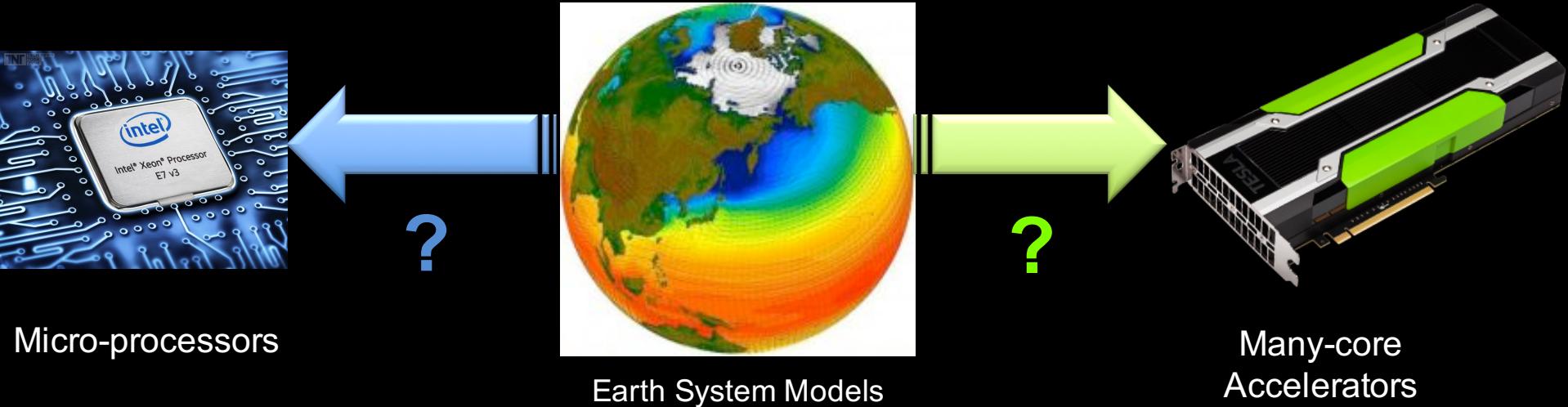
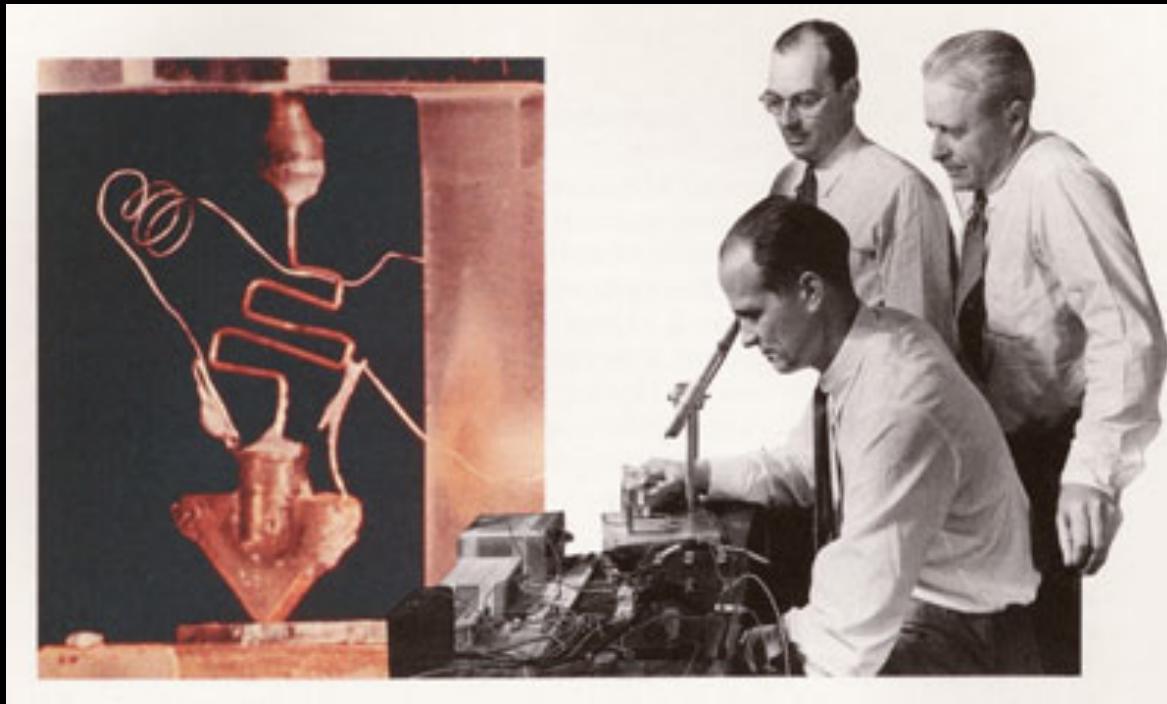


