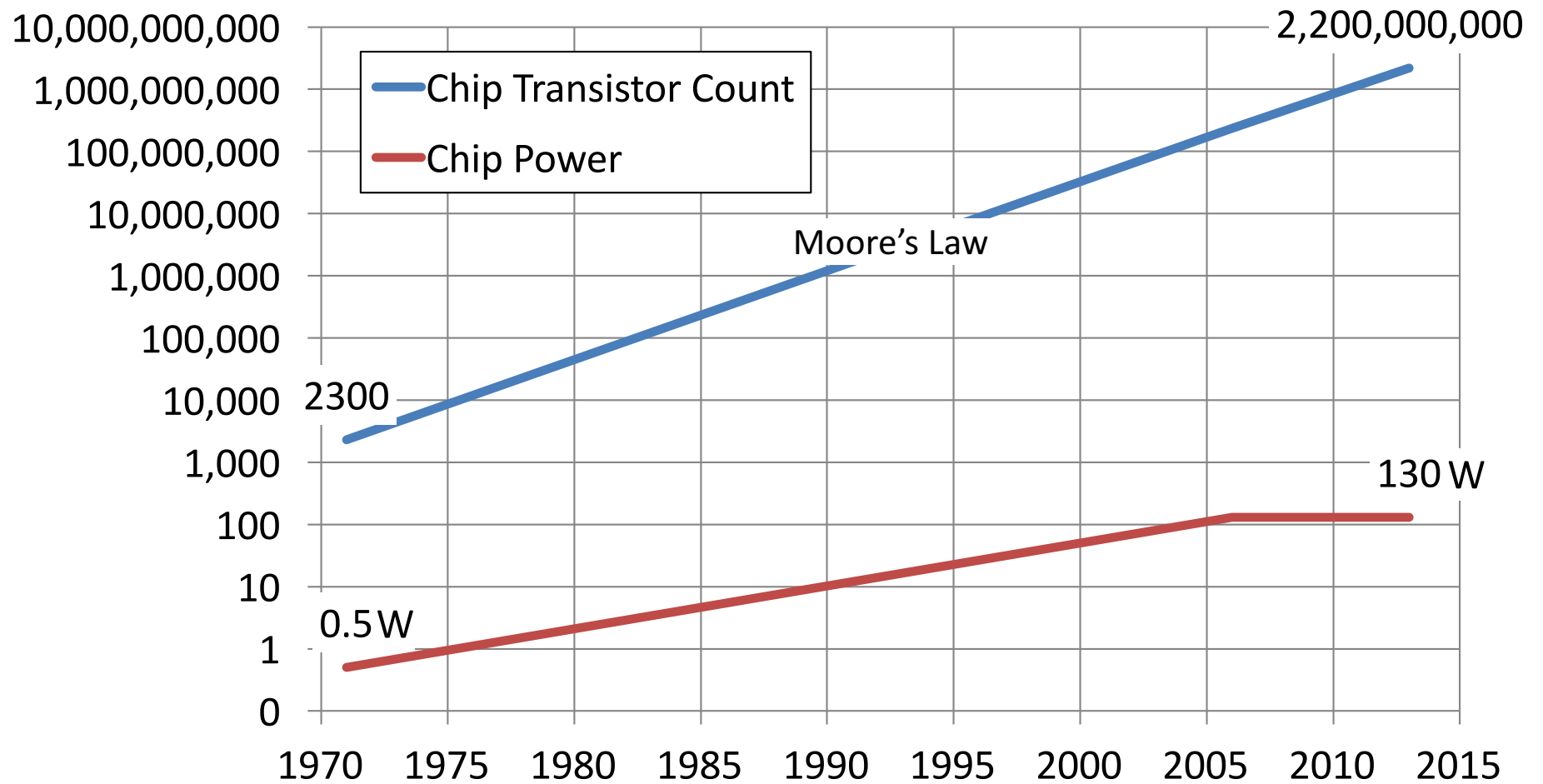# CSEP 548: Computer Systems Architecture

*Dark Silicon, Specialization, Systems for ML*

Luis Ceze, Spring 2017

(based on slides lifted from Me, Hadi Esmaeilzadeh, Michael Taylor, Carlo Del Mundo, Liang Luo and the interwebs at large)
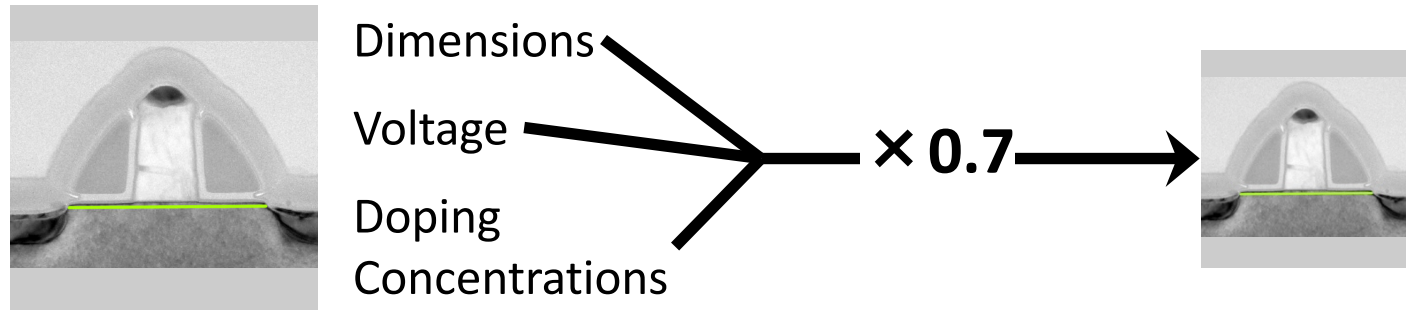
# What is the catch with Moore's law?



