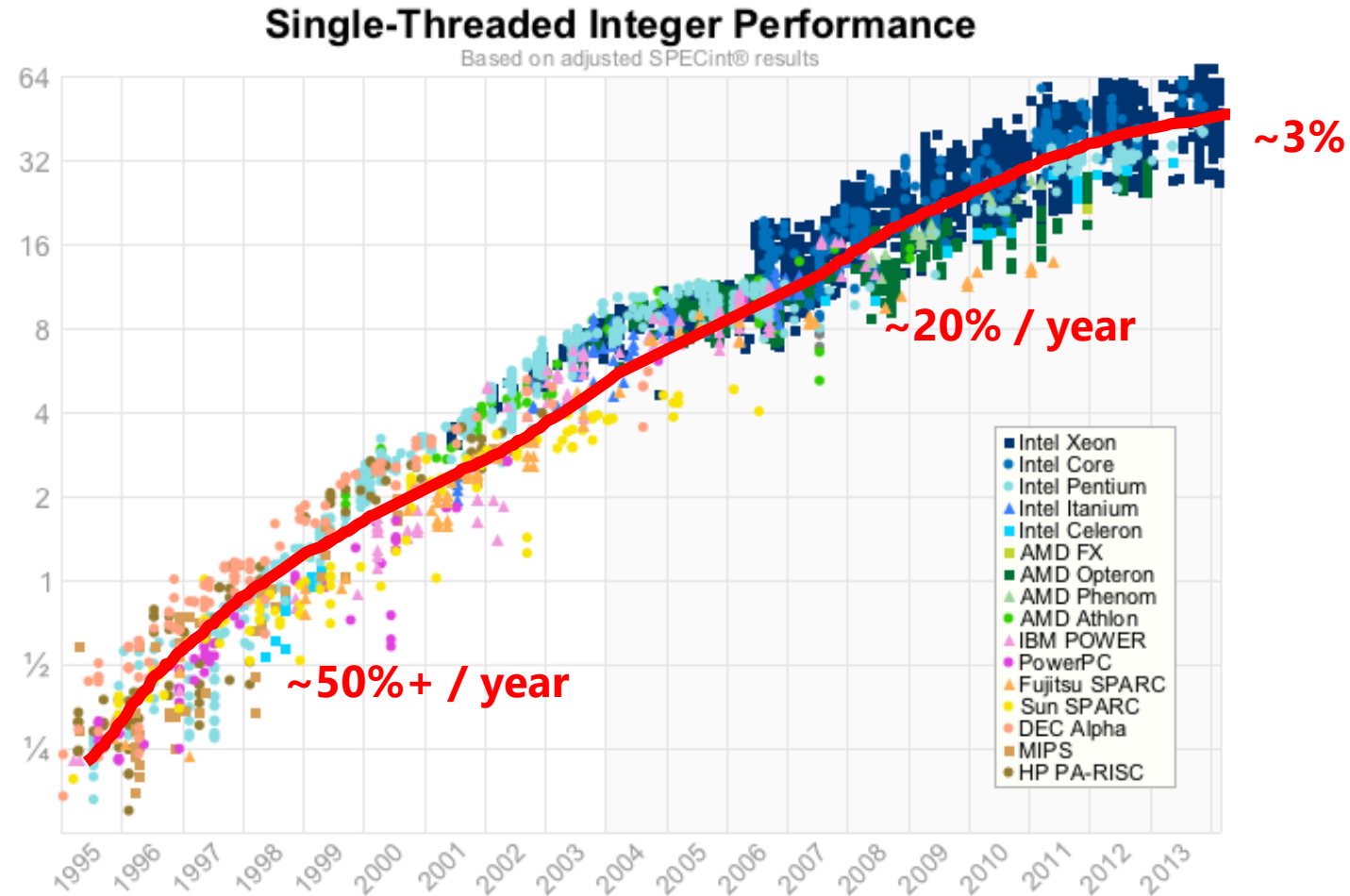Microsoft

# Heterogeneous Computing @ Microsoft

**Andrew Putnam – Azure Networking / Microsoft Research**
**Kalin Ovtcharov – AI & Architectures**
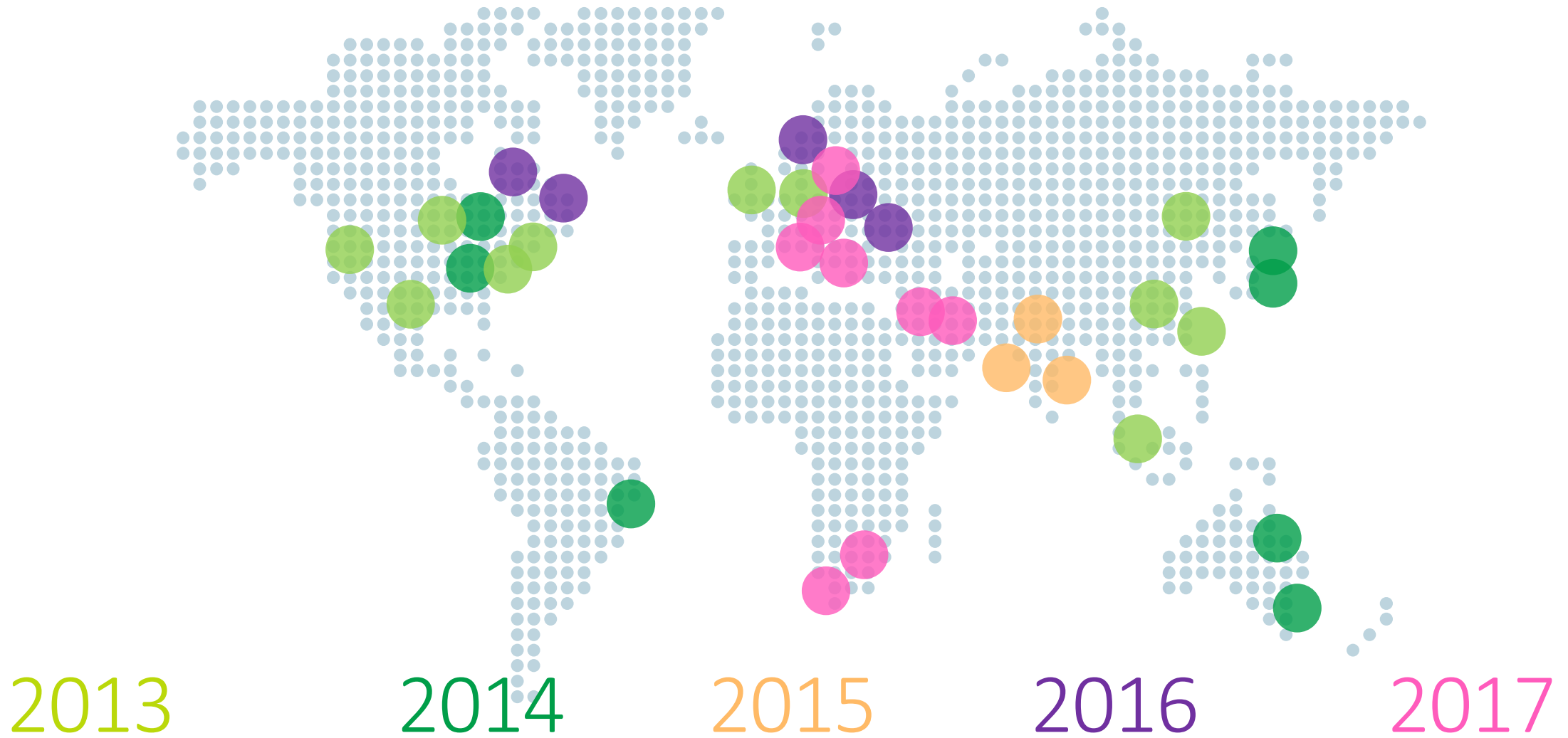September 11, 2019