Parallel Computing Trends

Dr. Richard D. Loft
Director, Technology Development Division
Computational and Information Systems Laboratory
National Center for Atmospheric Research
loft@ucar.edu

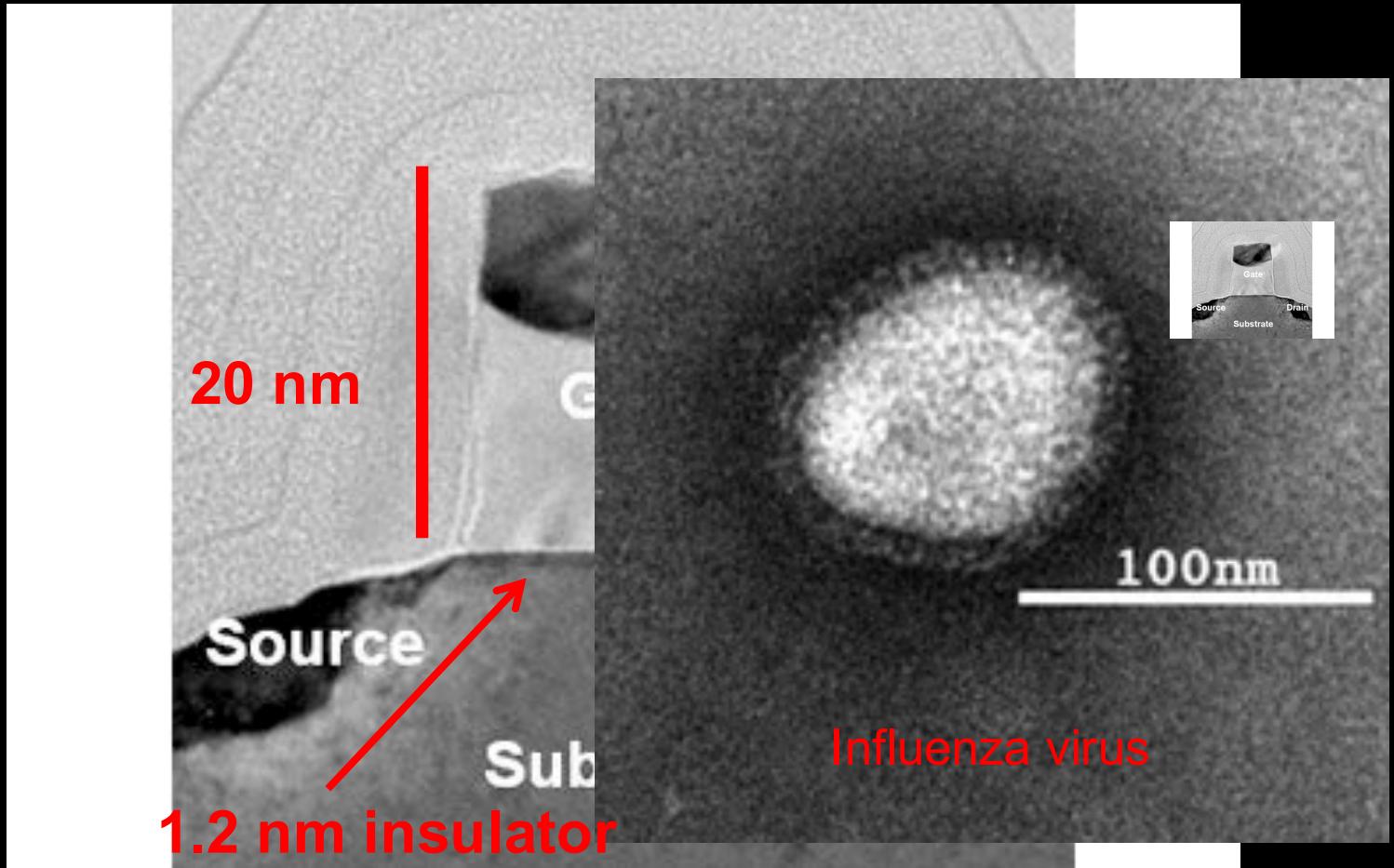


