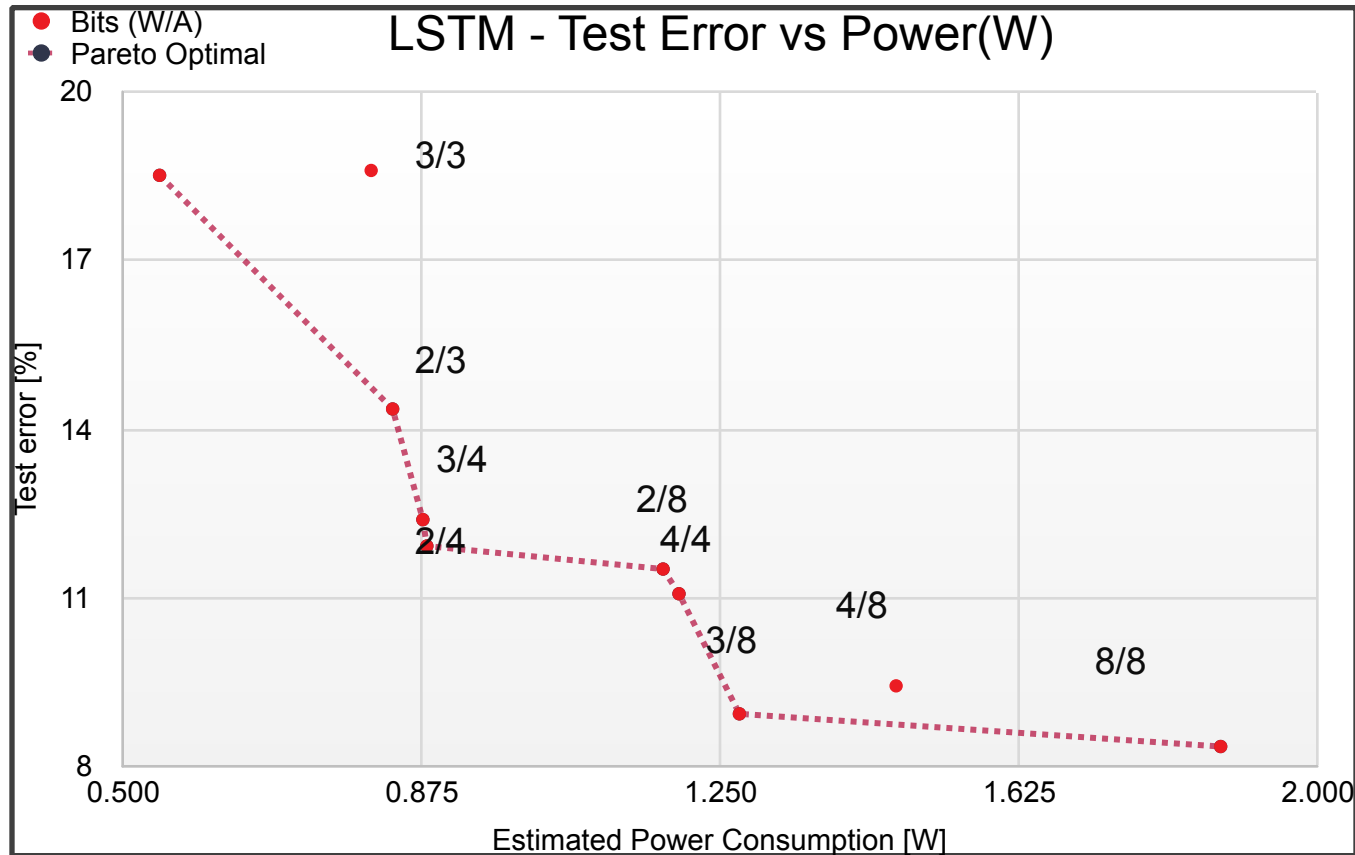> >> Adaptable HW a key benefit

> **Move to Lower Precision**
> >> ML inference moving to INT8 & lower
> >> Better Perf/W with similar accuracy
> >> Xilinx devices natively support variable precision

> **Compressed Networks**
> >> Higher performance with reduced compute / memory needs
> >> Pruning & load balancing to match network requirements