2

# Dennard scaling:
## Doubling the transistors; scale their power down

Transistor: 2D Voltage-Controlled Switch



Dimensions

Voltage

Doping Concentrations

$\times$ **0.7**

Area ——————— **0.5 $\times$ ↓** ———————→

Capacitance ——————— **0.7 $\times$ ↓** ———————→

Frequency ——————— **1.4 $\times$ ↑** ———————→

Power = Capacitance $\times$ Frequency $\times$ Voltage$^2$

Power ——————— **0.5 $\times$ ↓** ———————→

# Dark silicon
## What if you can't power them anymore?

Area ——————————— **0.5 × ↓** ——————————→
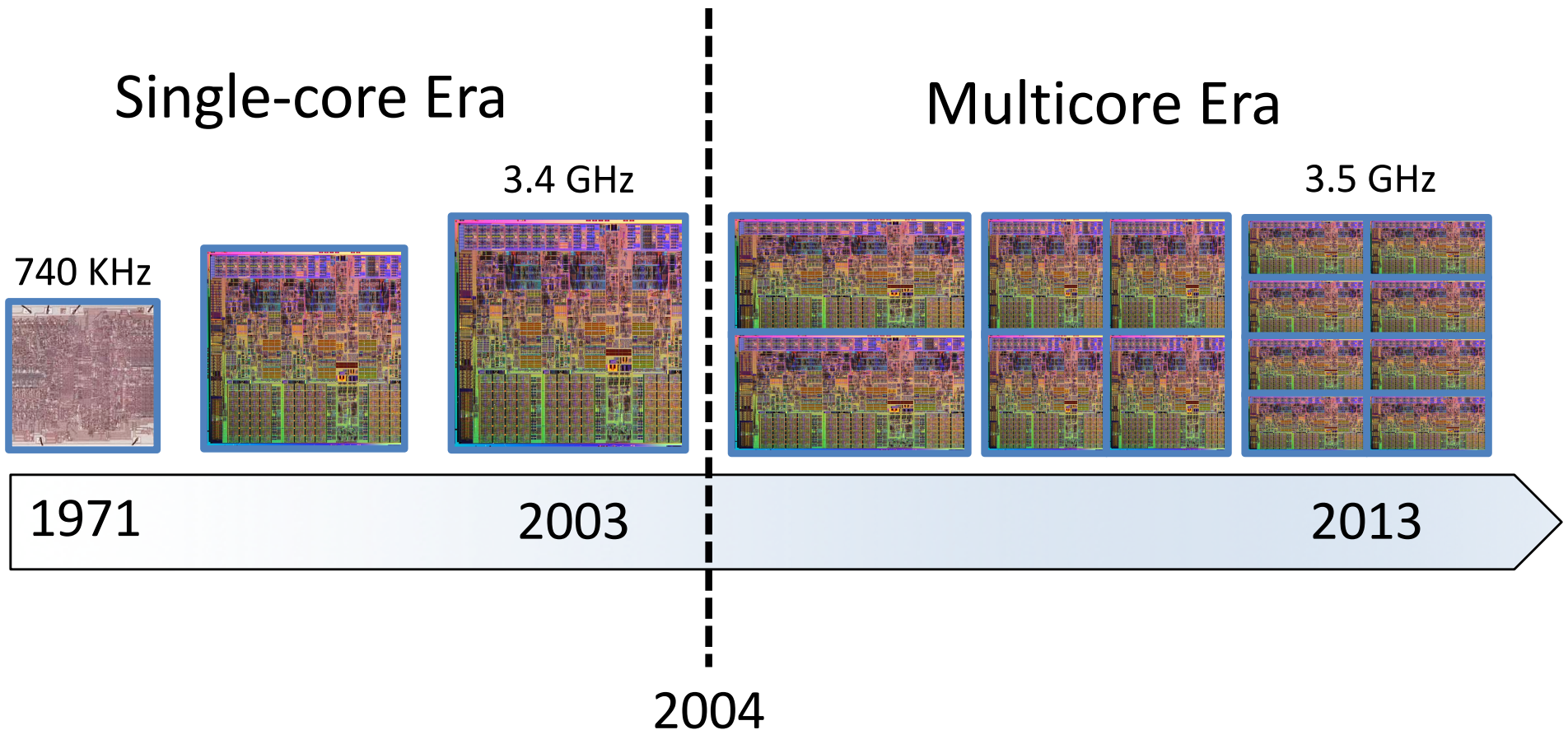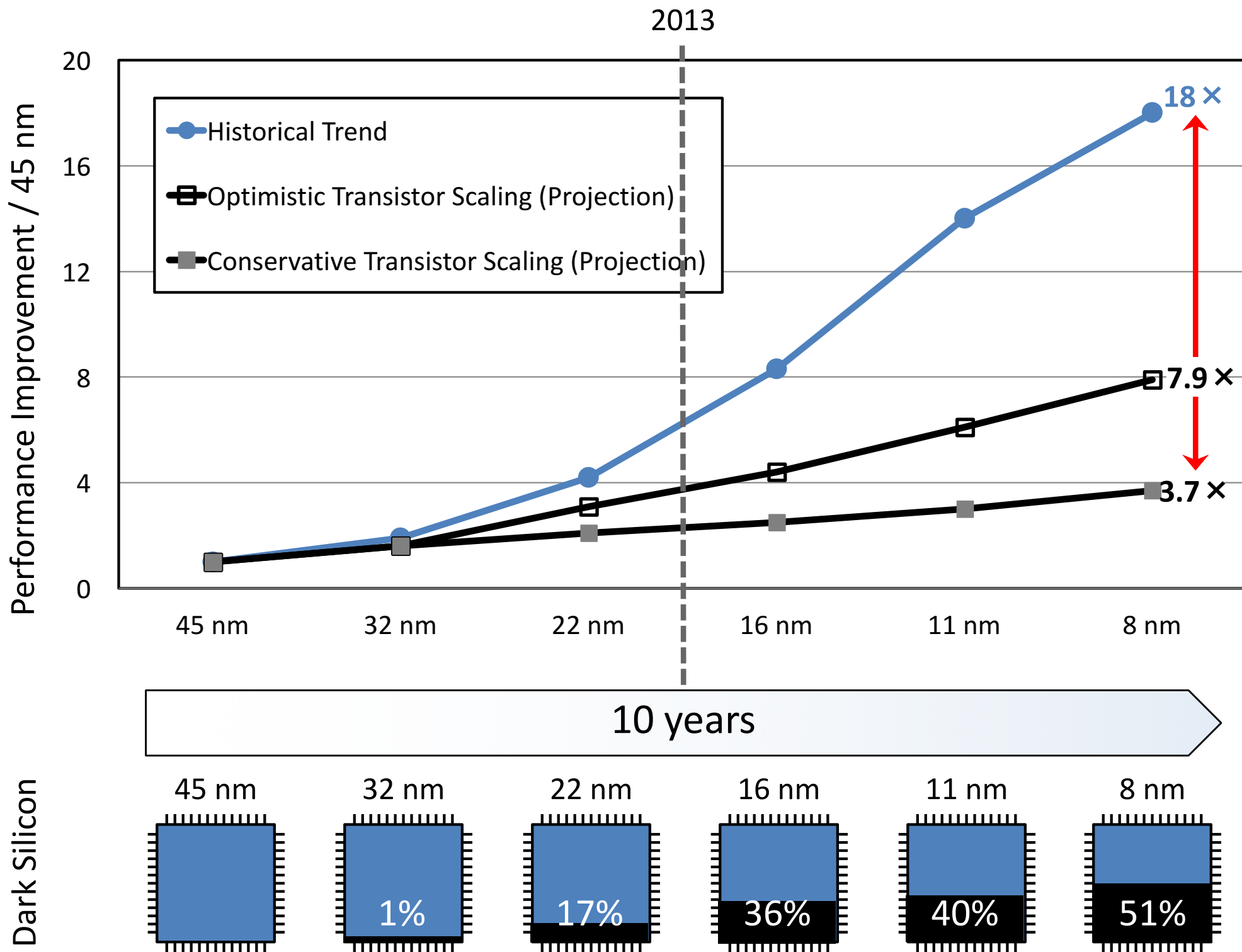
Power ——————————— **0.5 × ↓** ——————————→

⬇

### Dark Silicon

Can't turn all transistor on at the same time.
Part of the chip gets "dark".

# Looking back
## Evolution of processors

Single-core Era

Multicore Era

740 KHz

3.4 GHz

3.5 GHz

1971

2003

2013

2004

Is parallelism long-term solution?

2013

Performance Improvement / 45 nm

- Historical Trend
- Optimistic Transistor Scaling (Projection)
- Conservative Transistor Scaling (Projection)

18 ×
7.9 ×
3.7 ×

45 nm   32 nm   22 nm   16 nm   11 nm   8 nm

10 years

Dark Silicon

45 nm   32 nm   22 nm   16 nm   11 nm   8 nm
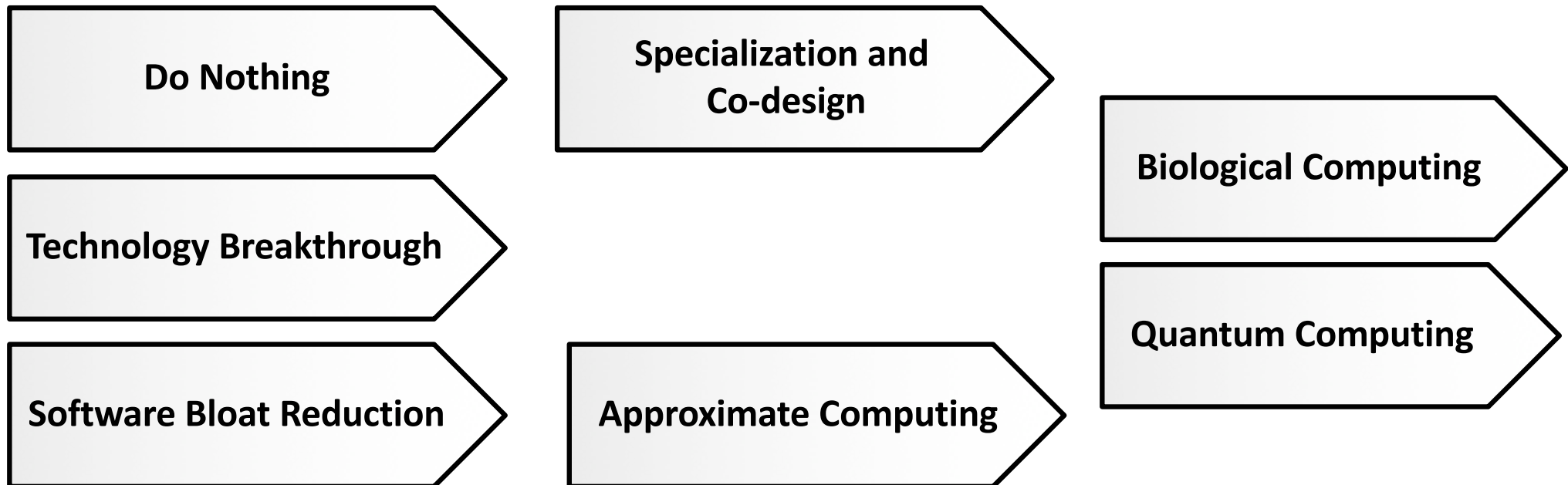
1%   17%   36%   40%   51%

# What now?

Unicore Era | Multicore Era | ?