# Technology Scaling

- **Interactive Cloud apps rely on single threaded performance**

- **CPUs aren't getting much faster. We just get a few more**

- **2x users requires ~2x the number of servers**

## Single-Threaded Integer Performance
Based on adjusted SPECint® results

**~3%**

**~20% / year**

**~50%+ / year**

Legend:
- Intel Xeon
- Intel Core
- Intel Pentium
- Intel Itanium
- Intel Celeron
- AMD FX
- AMD Opteron
- AMD Phenom
- AMD Athlon
- IBM POWER
- PowerPC
- Fujitsu SPARC
- Sun SPARC
- DEC Alpha
- MIPS
- HP PA-RISC

Jeff Preshing, Henk Poley, http://preshing.com/20120208/a-look-back-at-single-threaded-cpu-performance/

# Datacenter Scaling



2013    2014    2015    2016    2017

~100%+ Growth for the past 5 years

2016

2018

Tapped Branch

Herbert Dr

Herbert D

Microsoft Data Center Boydton

Microsoft Data Center

608

58

# Cloud Server Changes

| | 2012 | 2018 | Ratio |
|---|---|---|---|
| CPU Cores | 16 | 36 | 2.25x |
| Storage | 4 TB HDD | 7 TB SDD (M.2) (120TB HDD*) | 1.75x 30x |
| Network | 1Gb | 50Gb | 50x |

\* - Bing SKU

# Cloud Server Changes

| | 2012 | 2018 | Ratio |
|---|---|---|---|
| CPU Cores | 16 | 36 | 2.25x |
| Storage | 4 TB HDD | 7 TB SDD (M.2) (120TB HDD*) | 1.75x 30x |
| Network | 1Gb | 50Gb | 50x |

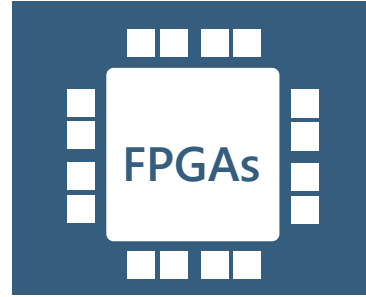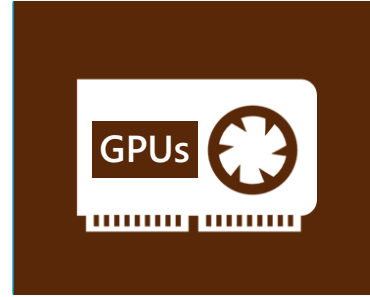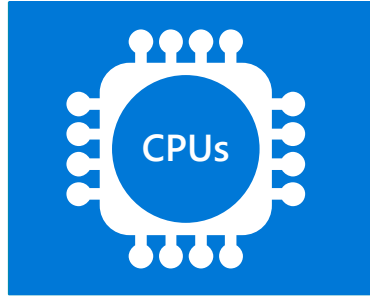* - Bing SKU

## Increasing Server Lifetimes

## More Data, Coming in Faster... and CPUs aren't keeping up

# Performance & Efficiency via Specialization



Source: Bob Broderson, Berkeley Wireless group
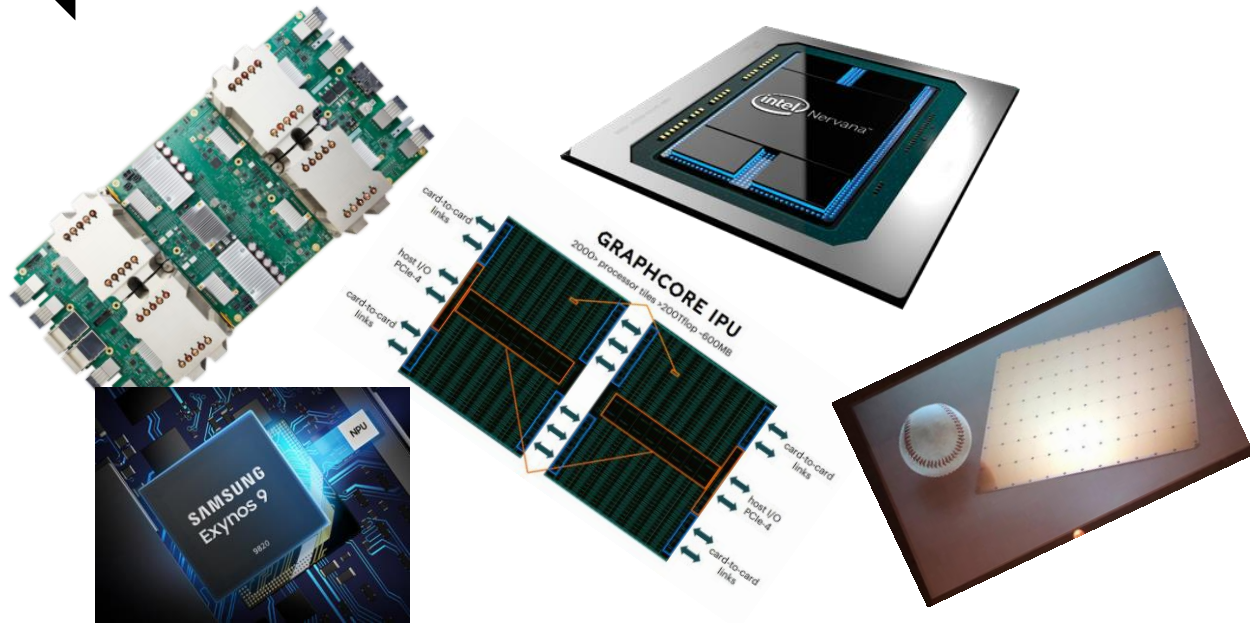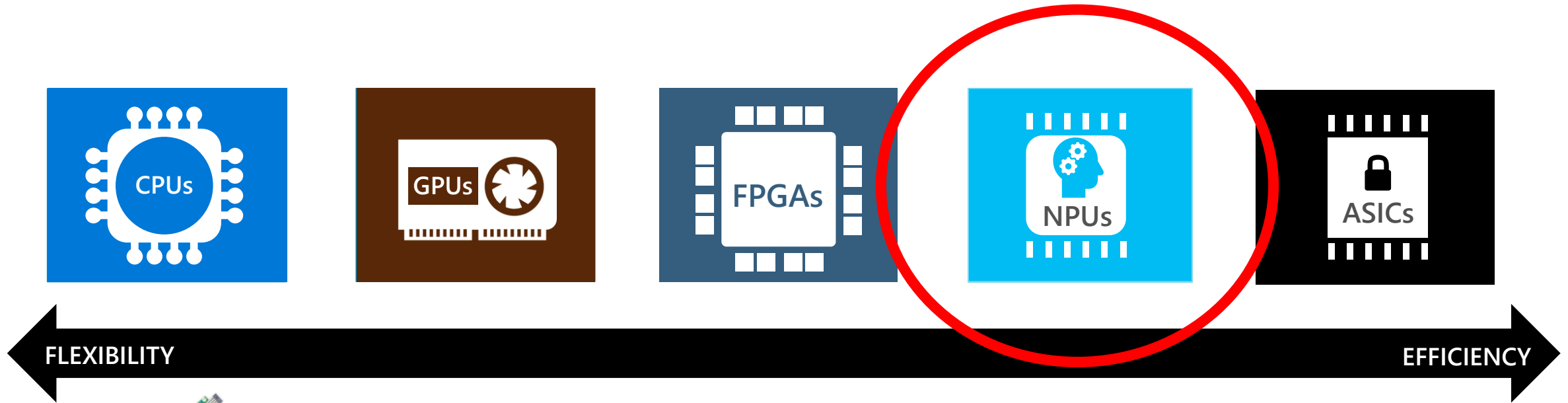
# Microsoft's Heterogeneous Cloud