XILINX

# Inference Architectures – Evolving Workloads



## Increasing, Evolving Workloads

- New acceleration needs & algorithms
- ML "infused" in many applications
- Adaptable HW a key benefit

## Move to Lower Precision

- ML inference moving to INT8 & lower
- Better Perf/W with similar accuracy
- Xilinx devices natively support variable precision

## Compressed Networks

- Higher performance with reduced compute / memory needs
- Pruning & load balancing to match network requirements

# Inference Architectures – Evolving Workloads



**Increasing, Evolving Workloads**

>> New acceleration needs & algorithms

>> ML "infused" in many applications

>> Adaptable HW a key benefit

**Move to Lower Precision**

>> ML inference moving to INT8 & lower

>> Better Perf/W with similar accuracy

>> Xilinx devices natively support variable precision

**Compressed Networks**

>> Higher performance with reduced compute / memory needs

>> Pruning & load balancing to match network requirements

XILINX

# Inference Architectures – Precision vs Power

**FPGA:**



**ASIC:**



*Source: Bill Dally (Stanford), Cadence Embedded Neural Network Summit, February 1, 2017*

Michaela Blott, Hot Chips 2018 Tutorial, "Overview of Deep Learning and Computer Architectures for Accelerating DNNs"

*Target Device ZU7EV ● Ambient temperature: 25 °C ● 12.5% of toggle rate ● 0.5 of Static Probability ● Power reported for PL accelerated block only*

*Rybalkin, V., Pappalardo, A., Ghaffar, M.M., Gambardella, G., Wehn, N. and Blott, M. "FINN-L: Library Extensions and Design Trade-off Analysis for Variable Precision LSTM Networks on FPGAs."*

XILINX.

# Xilinx Cloud Inference - ML Suite Overlays with xDNN

**Adaptable**
> AI algorithms are changing rapidly
> Adjacent acceleration opportunities

**Realtime**
> Lower latency than CPU and GPU
> Data flow processing

**Efficient**
> Performance/watt
> Low Power



- Built in Programmable Logic
- High Utilization, Thput or Latency Variants
- CPU offload for new layer exploration

xDNN w/ xfDNN Compiler



On-prem and cloud boards



https://github.com/Xilinx/ml-suite

# Xilinx Edge Inference - DeePhi



(2013)

(2016)

"Learning both Weights and Connections for Efficient Neural Networks", NeurIPS 2015

"EIE: Efficient Inference Engine on Compressed Deep Neural Network", ISCA 2016

"ESE: Efficient Speech Recognition Engine with Compressed LSTM on FPGA", FPGA 2017

XILINX.

# Xilinx Edge Inference - DeePhi



(2013)

(2016)

"Learning both Weights and Connections for Efficient Neural Networks", NeurIPS 2015

"EIE: Efficient Inference Engine on Compressed Deep Neural Network", ISCA 2016

"ESE: Efficient Speech Recognition Engine with Compressed LSTM on FPGA", FPGA 2017

Xilinx Announces the Acquisition of DeePhi Tech

Deal to Accelerate Data Center and Intelligent Edge Applications

Jul 17, 2018

PYNQ-Z1

ZCU104

UltraS6

ZCU102

# Cloud & Edge Integration

# Xilinx and Open Source

Cloud
Dynamic Workloads

Network
Massive Bandwith

Edge
Thermal & Safety

ML Suite

DeePhi

Datacenter Runtime

Edge Runtime

aws    NIMBIX    XILINX ALVEO.

SoCs    SCM    Cards

XILINX.

# Xilinx and Open Source

PYNQ



Quantized Neural Networks



Xilinx Runtime for PCIe Attached FPGAs



…More on www.github.com/Xilinx

XILINX

# Xilinx and Open Source

PYNQ

Quantized Neural Networks

Xilinx Runtime for PCIe Attached FPGAs



…More on www.github.com/Xilinx

XILINX

# Python is increasingly the Language of Choice

Top Programming Languages,
IEEE Spectrum, July'18



| Language Rank | Types | | | | Spectrum Ranking |
|---|---|---|---|---|---|
| 1. Python | 🌐 | | 🖥 | ▊ | 100.0 |
| 2. C++ | 📱 | | 🖥 | ▊ | 98.4 |
| 3. C | 📱 | | 🖥 | ▊ | 98.2 |
| 4. Java | 🌐 | 📱 | 🖥 | | 97.5 |
| 5. C# | 🌐 | 📱 | 🖥 | | 89.6 |
| 6. PHP | 🌐 | | | | 85.4 |
| 7. R | | | 🖥 | | 83.3 |
| 8. JavaScript | 🌐 | 📱 | | | 82.6 |
| 9. Go | 🌐 | | 🖥 | | 76.7 |
| 10. Assembly | | | | ▊ | 74.5 |