Need at least 18%-40% per generation from architecture alone without additional power

# Possible paths forward

Do Nothing

Technology Breakthrough

Software Bloat Reduction

Specialization and Co-design

Approximate Computing

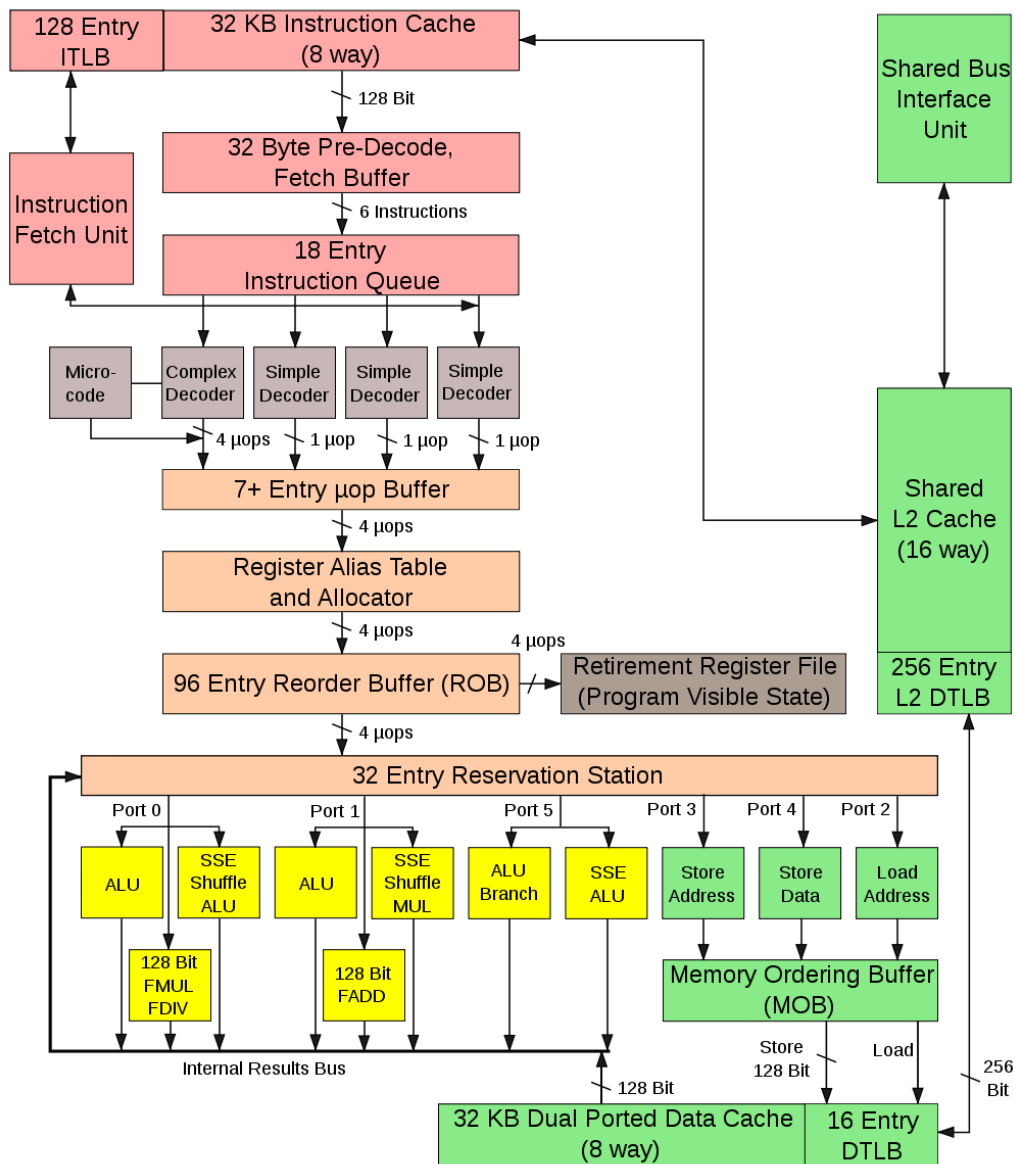Biological Computing

Quantum Computing

# Specialization and efficiency



Source: Bob Broderson, Berkeley Wireless group

Why?

# CPU



Intel Core 2 Architecture

# nVidia Fermi GPU

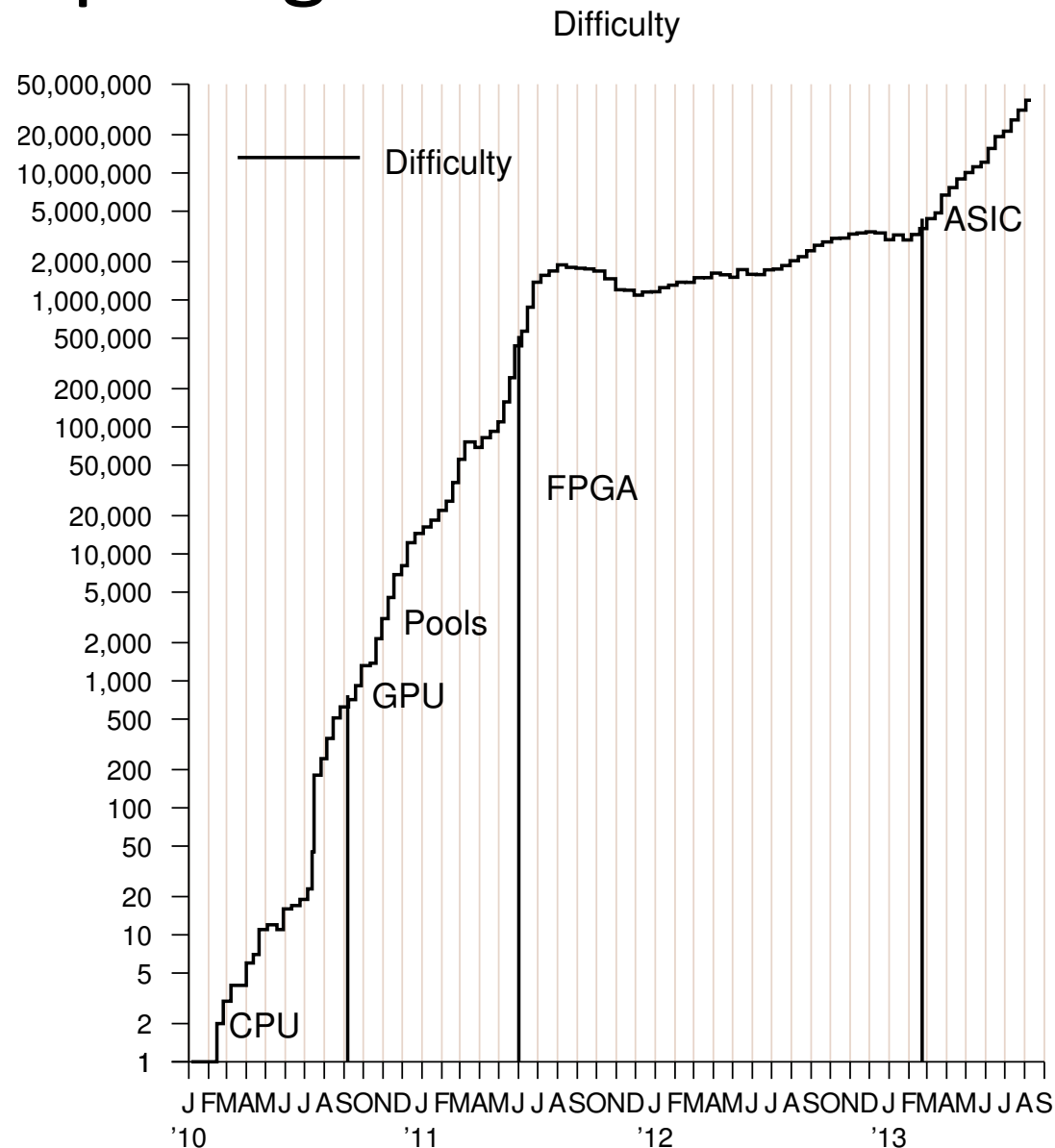# The Value of a Bitcoin



USD/BTC Exchange rate (in USD)

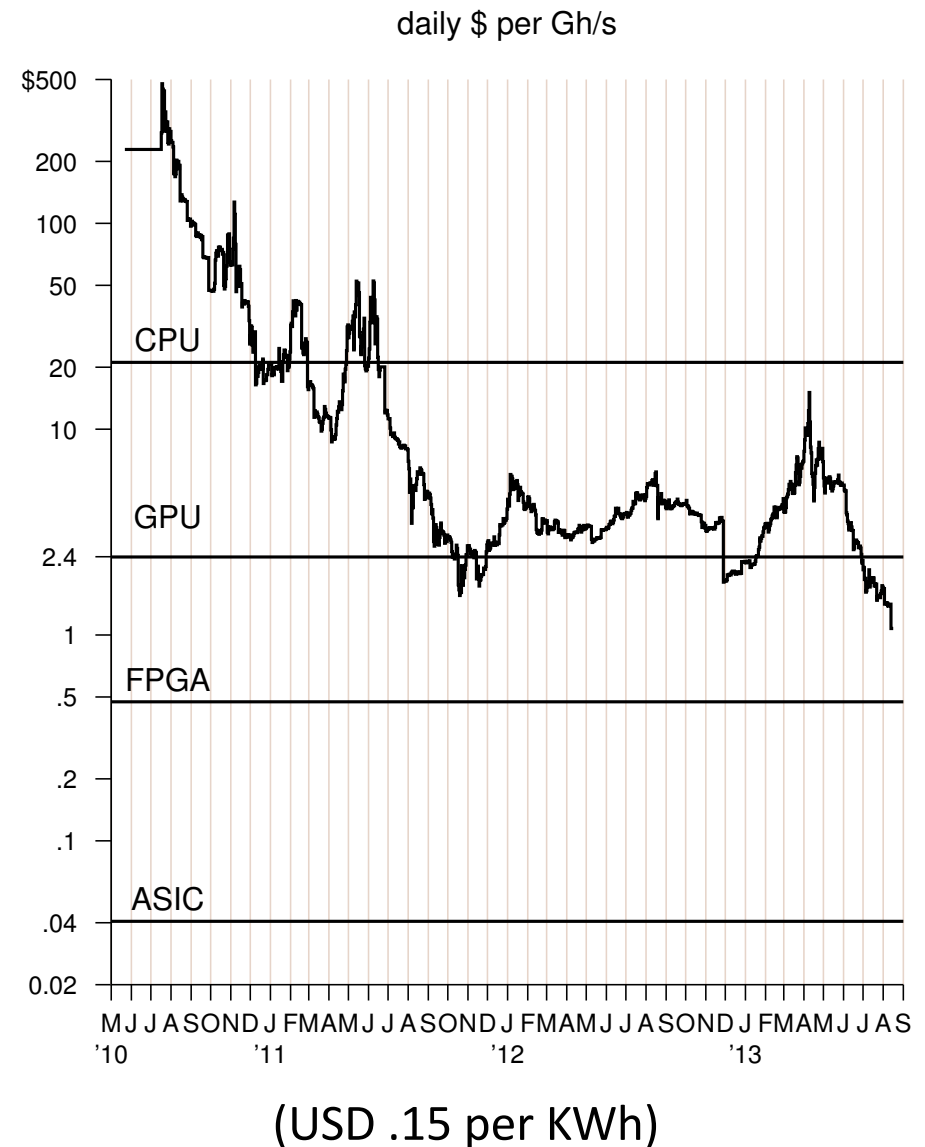# BTC Mining Computing Evolution

- ## CPU

- ## GPU

  – Portable OpenCL Imp

  – Completely unrolled double SHA256 hash

  – AMD >> Nvidia

  - instruction set match
  - microarch (VLIW) match
  - higher ALU density
  - memory BW not used