CPUs · GPUs · FPGAs · NPUs · ASICs

FLEXIBILITY ← → EFFICIENCY

# Microsoft's Heterogeneous Cloud

FLEX ◄                                    ► EFFICIENCY

**NPUs**

**ASICs**

Ideal for ASICs

100

% of servers

Cloud Applications

0

1  2  3  4  5  6

Time Workload is Stable (years)

- Slow time to market
- Unlikely to last 6+ years
- ...ficult

Pokémon GO

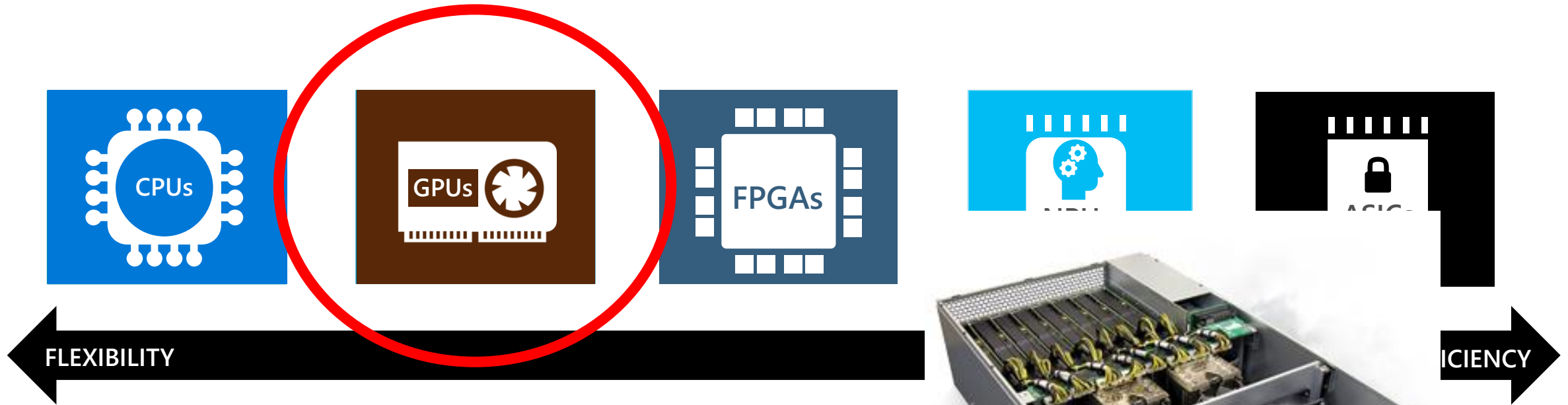# Microsoft's Heterogeneous Cloud

| CPUs | GPUs | FPGAs | NPUs | ASICs |

**FLEXIBILITY** ← → **EFFICIENCY**

- Very fast moving space
- Immature software ecosystem
- No clear path for development across generations
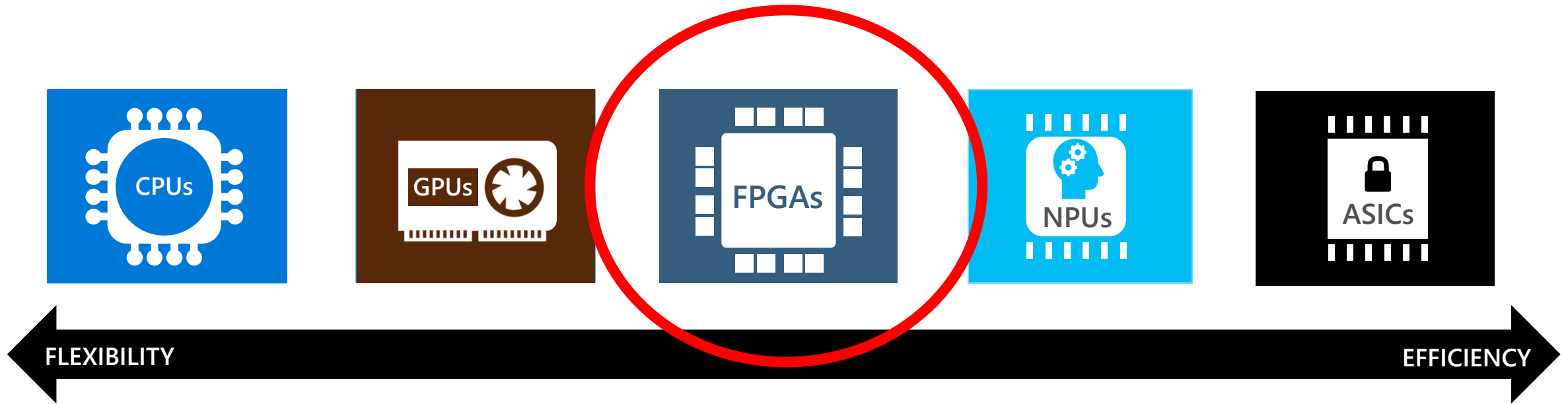
# Microsoft's Heterogeneous Cloud



**FLEXIBILITY** ⟷ **...ICIENCY**

- Great for HPC on batch data
- Stable software ecosystem allows easy adoption
- Power consumption limits deployment size
- Not ideal for latency-sensitive workloads