https://spectrum.ieee.org/at-work/innovation/the-2018-top-programming-languages



**To date**

https://stackoverflow.blog/2017/09/06/incredible-growth-python/

15

XILINX.

# Python is increasingly the Language of Choice

Top Programming Languages,
IEEE Spectrum, July'18

| Language Rank | Types | Spectrum Ranking |
|---|---|---|
| 1. Python | 🌐 🖥️ ▪️ | 100.0 |
| 2. C++ | 📱 🖥️ ▪️ | 98.4 |
| 3. C | 📱 🖥️ ▪️ | 98.2 |
| 4. Java | 🌐 📱 🖥️ | 97.5 |
| 5. C# | 🌐 📱 🖥️ | 89.6 |
| 6. PHP | 🌐 | 85.4 |
| 7. R | 🖥️ | 83.3 |
| 8. JavaScript | 🌐 📱 | 82.6 |
| 9. Go | 🌐 🖥️ | 76.7 |
| 10. Assembly | ▪️ | 74.5 |

Python is listed as an embedded language for the first time

https://spectrum.ieee.org/at-work/innovation/the-2018-top-programming-languages



**Growth of major programming languages**
Based on Stack Overflow question views in World Bank high-income countries

**To date**

% of overall question views each month

https://stackoverflow.blog/2017/09/06/incredible-growth-python/

XILINX.

# Python is increasingly the Language of Choice

Top Programming Languages,
IEEE Spectrum, July'18



https://spectrum.ieee.org/at-work/innovation/the-2018-top-programming-languages



To date

https://stackoverflow.blog/2017/09/06/incredible-growth-python/

**Python is the fastest growing language: driven by data science, AI, ML and academia**

# Python Productivity for Zynq

# Python Productivity for Zynq

Jupyter notebooks, browser-based interface

| Jupyter web server | |
|---|---|
| IPython kernel | |
| Ubuntu-based Linux | Overlays/designs |
| ARM A9 / A53 | ZU+ Fabric |

# Python Productivity for Zynq

Jupyter notebooks,
browser-based interface

PYNQ enables JupyterLab
on Zynq and ZU+

| Jupyter web server | |
| --- | --- |
| IPython kernel | |
| Ubuntu-based Linux | Overlays/designs |
| ARM A9 / A53 | ZU+ Fabric |

# Python Productivity for Zynq

Jupyter notebooks, browser-based interface

PYNQ enables JupyterLab on Zynq and ZU+

FPGA designs delivered as Python packages

| Jupyter web server |  |
| --- | --- |
| IPython kernel |  |
| Ubuntu-based Linux | Overlays/designs |
| ARM A9 / A53 | ZU+ Fabric |

XILINX

# Python Productivity for Zynq

Jupyter notebooks, browser-based interface

PYNQ enables JupyterLab on Zynq and ZU+

FPGA designs delivered as Python packages

| | |
|---|---|
| Jupyter web server | |
| IPython kernel | |
| Ubuntu-based Linux | Overlays/designs |
| ARM A9 / A53 | ZU+ Fabric |

Delivered as SD Card image

# PYNQ Community – ML, Non-ML & Academic Partners

**XILINX**

# PYNQ Community – ML, Non-ML & Academic Partners

© Copyright 2018 Xilinx

# Xilinx open source engagements related to today's TVM meeting



MicroPython

# Xilinx open source engagements related to today's TVM meeting



University of Washington

Xilinx Research

MicroPython

UC San Diego

UC Berkeley

# Finally, Xilinx & building new open source communities...

**Cloud Free Trials**



Get started with Alveo accelerator card applications today on the Nimbix Cloud

Nimbix has partnered with Xilinx to provide developers and engineers a trial account that provides up to 100 hours of free time on the Nimbix Cloud using Xilinx Tools and Accelerators.