- ## FPGA

  – verilog

  – "gateway drug to ASIC": boards, protocols, thermals, verilog

- ## ASIC

Difficulty

# Energy Costs and USD/BTC
# Say when to unplug/plug HW

- daily $ per Gh/s falls as technology advances and more machines deployed

- daily $/GH/s rises if USD/BTC rises.

- Today, CPUs, GPUs, and even FPGAs do not recoup energy costs

- Rising USD/BTC: old machines get fired up.

- Steady state: cheap energy wins (Iceland?)

daily $ per Gh/s

$500
200
100
50
CPU
20
10
GPU
2.4
1
FPGA
.5
.2
.1
ASIC
.04
0.02

M J J A S O N D J F M A M J J A S O N D J F M A M J J A S O N D J F M A M J J A S
'10            '11                '12                '13

(USD .15 per KWh)

amazon
amazon
amazon
amazon
amazon

Shop by
Department

Home & Kitchen    Best Sellers    M

All

Prime

Prime

PrimeFresh

PrimePlus

Prime

PrimeFresh

Beddin

☑ AMD Sempron 145 Processor (SDX145HBGMBOX) $36.98

Pre-order

Pre-order t

Pre-order item

Trade in

dd-on Item

Buy
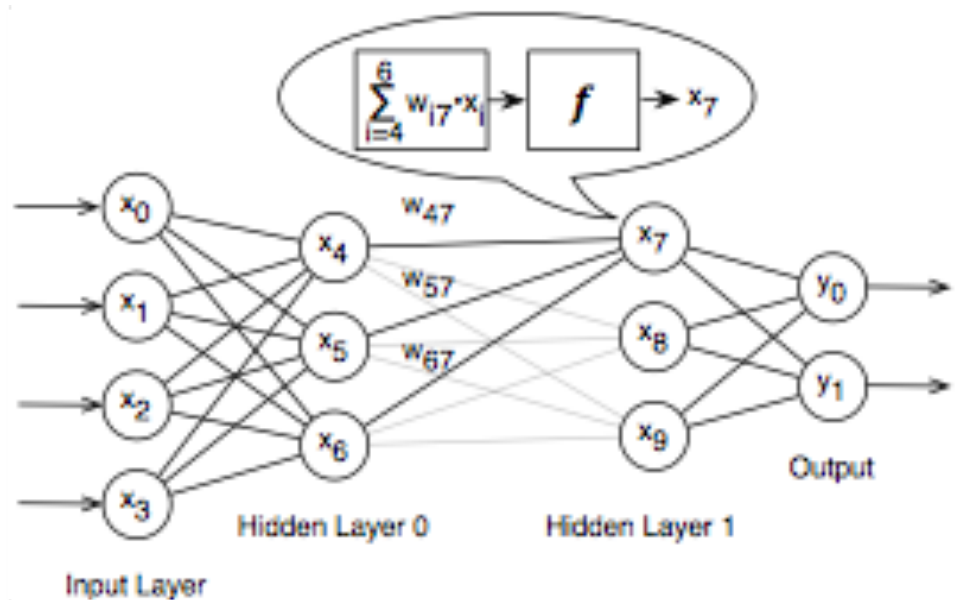
Check Ou

nVidia SLI/ ..

$4.98

ck®

# HW design in one slide

- Declare compute components, memory elements, interconnection

- "Place and route" distributes those in space
  - And checks is timing works --- i.e., all signals can be stable for a target clock frequency
  - Assess HW resource utilization, power consumption, etc.

```verilog
//----------------------------------------------------
// Design Name : parity_using_assign
// File Name   : parity_using_assign.v
// Function    : Parity using assign
// Coder       : Deepak Kumar Tala
//----------------------------------------------------
module parity_using_assign (
data_in     , //  8 bit data in
parity_out    //  1 bit parity out
);
output  parity_out ;
input [7:0] data_in ;

wire parity_out ;

assign parity_out =  (data_in[0] ^ data_in[1]) ^
                     (data_in[2] ^ data_in[3]) ^
                     (data_in[4] ^ data_in[5]) ^
                     (data_in[6] ^ data_in[7]);

endmodule
```
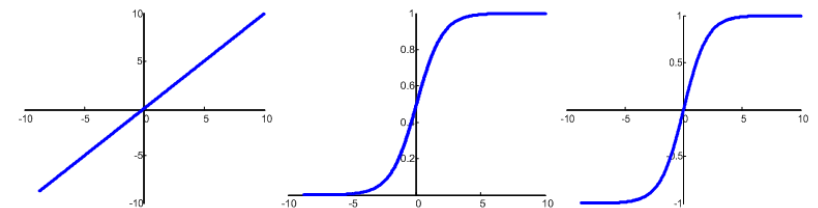
# Neural networks



neural network



computing a single layer

$$\begin{bmatrix} x_7 \\ x_8 \\ x_9 \end{bmatrix} = f\left( \begin{bmatrix} W_{67} & W_{57} & W_{47} \\ W_{68} & W_{58} & W_{48} \\ W_{69} & W_{59} & W_{49} \end{bmatrix} \begin{bmatrix} x_6 \\ x_5 \\ x_4 \end{bmatrix} \right)$$
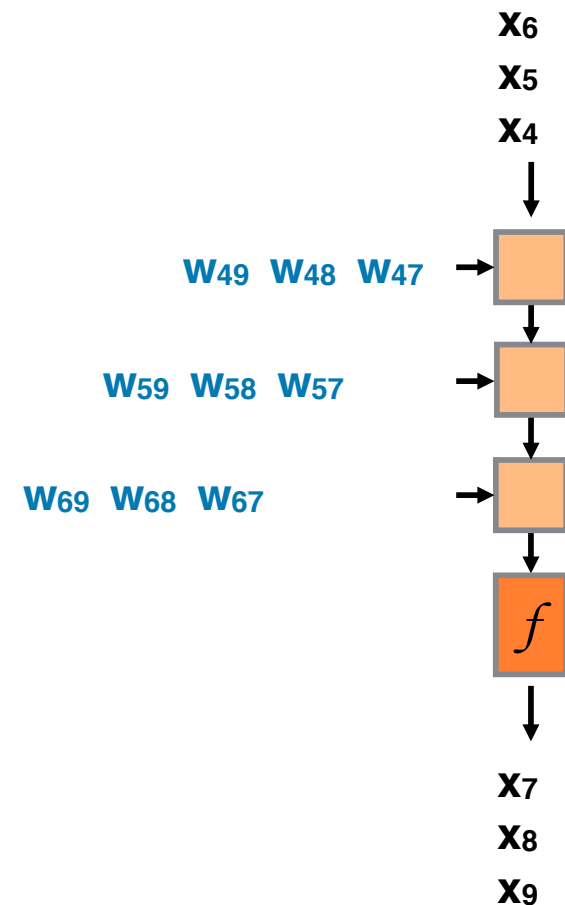
activation function $f$

# Systolic Arrays

*computing a single layer*
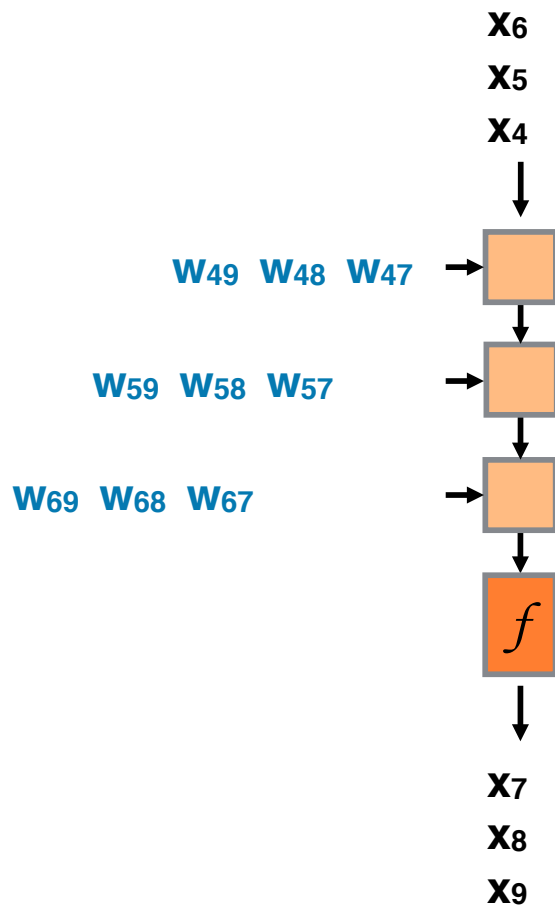
$$\begin{bmatrix} X_7 \\ X_8 \\ X_9 \end{bmatrix} = f\left( \begin{bmatrix} W_{67} & W_{57} & W_{47} \\ W_{68} & W_{58} & W_{48} \\ W_{69} & W_{59} & W_{49} \end{bmatrix} \begin{bmatrix} X_6 \\ X_5 \\ X_4 \end{bmatrix} \right)$$