EnterPrise-XR-P

# Microsoft's Heterogeneous Cloud

CPUs | GPUs | FPGAs | NPUs | ASICs

FLEXIBILITY ⟵———————————————————⟶ EFFICIENCY

# What are FPGAs?

**<u>F</u>ield <u>P</u>rogrammable <u>G</u>ate <u>A</u>rray**

**FPGAs** are a sea of generic logic and interconnect

"Silicon Legos" – build them into exactly the right circuit for each task

Special-purpose hardware (FPGAs) is faster and more efficient than general-purpose hardware (CPUs)
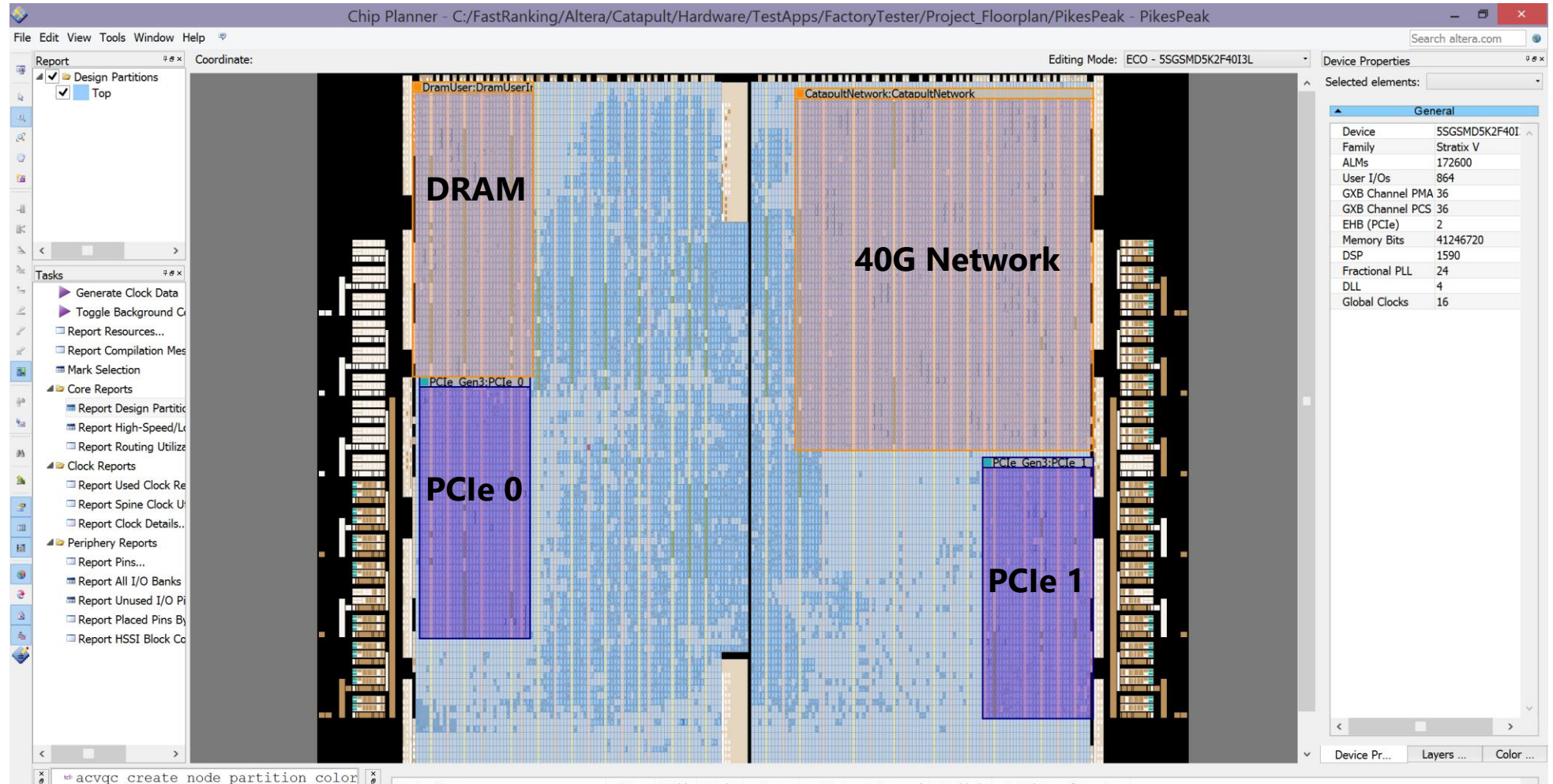
Low power compared to CPU/GPU

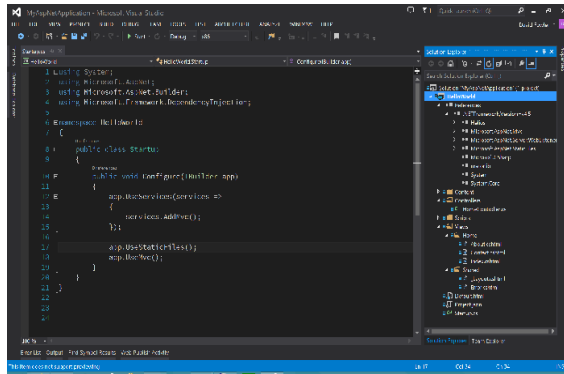**Change the hardware anytime!**

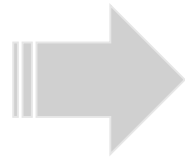100 ms to 1 second reconfiguration time

# FPGA Physical Layout
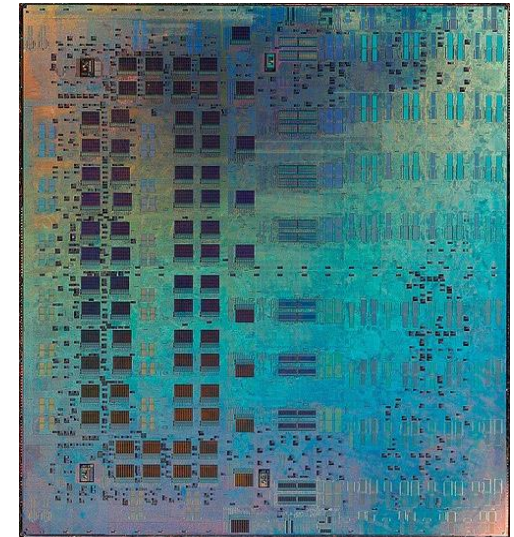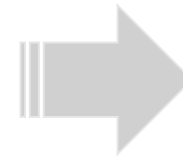


**Customize both the processing logic *and* the I/O**

# Gradual Migration to ASIC



Software

FPGA

ASIC

ASIC Integration via Chiplets

Azure

# Why is the FPGA a good choice as an accelerator?

- **Greater Performance and Efficiency than CPU, more general purpose than ASIC**

- **Many applications aren't about throughput or double-precision floating point**

  - AI/ML, Bioinformatics, text processing, financial services...