**pynq.io/community**



**DAC2019 Design Contest**



**OpenHW Design Contest**

# Summary

Xilinx
Great for exploring and deploying inference


Xilinx Open Source
We're actively engaging with TVM and other communities


Email: graham.schelle@xilinx.com
Visit: Boulder, Colorado

# Adaptable.
# Intelligent.

# Edge to Cloud Inference – Automotive

ADAS/AD Central Module

# Edge to Cloud Inference – Automotive



Surround-View Camera Back

Short-Range Radar

Forward-Looking Camera

Drive Monitor Camera

Surround-View Camera Left

Short-Range Radar

Long-Range Lidar

Surround-View Camera Right

ADAS/AD Centra Module

Surround-View Camera Front

Short-Range Radar

# Edge to Cloud Inference – Xilinx Platforms



ZCU104

PYNQ-Z1

Ultra96

ZCU102

F1

amazon
web services

## Edge Devices
Custom I/O, ARM CPUs

## Cloud Platforms
Power Efficient, PCIe, Networking

XILINX.

# Edge to Cloud Inference – IIoT Latency/Data Example
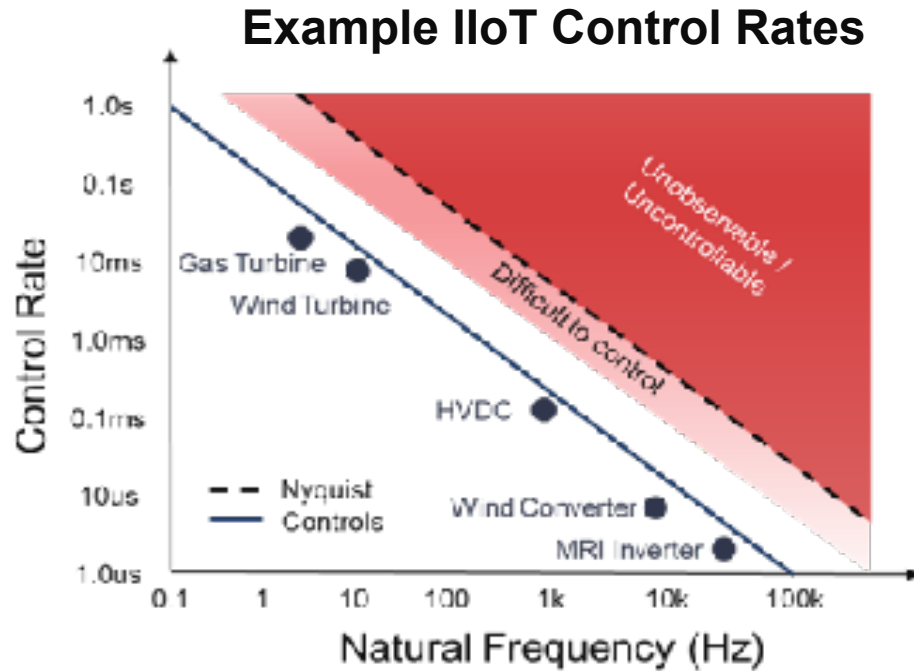


Example IIoT Control Rates

XILINX.

# Edge to Cloud Inference – IIoT Latency/Data Example

**Example IIoT Control Rates**



Distance NYC to LA: 2,800 miles
Speed of light: 186,000 miles/s
**Round trip: 2*2800/186000 = 30ms**
<span style="color:red">**Required Control Rate = 10ms**</span>