*systolic array*

# Making it fast in HW

*systolic array*

*processing unit*

X6
X5
X4

W49  W48  W47 →

W59  W58  W57 →

W69  W68  W67 →

$f$

X7
X8
X9

**PU**

**control**

Storage

PE

PE

PE

PE

$f$

1 - processing elements in hardwired logic



2 - local storage for synaptic weights

3 - sigmoid unit implements non-linear activation functions



4 - vertically micro-coded sequencer

# Scaling it up

# Google's Tensor Processing Unit (TPU)

- **30-80x** TOPS/watt vs. 2015 CPUs and GPUs.
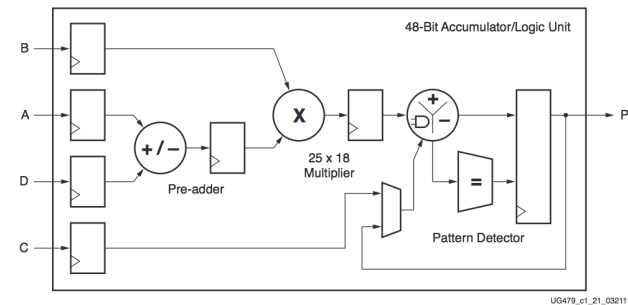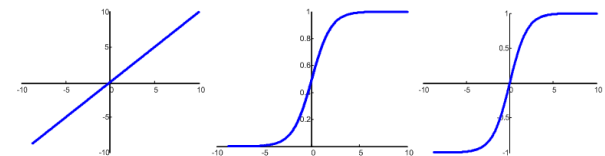
- 8 GiB DRAM.

- 8-bit fixed point.

- 256x256 MAC unit.

- Support for data reordering, matrix multiply, activation, pooling, and normalization.



**Figure 3.** TPU Printed Circuit Board. It can be inserted in the slot for an SATA disk in a server, but the card uses PCIe Gen3 x16.

# TPU Block Diagram & Floor Plan



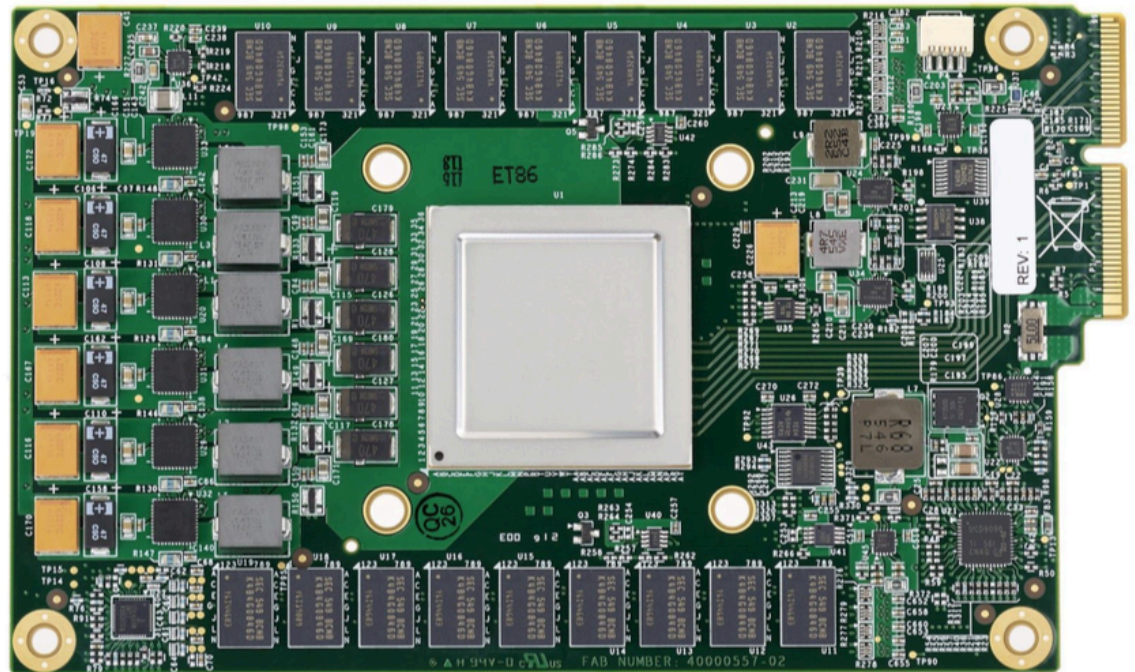**Figure 1.** TPU Block Diagram. The main computation part is the yellow Matrix Multiply unit in the upper right hand corner. Its inputs are the blue Weight FIFO and the blue Unified Buffer (UB) and its output is the blue Accumulators (Acc). The yellow Activation Unit performs the nonlinear functions on the Acc, which go to the UB.



**Figure 2.** Floor Plan of TPU die. The shading follows Figure 1. The light (blue) data buffers are 37% of the die, the light (yellow) compute is 30%, the medium (green) I/O is 10%, and the dark (red) control is just 2%. Control is much larger (and much more difficult to design) in a CPU or GPU

# Experimental Testbed

| Model | Die | | | | | | | | | | Benchmarked Servers | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $mm^2$ | nm | MHz | TDP | Measured Idle | Measured Busy | TOPS/s 8b | TOPS/s FP | GB/s | On-Chip Memory | Dies | DRAM Size | TDP | Measured Idle | Measured Busy |
| Haswell E5-2699 v3 | 662 | 22 | 2300 | 145W | 41W | 145W | 2.6 | 1.3 | 51 | 51 MiB | 2 | 256 GiB | 504W | 159W | 455W |
| NVIDIA K80 (2 dies/card) | 561 | 28 | 560 | 150W | 25W | 98W | -- | 2.8 | 160 | 8 MiB | 8 | 256 GiB (host) + 12 GiB x 8 | 1838W | 357W | 991W |
| TPU | NA* | 28 | 700 | 75W | 28W | 40W | 92 | -- | 34 | 28 MiB | 4 | 256 GiB (host) + 8 GiB x 4 | 861W | 290W | 384W |

**Table 2.** Benchmarked servers use Haswell CPUs, K80 GPUs, and TPUs. Haswell has 18 cores, and the K80 has 13 SMX processors. Figure 10 has measured power. The low-power TPU allows for better rack-level density than the high-power GPU. The 8 GiB DRAM per TPU is Weight Memory. GPU Boost mode is not used (Sec. 8). SECDEC and no Boost mode reduce K80 bandwidth from 240 to 160. No Boost mode and single die vs. dual die performance reduces K80 peak TOPS from 8.7 to 2.8. (*The TPU die is ≤ half the Haswell die size.)