- **Exploits different forms of parallelism than other accelerators**

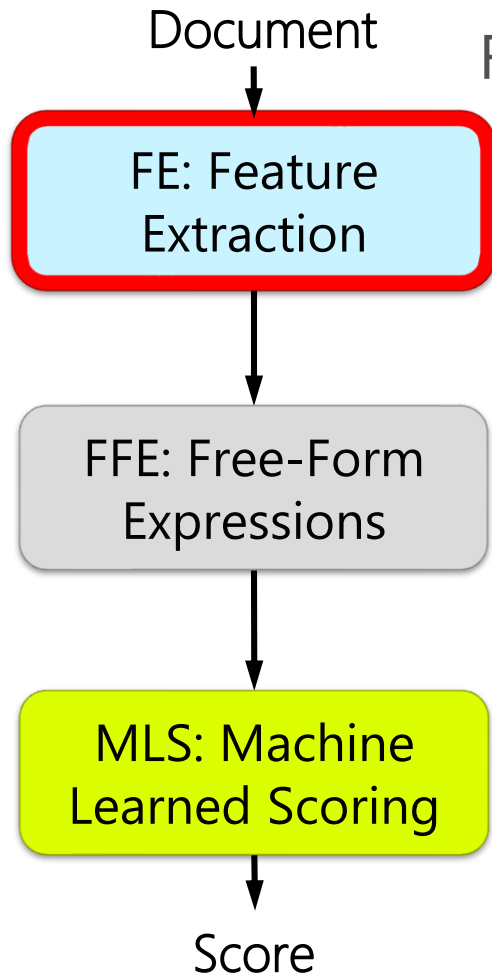**Multiple instruction streams, single data stream (MISD)**

Main article: MISD

[1] Multiple instructions operate on one data stream. This is an uncommon architecture which is generally used for fault tolerance. Heterogeneous systems operate on the same data stream and must agree on the result. Examples include the Space Shuttle flight control computer. [5]

[27]

**# Instruction Streams**

| | Single | Multiple |
|---|---|---|
| **Single** | **SISD** *No Parallelism* CPU | **SIMD** *Same thing to lots of data* GPUs (FP) FPGAs (Int) |
| **Multiple** | **MISD** *Different ops to same data* FPGAs | **MIMD** *Embarrassingly Parallel* Cluster |

**# Data streams**

© M                                                                                                    Azure

# FE: Feature Extraction

**Query: "FPGA Configuration"**

Document

Features:

| NumberOfOccurrences_0 = 7 | NumberOfOccurrences_1 = 4 | NumberOfTuples_0_1 = 1 |
|---|---|---|

FE: Feature Extraction

FFE: Free-Form Expressions

MLS: Machine Learned Scoring

Score

# Feature Extraction Accelerator

Document → Features

SM1
SM2 AllSpans
SM3 NumberOfOccurrences
SM4
SM5
SM6
SM7
SM8
SM9 NumberPerfectMatches
SM10

# FPGAs in Physics Applications



SETI

# FPGAs in Cosmology





EOR Science can be done with a paperclip and a supercomputer
-- Don C. Backer

Cosmologists often refer to their telescopes as "software telescopes"

# Processing Pipeline

Scientific Computing

HLS, ML

Detection,
Simple classification,
Triggering

Microsoft Catapult

ML at the source

Persistent DNNs,
Complex Learning,
In-Network Processing

Cloud:     IoT              Edge          Fog           Cloud

HPC:     Sensors        Detectors    Clusters     Supercomputer

# Catapult: Long, Fruitful FPGA Investment

**Catapult v1: Mt Granite**
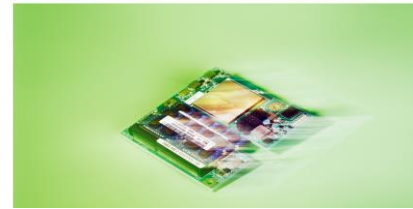Distributed solution
Integrated with WCS (OCP) 1.0

**v2: Pikes Peak**
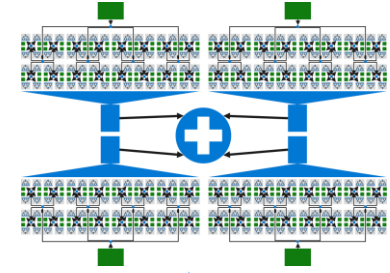Integrated Bing + Azure design
Bump-in-the-wire introduced

**Azure AccelNet Unveiled**
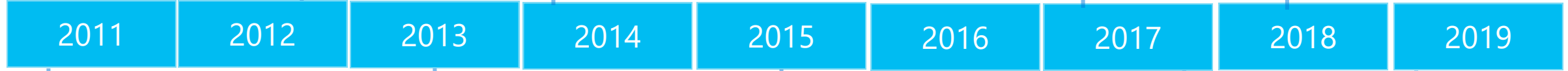Azure production launch
AI Supercomputer demo

**Project BrainWave / Storm Peak**
Real-Time AI
First 3rd Party FPGA Service

**Pre-History:**
May 2009: Bing Launched
Feb 2010: Azure Launched
Dec 2010: Catapult concept

MICROSOFT BETS ITS FUTURE ON A REPROGRAMMABLE COMPUTER CHIP

...

| 2011 | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 | 2018 | 2019 |

...