Transistor in 1947: ~ 5 cm across



William Shockley (1910-1989), John Bardeen (1908-1991),
Walter Brattain (1902-1987)

The modern transistor: ~200 atoms across



6/18/16

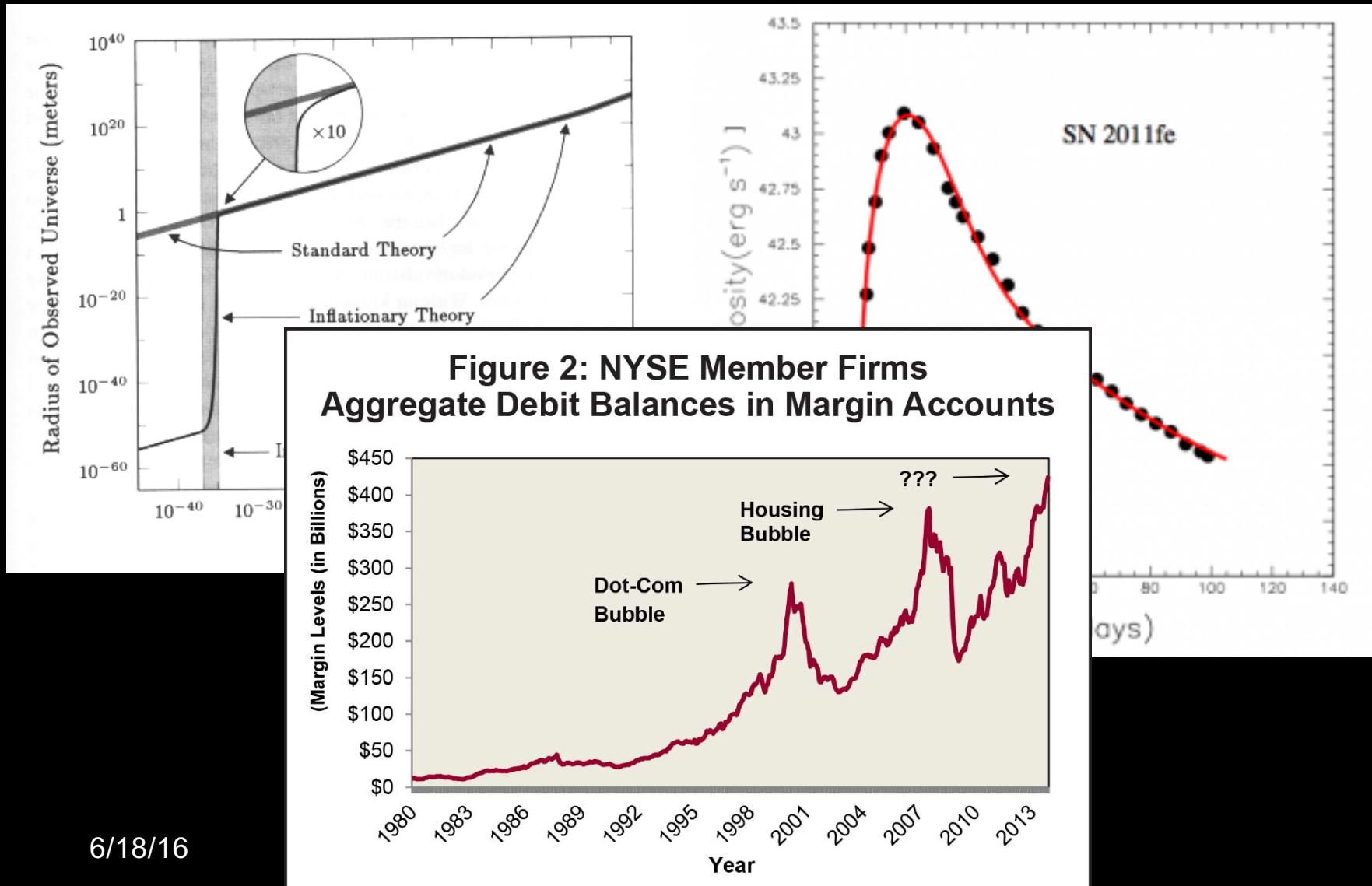
So small that quantum tunneling is a problem!