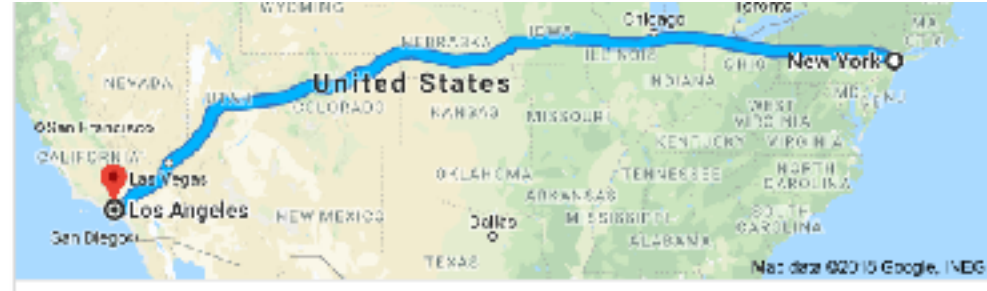**XILINX**

# Edge to Cloud Inference – IIoT Latency/Data Example

## Example IIoT Control Rates
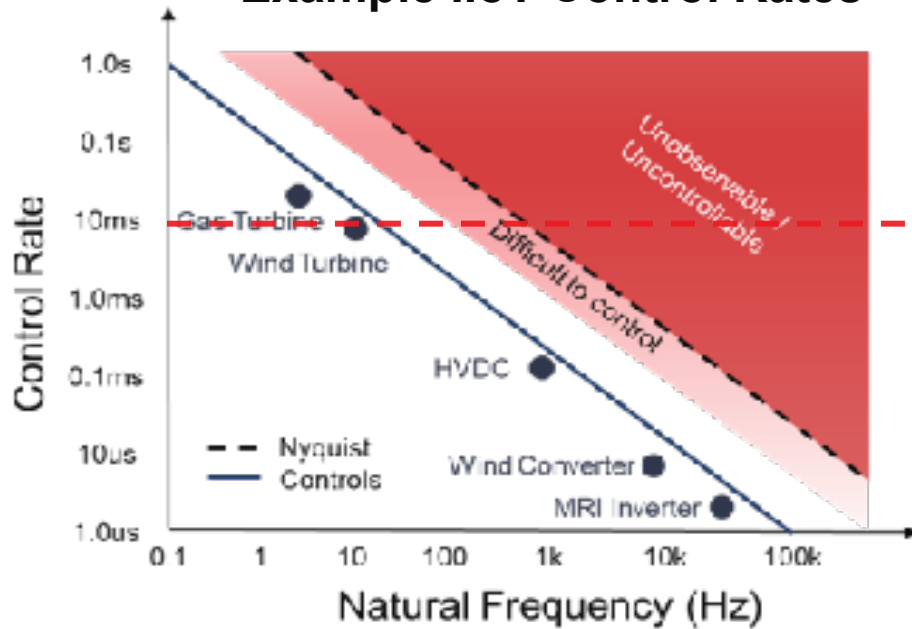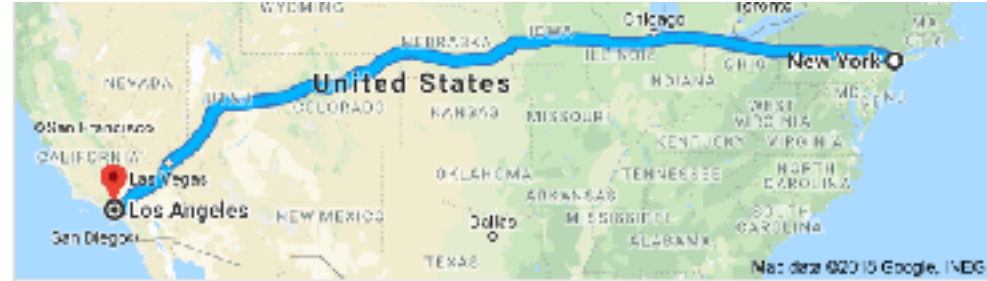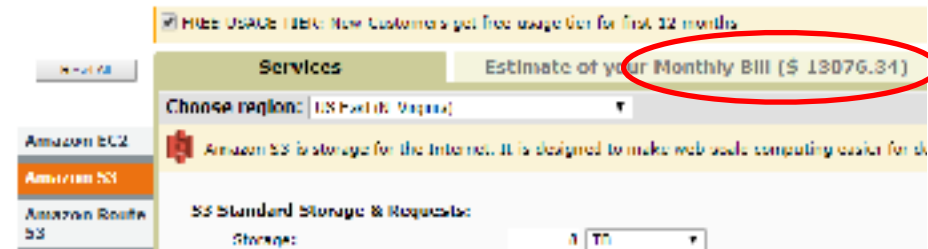


Distance NYC to LA: 2,800 miles
Speed of light: 186,000 miles/s
**Round trip: 2*2800/186000 = 30ms**
**Required Control Rate = 10ms**

E.g. Power Plant @ 8TB/Month

© Copyright 2018 Xilinx