**v1: Scale Pilot**
1632 servers deployed
Bing IndexServe accelerated

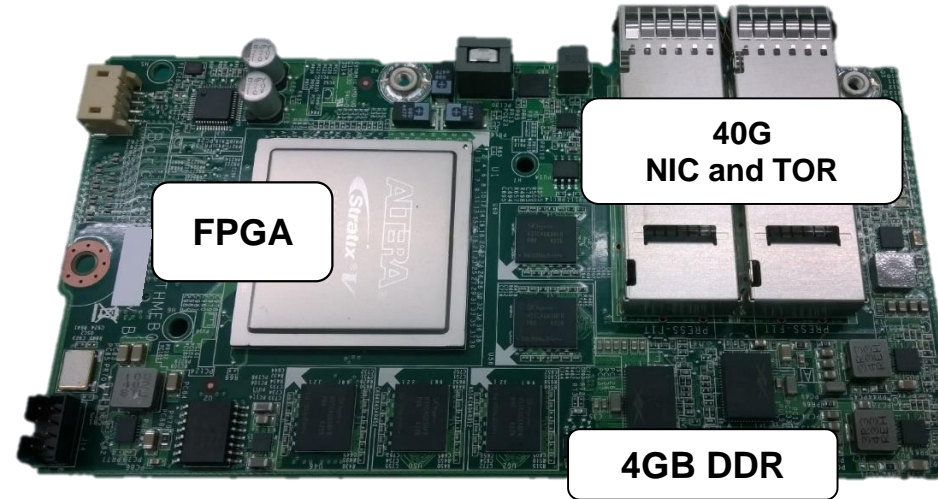MICROSOFT SUPERCHARGES BING SEARCH WITH PROGRAMMABLE CHIPS

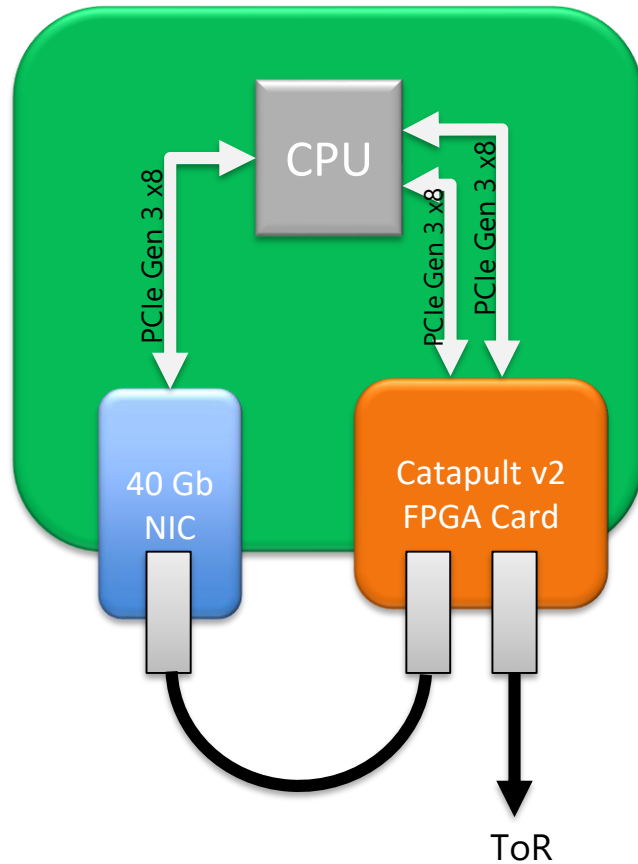**Catapult v3: Longs Peak**
DNN Platform for Bing
50Gb w/ integrated NIC

**Azure Databox Edge**
On-Site Inference

**v0: Research POC**
Built v0 board w/6 Xilinx FPGAs
30k lines of Bing code on FPGA

**v2 Production and ramp**
FPGAs reach production
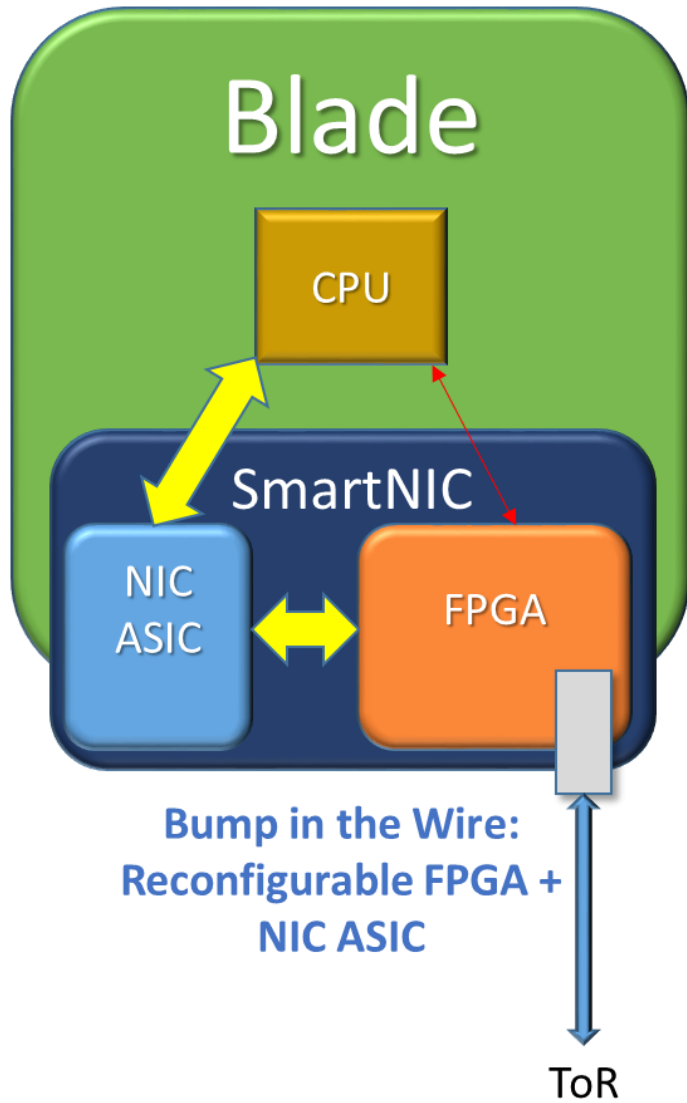Deployed in all new servers

Azure

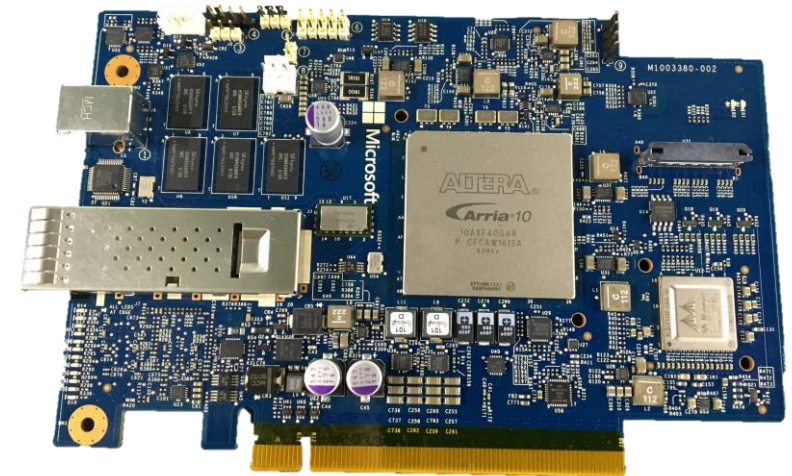# Catapult v2 – Bump in the Wire



Pikes Peak

Storey Peak

# Catapult v3 – Converged Boards



Dragontail Peak

Longs Peak

# Accelerator Integration



Traditional Integration



Bump in the Wire

Azure

# Basic Catapult Architecture



0.5m QSFP cable from NIC to FPGA

NIC

FPGA

~3m QSFP cable from FPGA to TOR

- Bump-in-the-wire architecture
- One FPGA in every server Microsoft has deployed since 2015

# New Generation Server Chassis



(6) N+2 Fans

Optional Remote Heatsink for high wattage CPUs

Up to (8) M.2 NVMe SSDs

Dual 3Φ PSU with Battery

945 mm

DDR4 DIMMs

441 mm

50G Networking

Next Gen CPUs

Up to (3) FHHL PCIe x16 Cards

Universal Motherboard

# Bump-in-the-wire Architectue

Azure ML
BrainWave

Azure
AccelNet

Hardware as a Service
Network Acceleration
Compute Acceleration

NIC

40G Ethernet

CPU

PCIe Gen 3

PCIe Gen 3

bing

© Microsoft Corporation

Azure

# Global-Scale FPGA



7 FPGAs Involved

2 FPGAs Involved

© Microsoft Corporation

Azure

Configurable Cloud

CPU compute layer

Reconfigurable compute layer

Converged network

# Network Latencies



- Extremely low latency (Similar to Infiniband)
- Global-scale FPGA

# HPC with the Cloud?

- The idea *sounds* great

- Pay for compute only when you use it

- When it breaks, it's someone else's problem

- No need to call the realtor / utility company when you want a bigger machine

- New hardware just shows up. No retrofits needed.