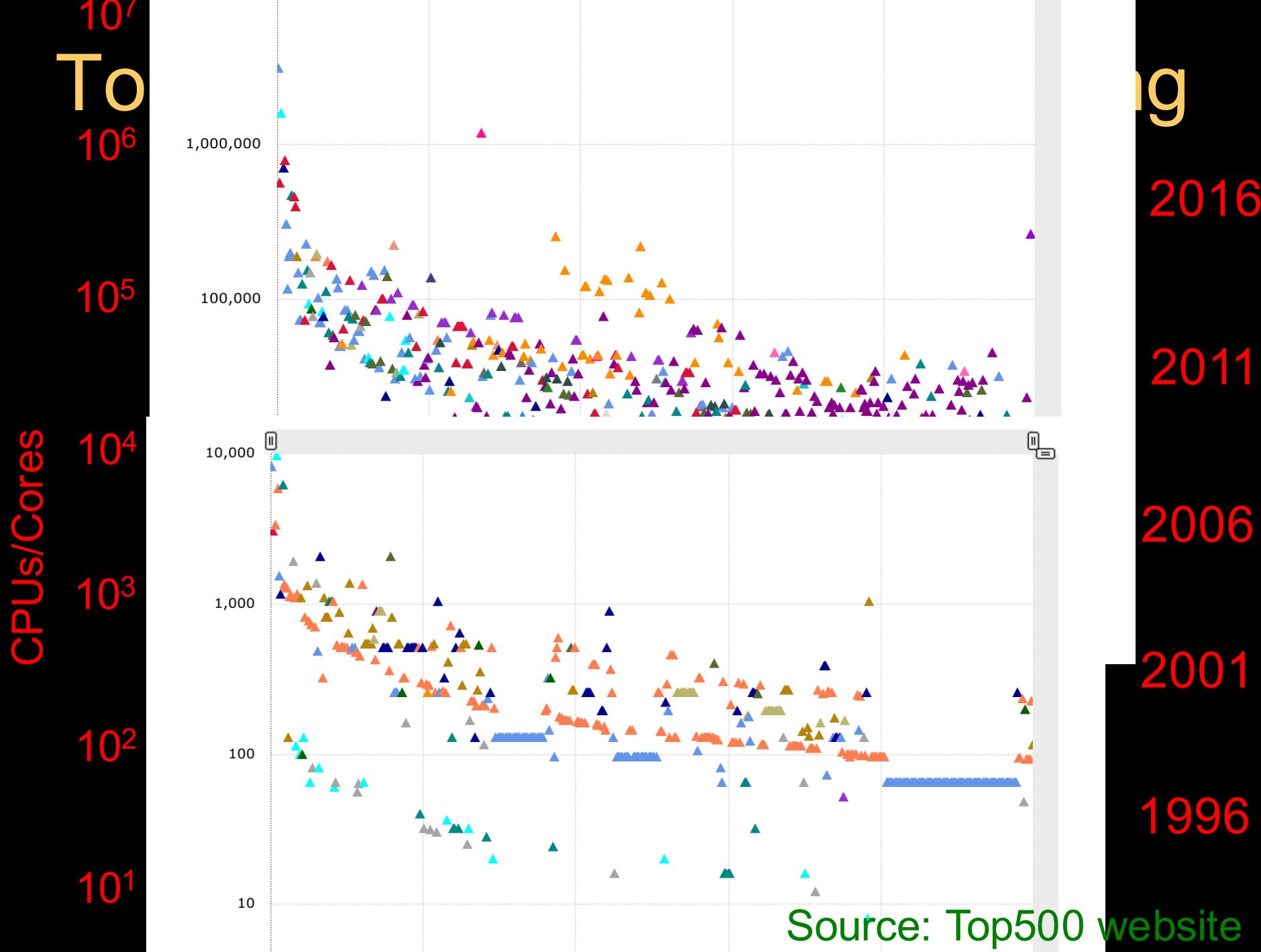
The exponential increase in the speed of computers. Looks relentless.