**8x K80 GPUs**



**Figure 3.** TPU Printed Circuit Board. It can be inserted in the slot for an SATA disk in a server, but the card uses PCIe Gen3 x16.
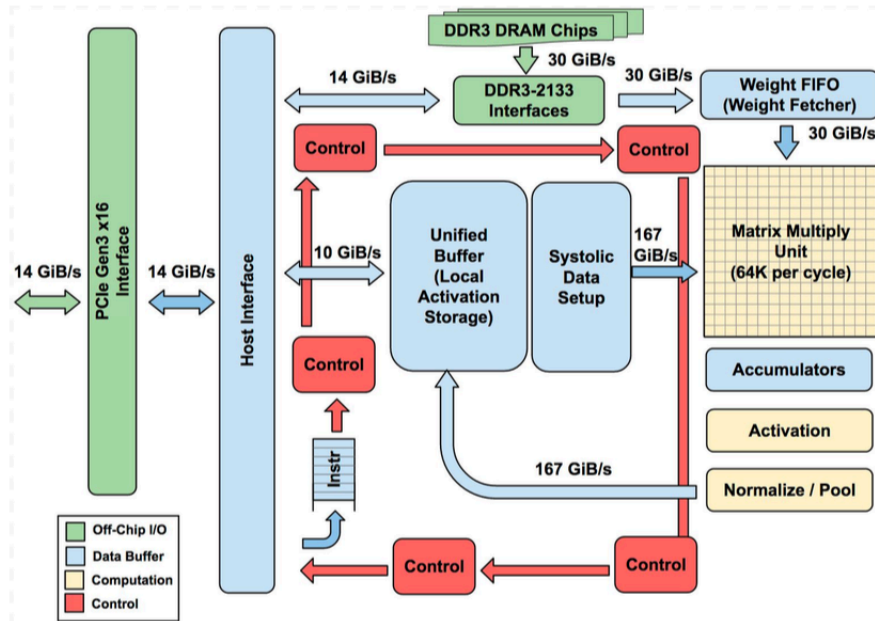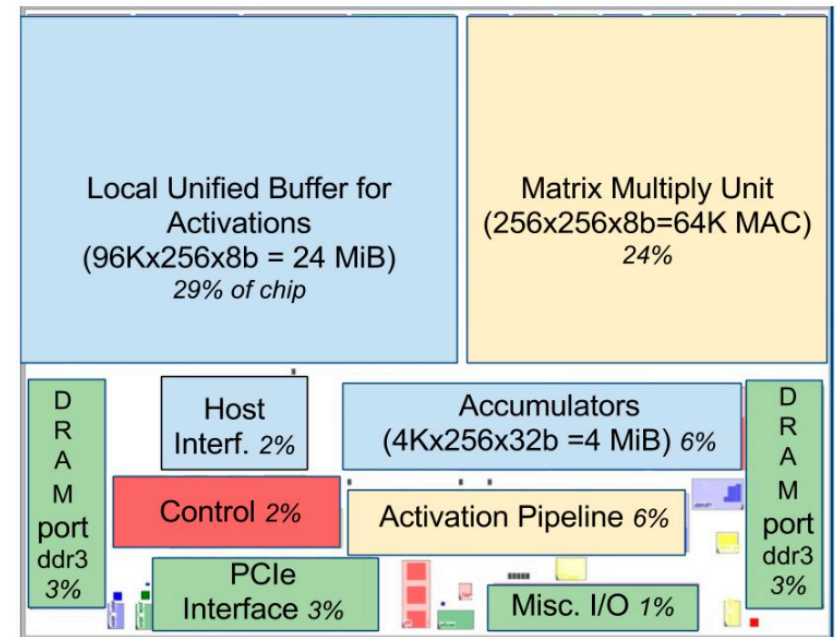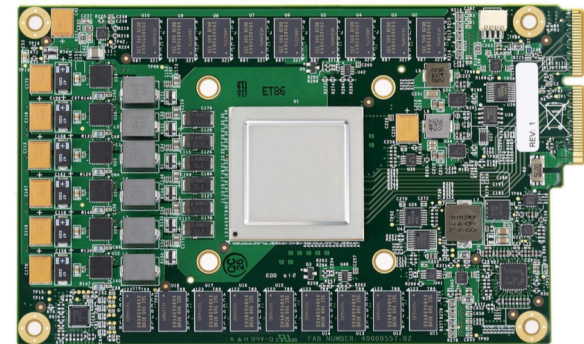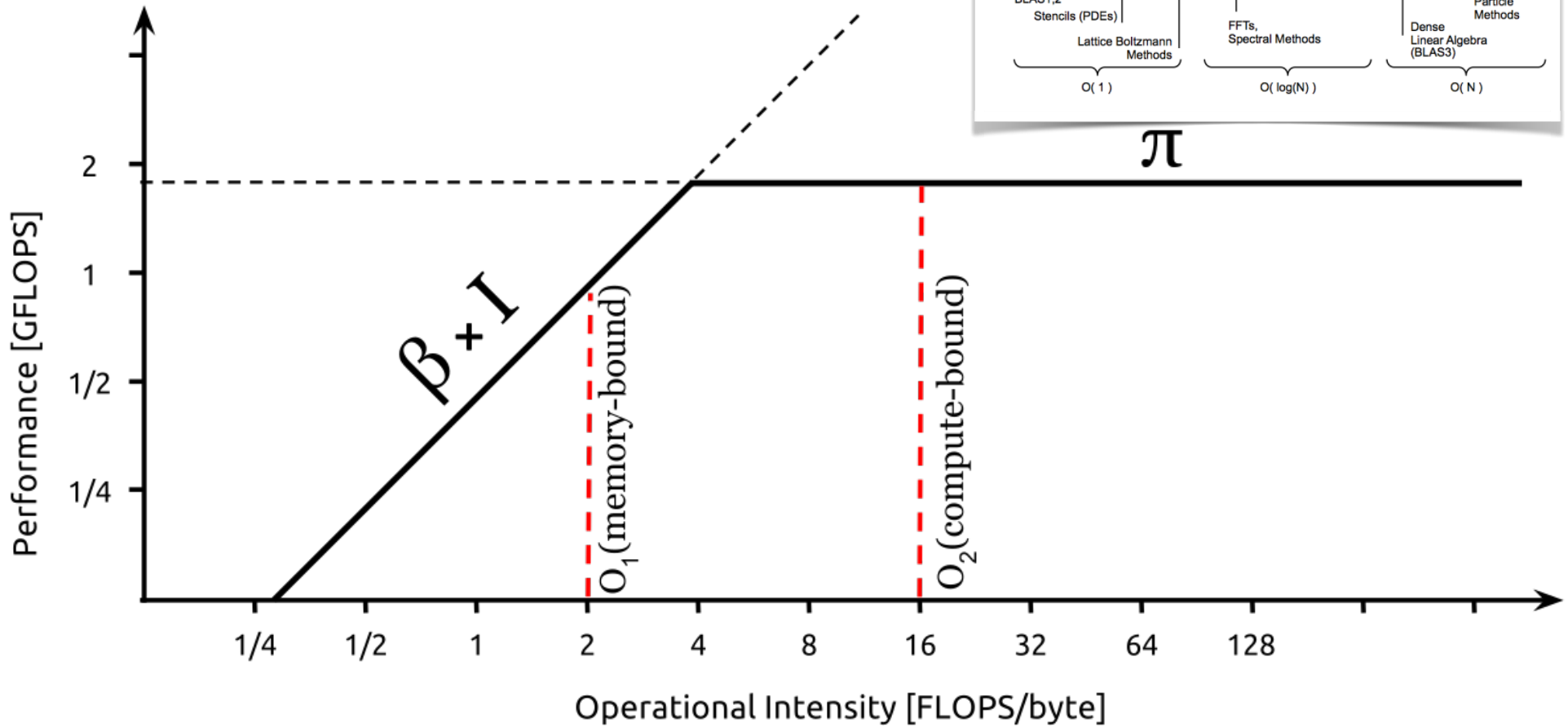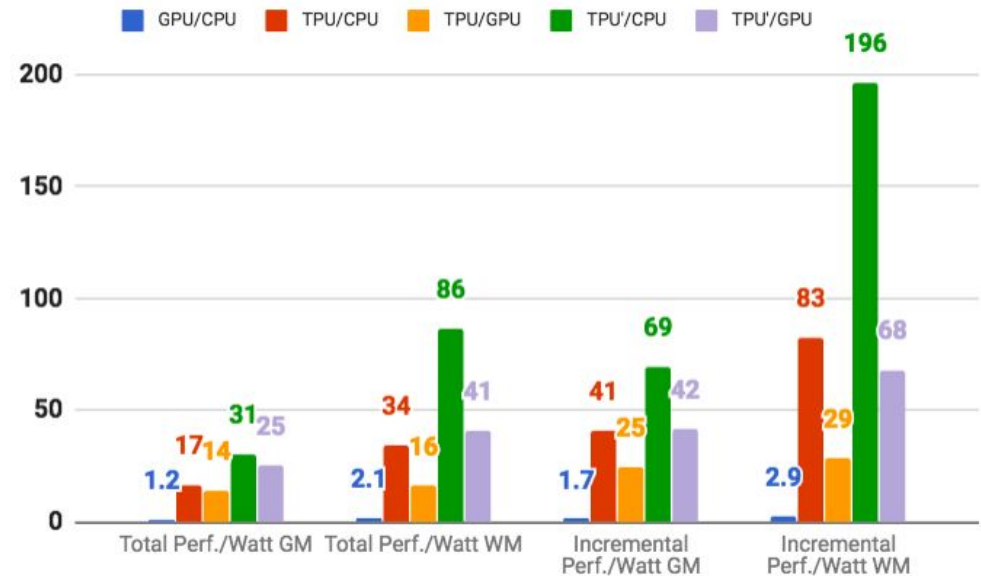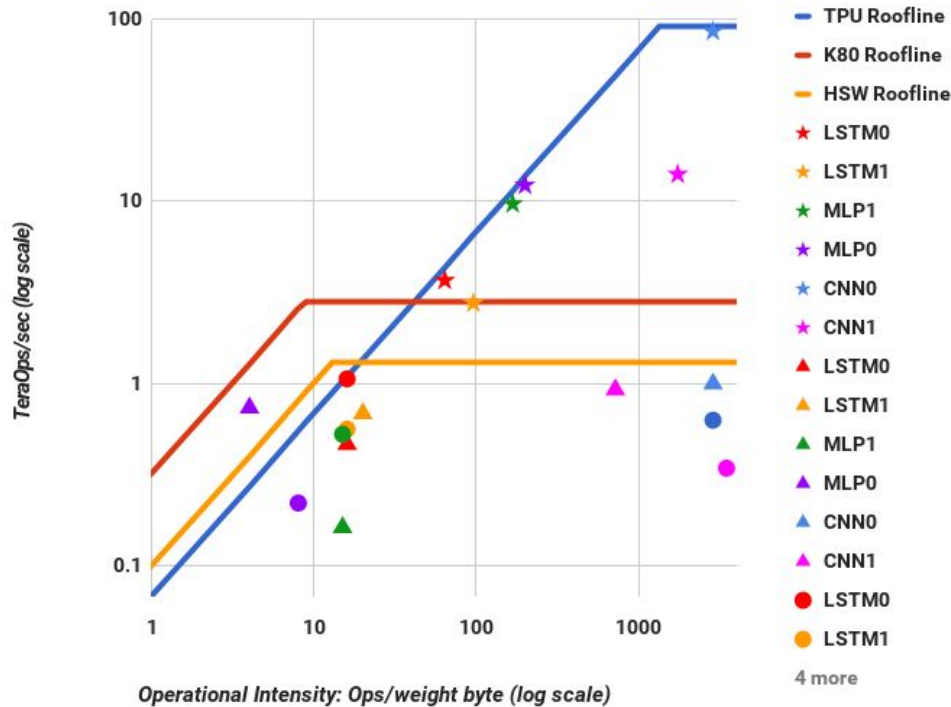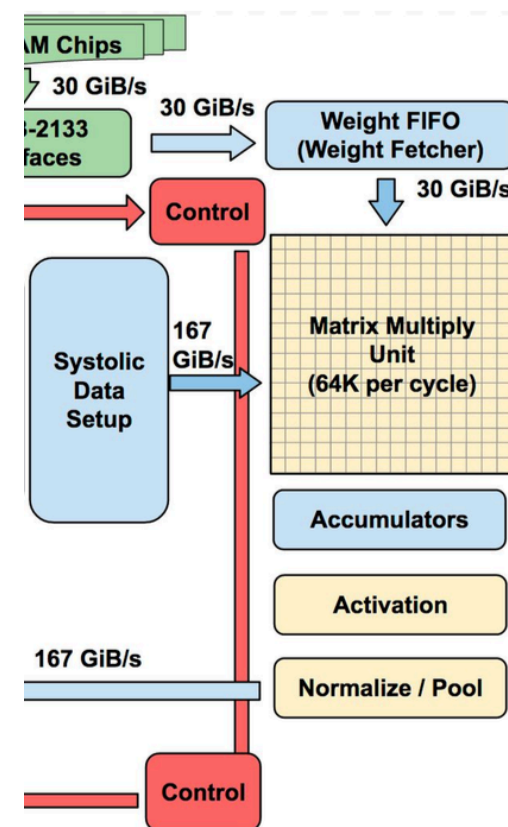
# The Roofline Model