# Why hasn't Supercomputing moved to the Cloud?

❑ CPUs look largely the same

❑ Networks are highly specialized, tuned for low-latency, high bandwidth

    ❑ Proved this works FPGA-to-FPGA, but what about software (CPU-to-CPU)?

❑ Supercomputers include specialized accelerators (especially GPUs)

❑ Won't running virtual machines kill performance?

# Network Acceleration – Azure Accelerated Networking

| VM / Container (10.1.1.1) |
| OS |
| NIC |

| VM / Container (10.1.1.2) |
| OS |
| NIC |

10.1.1.2

# Virtualization Overhead – Standard Virtual Machines



VM
(10.1.1.1)

Guest OS

Hypervisor

SDN Stack

NIC

10.1.1.2

10.1.1.1

157.2.21.4

157.2.21.4

Virtualization Overhead

SDN Stack Functions
SLB: Software Load Balancers
ACL: Access Control Lists
NAT: Network Address Translation
VNET: Virtual Networks
Metering

# Virtualization Overhead – SRIOV NICs

# Virtualization Overhead – Azure SmartNIC w/ FPGAs

# Virtualization Overhead – Azure SmartNIC w/ FPGAs & DPDK

# AccelNet Performance



Lowest latency, highest-bandwidth network in the Cloud... for a while

Same approach can work for storage I/O

# Distributed Hardware-as-a-Service

- **Use multiple FPGAs to perform a complex task**

- **Head node is the gateway**

- **Use as many nodes as necessary – calling function doesn't need to know**

Image

Result

Pretrained DNN Model
(CNTK, TensorFlow, etc…)

BrainWave Compiler & Mapper

L1

Network switches

L0

L0

FPGAs

F F F F F F

Scalable DNN Hardware
Microservice

Azure

# Project BrainWave Concept



Traditional Approach

Project Brainwave

# Brainwave System at Cloud-Scale

# Project Brainwave – Fastest Image Classification

- 8 billion operations per image

- 1.5 ms on Resnet-50 with batch size of 1 (Realtime AI)

- Available for anyone to upload their models and run



10 CPU cluster 40 images/sec
(Single CPU 4 images/sec)

Single FPGA chip
550 images/sec

137 times faster than single CPU

Azure

https://github.com/Azure/aml-real-time-ai

# Bing Intelligent Search Backed By Project Brainwave

Stratix V RNN-optimized NPUs
in Scale Production

DECEMBER 13 2017

Bing launches new intelligent search features, powered by AI

Today we announced new Intelligent Search features for Bing, powered by AI, to give you answers faster, give you more comprehensive and complete information, and enable you to interact more naturally with your search engine.
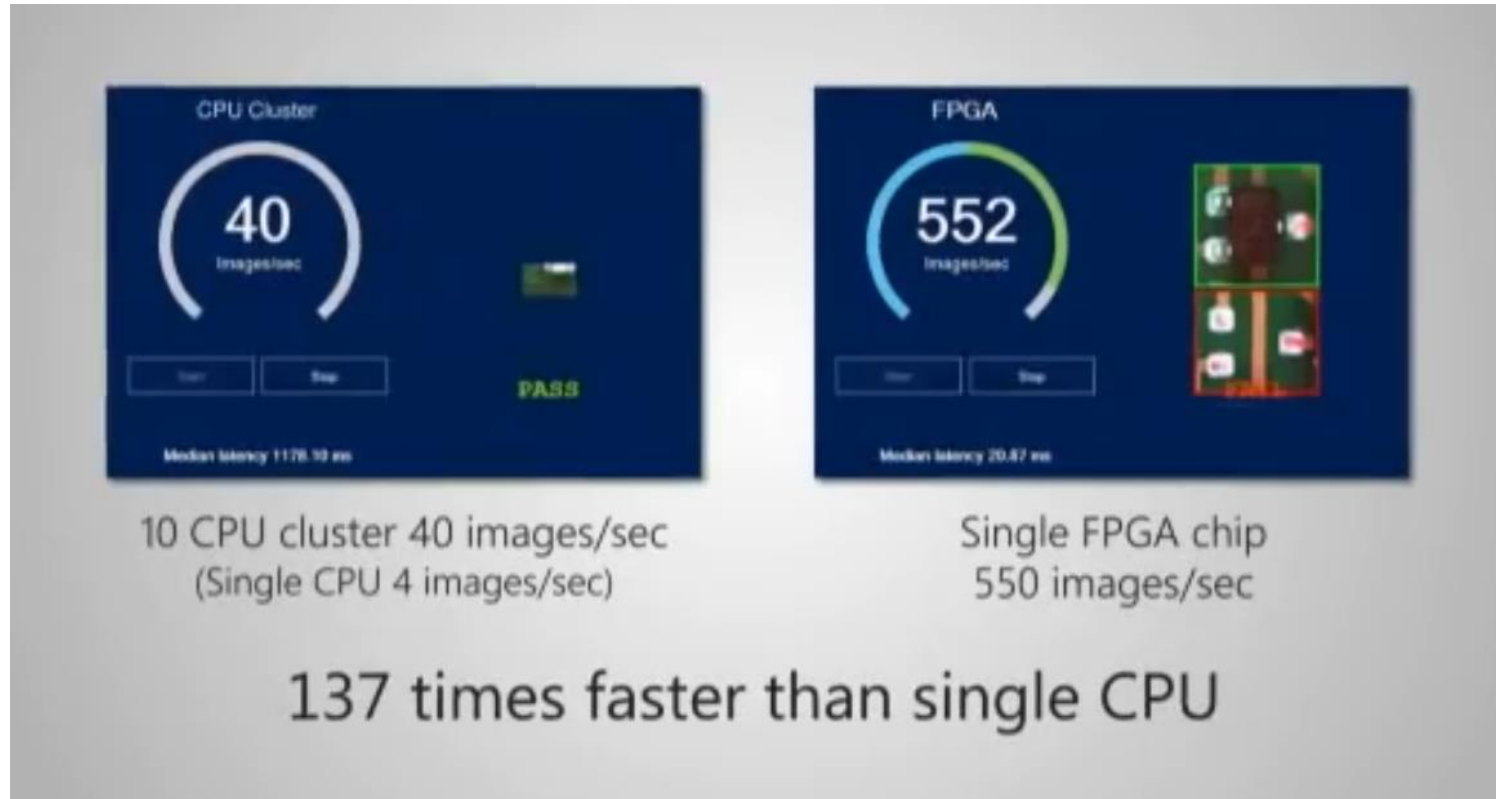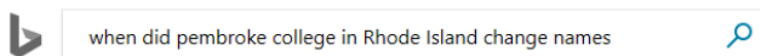
Intelligent answers:

Intelligent answers leverage the latest state of the art machine reading comprehension, backed by Project Brainwave running on Intel's FPGAs, to read and analyze billions of documents to understand the web and help you more quickly and confidently get the answers you need.