... This laptop would have been on the top 500 list in 2000

No exponential lasts forever!





Path Forward is through Parallelism

The largest supercomputers now have $O(10^5 - 10^6)$ compute cores . . .



Tianhe 2: 3.12 Million cores

Mira: 786,432 cores

Blue Waters: 362,240 cores

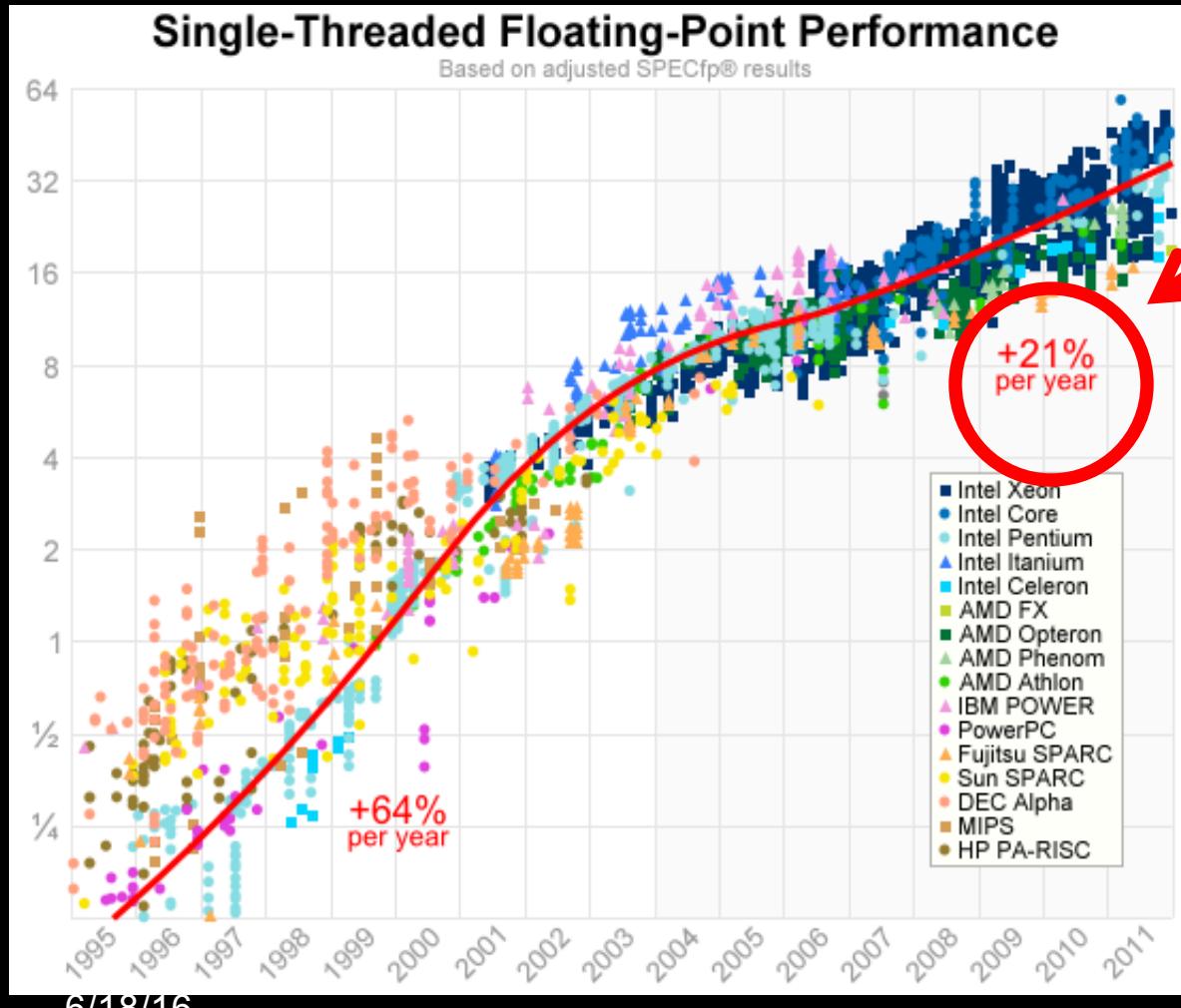
6/18/16

Processor Trends

What happened?

- Around 2000 chip designers realized:
 - Power wall
 - Clock speed couldn't continue to rise due to
 - power density
 - Memory wall