# Performance



Stars are for the TPU, triangles are for the K80, and circles are for Haswell. All TPU stars are at or above the other 2 rooflines.

# App breakdown by Performance Counters

| Application | MLP0 | MLP1 | LSTM0 | LSTM1 | CNN0 | CNN1 | Mean | Row |
|---|---|---|---|---|---|---|---|---|
| Array active cycles | 12.7% | 10.6% | 8.2% | 10.5% | 78.2% | 46.2% | 28% | 1 |
| Useful MACs in 64K matrix (% peak) | 12.5% | 9.4% | 8.2% | 6.3% | 78.2% | 22.5% | 23% | 2 |
| Unused MACs | 0.3% | 1.2% | 0.0% | 4.2% | 0.0% | 23.7% | 5% | 3 |
| Weight stall cycles | 53.9% | 44.2% | 58.1% | 62.1% | 0.0% | 28.1% | 43% | 4 |
| Weight shift cycles | 15.9% | 13.4% | 15.8% | 17.1% | 0.0% | 7.0% | 12% | 5 |
| Non-matrix cycles | 17.5% | 31.9% | 17.9% | 10.3% | 21.8% | 18.7% | 20% | 6 |
| RAW stalls | 3.3% | 8.4% | 14.6% | 10.6% | 3.5% | 22.8% | 11% | 7 |
| Input data stalls | 6.1% | 8.8% | 5.1% | 2.4% | 3.4% | 0.6% | 4% | 8 |
| TeraOps/sec (92 Peak) | 12.3 | 9.7 | 3.7 | 2.8 | 86.0 | 14.1 | 21.4 | 9 |

Stalls due 2 Memory

Low utilization

**Table 3.** Factors limiting TPU performance of the NN workload based on hardware performance counters. Rows 1, 4, 5, and 6 total 100% and are based on measurements of activity of the matrix unit. Rows 2 and 3 further break down the fraction of 64K weights in the matrix unit that hold useful weights on active cycles. Our counters cannot exactly explain the time when the matrix unit is idle in row 6; rows 7 and 8 show counters for two possible reasons, including RAW pipeline hazards and PCIe input stalls. Row 9 (TOPS) is based on measurements of production code while the other rows are based on performance-counter measurements, so they are not perfectly consistent. Host server overhead is excluded here. The MLPs and LSTMs are memory-bandwidth limited but CNNs are not. CNN1 results are explained in the text.

# Latency Results (99%ile)
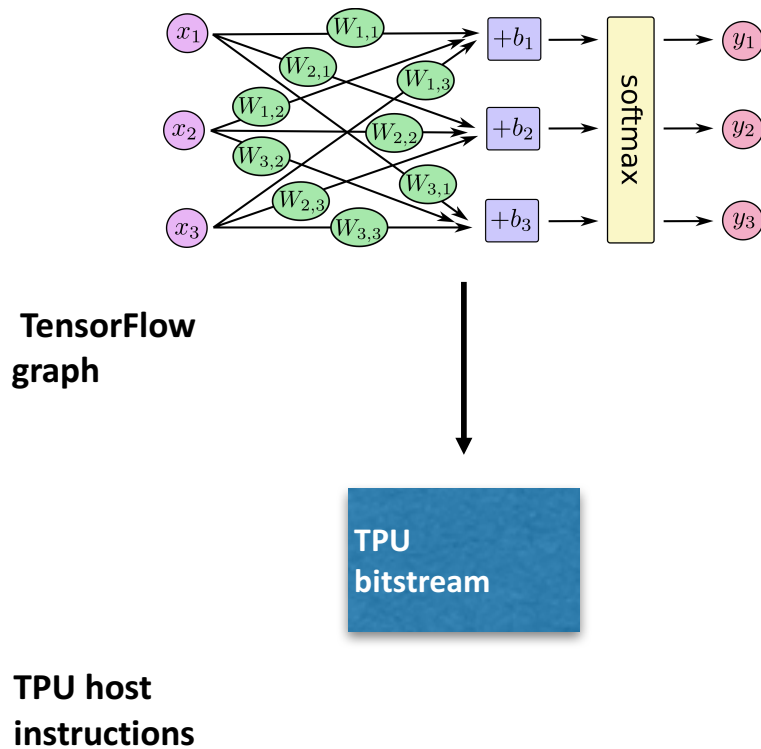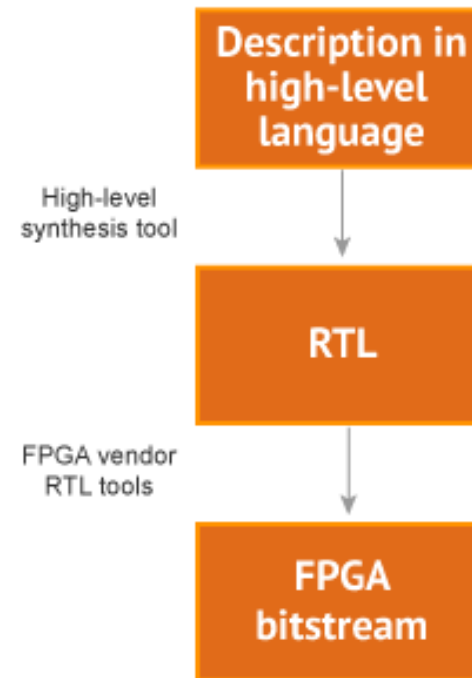
| Type | Batch | 99th% Response | Inf/s (IPS) | % Max IPS |
|------|-------|----------------|-------------|-----------|
| CPU | 16 | 7.2 ms | 5,482 | 42% |
| CPU | 64 | 21.3 ms | 13,194 | 100% |
| GPU | 16 | 6.7 ms | 13,461 | 37% |
| GPU | 64 | 8.3 ms | 36,465 | 100% |
| TPU | 200 | 7.0 ms | 225,000 | 80% |
| TPU | 250 | 10.0 ms | 280,000 | 100% |

**Table 4.** 99-th% response time and per die throughput (IPS) for MLP0 as batch size varies for MLP0. The longest allowable latency is 7 ms. For the GPU and TPU, the maximum MLP0 throughput is limited by the host server overhead. Larger batch sizes increase throughput, but as the text explains, their longer response times exceed the limit, so CPUs and GPUs must use less-efficient, smaller batch sizes (16 vs. 200).

# Programming the TPU

## Programming FPGAs

# NVIDIA's Rebuttal to the TPU

|  | K80 2012 | TPU 2015 | P40 2016 |
|---|---|---|---|
| Inferences/Sec <10ms latency | $\frac{1}{13}$X | 1X | 2X |
| Training TOPS | 6 FP32 | NA | 12 FP32 |
| Inference TOPS | 6 FP32 | 90 INT8 | 48 INT8 |
| On-chip Memory | 16 MB | 24 MB | 11 MB |
| Power | 300W | 75W | 250W |
| Bandwidth | 320 GB/S | 34 GB/S | 350 GB/S |

https://blogs.nvidia.com/blog/2017/04/10/ai-drives-rise-accelerated-computing-datacenter/
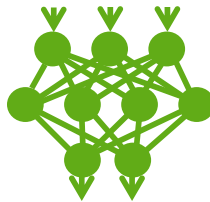
# Interesting quote

"CNNs constitute only about 5% of the representative NN workload for Google. More attention should be paid to MLPs and LSTMs. Repeating history, **it's similar to when many architects concentrated on floating-point performance when most mainstream workloads turned out to be dominated by integer operations.**"

# Neural acceleration

**Find** an approximate program component

**Compile** the program and train a neural network

Program

# Neural acceleration

[Esmaeilzadeh et al.]