Bing now uses deep neural networks to validate answers by aggregating across multiple reputable sources, rather than just one, so you can feel more confident about the answer you're getting.

when did pembroke college in Rhode Island change names

All    Images    Videos    Maps    News    Shop    |    My saves

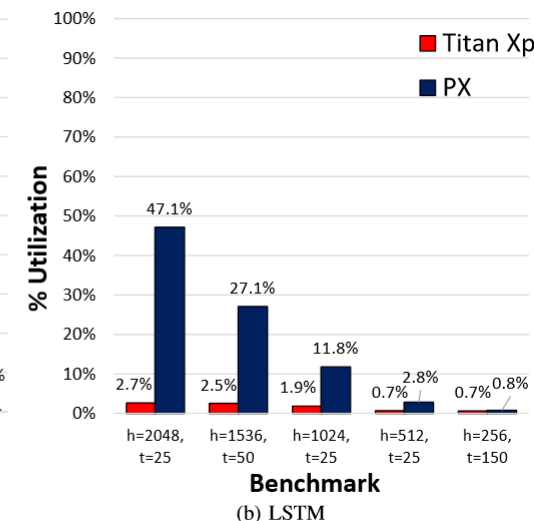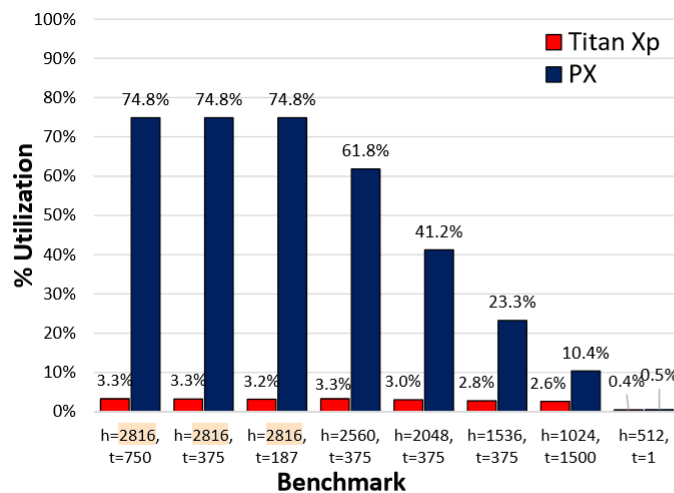281,000 Results    Any time ▾

1928
Consolidated from multiple sources

In 1928, the Women's College was renamed "Pembroke College in Brown University" in honor of Pembroke College at the University of Cambridge in England. Roger Williams, one of the founders of Rhode Island, was an alumnus of Cambridge's Pembroke.

Pembroke College in Brown University - Wikipedia
en.wikipedia.org

Similar answer at: brown.edu

| Bing TP1 | | | |
|---|---|---|---|
| | **CPU-only** | **Brainwave-accelerated** | **Improvement** |
| Model details | GRU 128x200 (x2) + W2Vec | LSTM 500x200 (x8) + W2Vec | Brainwave-accelerated model is > 10X larger and > 10X lower latency |
| End-to-end latency per Batch 1 request at 95% | 9 ms | 0.850 ms | |

| Bing DeepScan | | | |
|---|---|---|---|
| | **CPU-only** | **Brainwave-accelerated** | **Improvement** |
| Model details | 1D CNN + W2Vec (RNNs removed) | 1D CNN + W2Vec + GRU 500x500 (x4) | Brainwave-accelerated model is > 10X larger and 3X lower latency |
| End-to-end latency per Batch 1 request at 95% | 15 ms | 5 ms | |



(a) GRU

(b) LSTM

# Attend the Turotial tomorrow for more detail on BrainWave

**THURSDAY, 12 SEPTEMBER**

**09:00** → 12:15 **Tutorial** 📍 1 West

09:00 **Docker and Kubernetes** 🕐 30m

09:30 **Azure** 🕐 1h 15m
Speaker: Ted Way (Microsoft)

10:45 **Coffee** 🕐 30m

11:15 **Galapagos** 🕐 1h
Speakers: Dylan Sheldon Rankin (Massachusetts Inst. of Technology (US)), Naif Tarafdar (University of Toronto)
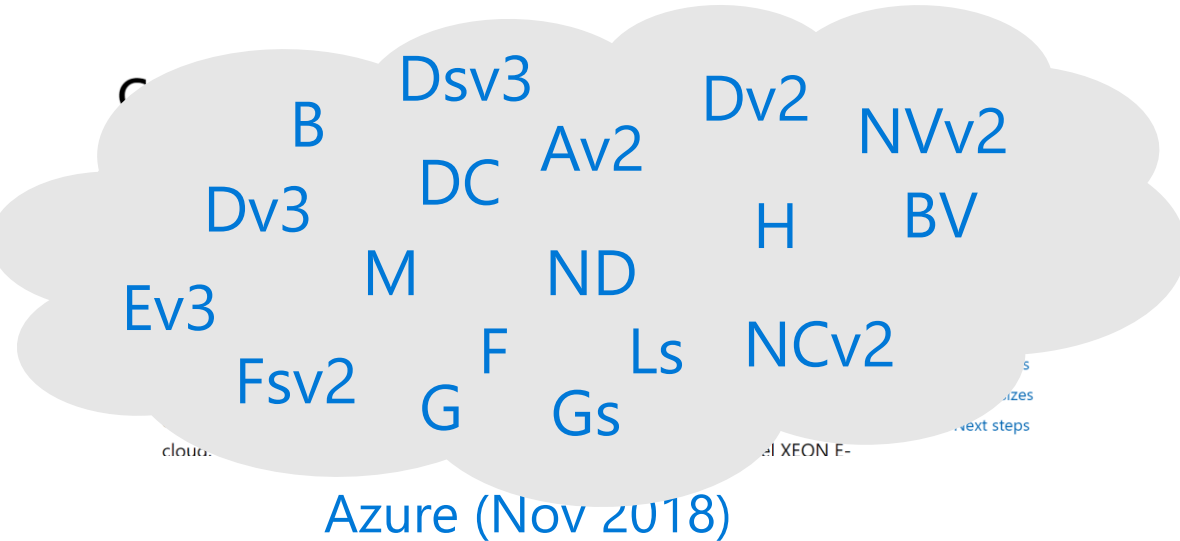
**14:00** → 17:00 **Tutorials, office hours** 📍 1 West

14:00 **hls4ml** 🕐 1h 15m
Speakers: Javier Mauricio Duarte (Univ. of California San Diego (US)), Sioni Paris Summers (CERN), Zhenbin Wu (University of Illinois at Chicago (US))

15:15 **Coffee** 🕐 30m

15:45 **hls4ml** 🕐 1h 15m

Azure

# Why will FPGAs still be here tomorrow?

- Performance? Low Power? Ubiquity? Killer features? AI/ML?

Azure (Nov 2018)

Dsv3, B, Dv2, NVv2, Av2, DC, H, BV, Dv3, M, ND, Ev3, F, Ls, NCv2, Fsv2, G, Gs

Amazon AWS (Nov 2018)

A1, D2, X1e, F1, M5, P3, T2, H1, C5, R4, I3

- Every time we change hardware, we need to make a new VM type

- FPGAs allow us to back new features into existing VM types

**Developers (Customers) don't want to spend time fighting old battles**