 - Memory system speed couldn't keep up with processor speed
 - ILP wall
 - Benefit of adding processor "smarts" reached point of diminishing returns
- Solution:
 - Break up monolithic single processor into multiple computing regions: cores

Clock speeds stopped improving and single thread improvements slowed



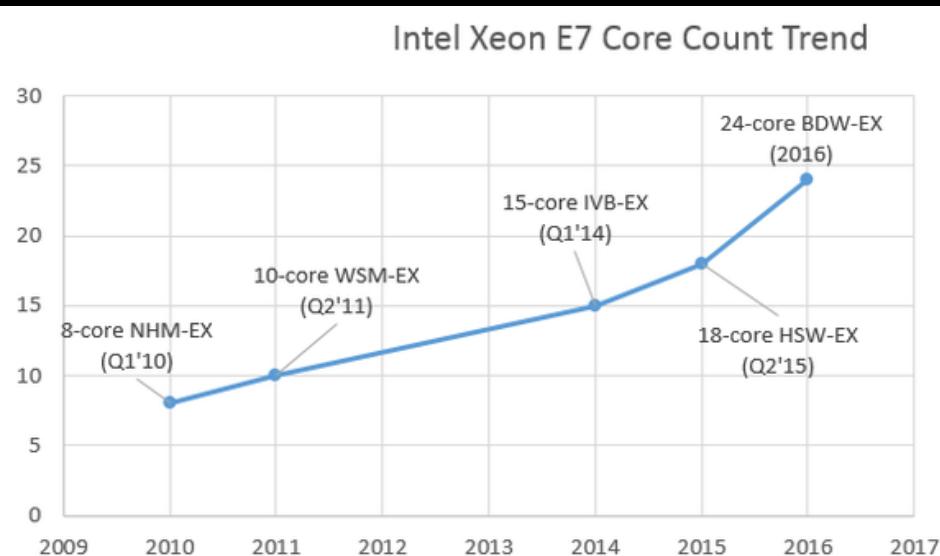
This translates to 5% Model resolution improvement per year

Bias corrected SpecFP results Preshing.com

Original Source-
Sutter, H. "The
Concurrency Revolution"
Dr. Dobb's Journal, 30(3),
March 2005.