Program

**Find** an approximate program component
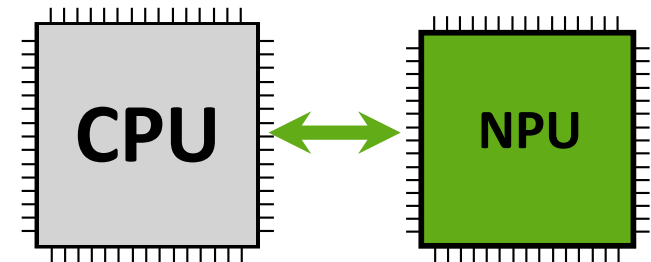
**Compile** the program and train a neural network
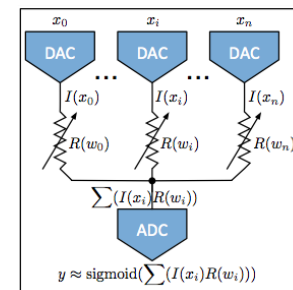
**Execute** on a fast Neural Processing Unit (NPU)

**CPU** ⟷ **NPU**

# Summary of NPU results

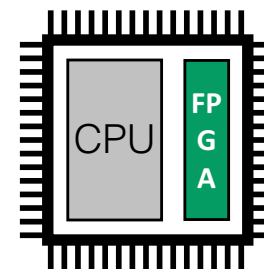| application | domain | error metric |
|---|---|---|
| blackscholes | option pricing | MSE |
| fft | DSP | MSE |
| inversek2j | robotics | MSE |
| jmeint | 3D-modeling | miss rate |
| jpeg | compression | image diff |
| kmeans | ML | image diff |
| sobel | vision | image diff |



0.9x - 24x (3.7x mean) speedup

1.5x - 51x (6.8x mean) energy red.



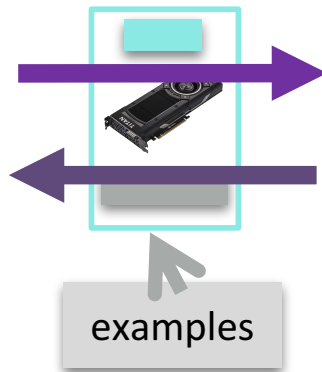0.8x - 11.1x (3x mean) speedup

1.1x - 21x (3x mean) energy red.


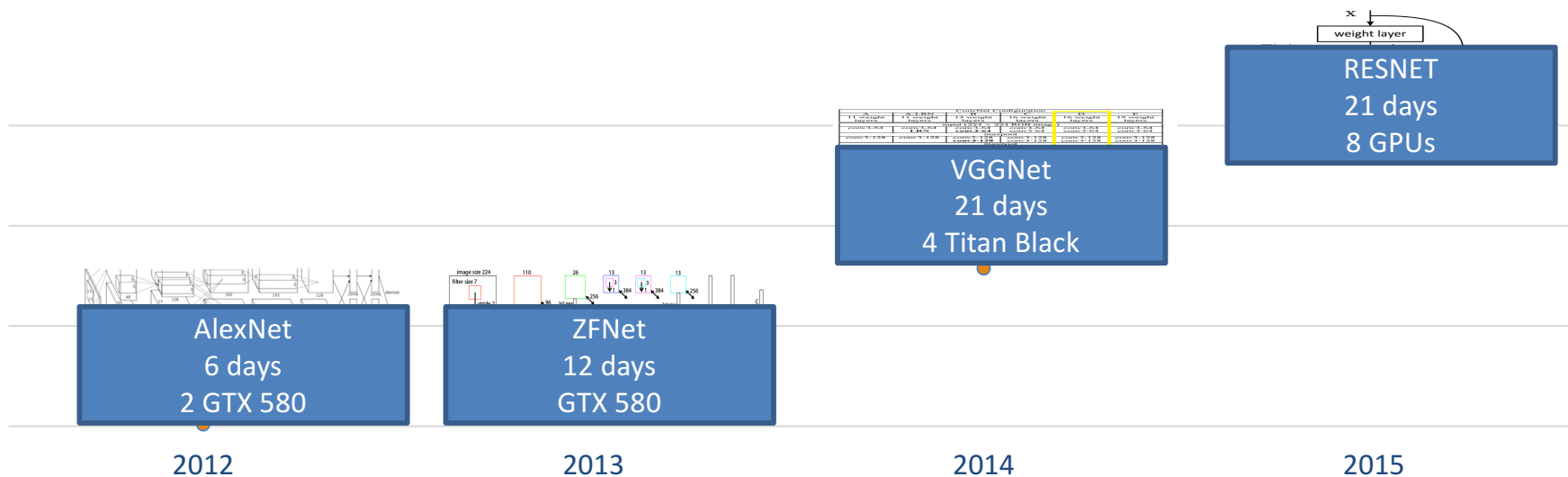
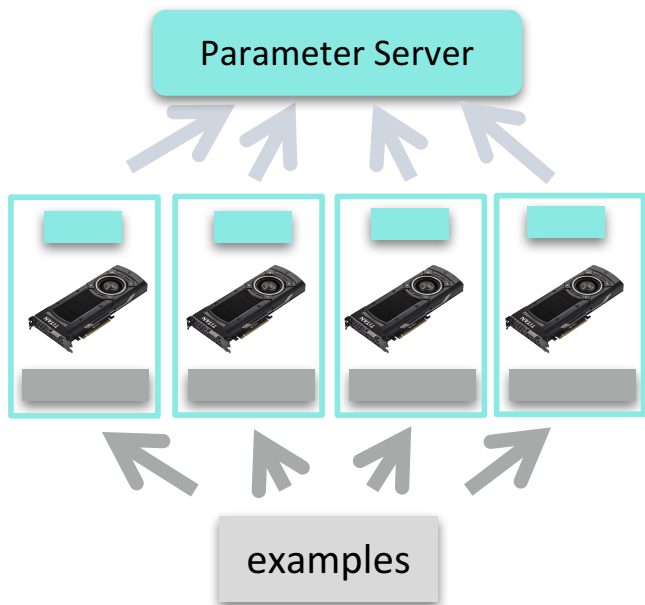1.3x - 38x (3.8x mean) speedup

0.9x - 28x (2.8x mean) energy red.

# DNN Training Time



examples

## Batches, foreward and backward propagation

1. A batch of samples are loaded into GPU.

2. The batch of samples does forward propagation and prediction error is derived.

3. The batch of samples undergoes backward propagation.

4. The model is updated and used for subsequent traning.



RESNET
21 days
8 GPUs

VGGNet
21 days
4 Titan Black

AlexNet
6 days
2 GTX 580

ZFNet
12 days
GTX 580

| 2012 | 2013 | 2014 | 2015 |

**Parameter Server**

examples

## Data Parallelism

1. Each device sees different parts of the data set. Devices work independently of each other.

2. Local gradient is calculated per device, and are communicated with parameter server during each batch.

# Distributed DNN Training (MXNET, TENSORFLOW…)

**Parameter Server**

1. Combine updates from all workers
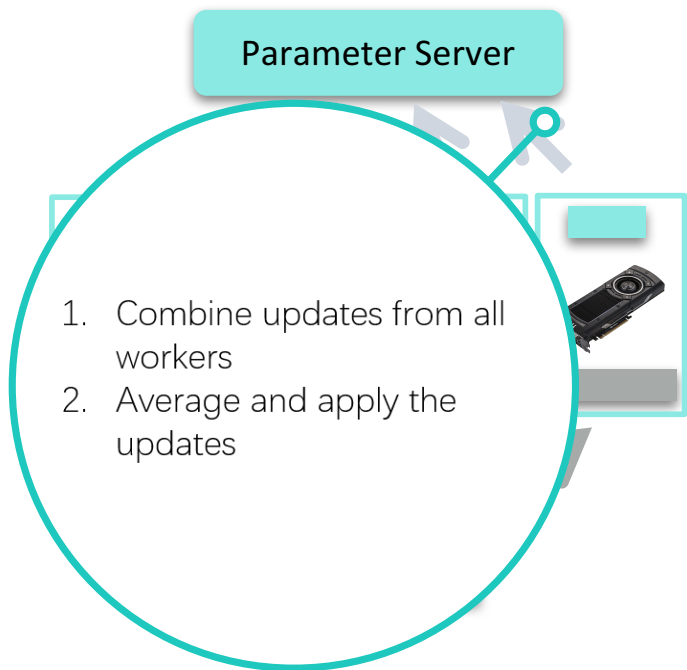2. Average and apply the updates

Data Parallelism

1. Each device sees different parts of the data set. Devices work independently of each other.

2. Local gradient is calculated per device, and are communicated with parameter server during each batch.

3. The parameter aggregates all updates and apply changes to the next model.

Where is the bottleneck? How do we improve it?