Many-core: evolutionary approach

- Hold clock speed ~constant
- Hold core complexity ~constant
- **Gradually** add more cores on the die as transistors shrink
- **Gradually** add SIMD features (vector units)

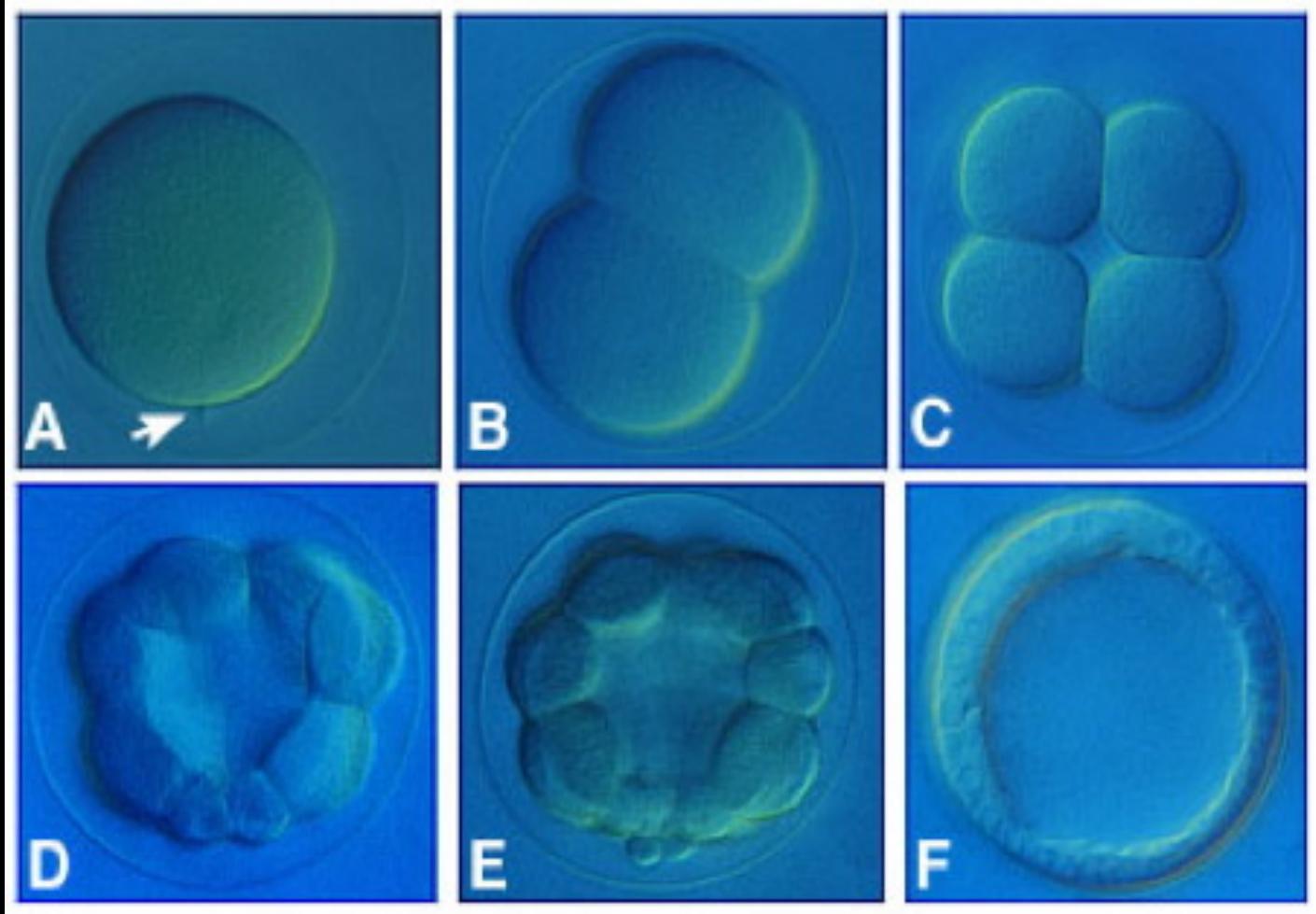


Source: Ashraf Eassa, The Motley Fool

SIMD =
Single Instruction Multiple Data



The Processor Zygote: Speed-up through increasing core parallelism on processor



6/18/16

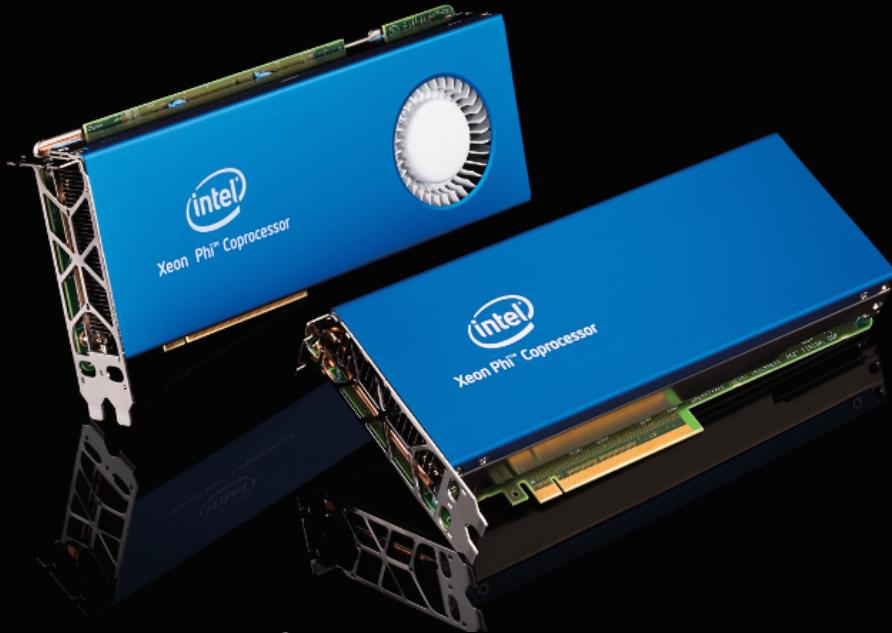
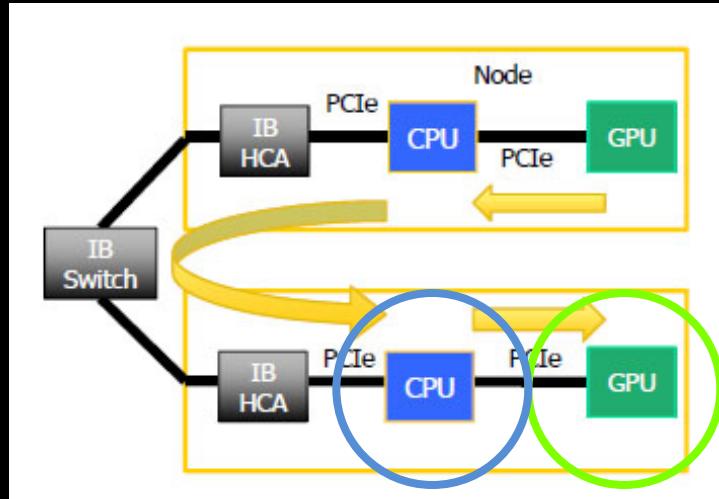
What are processors developing into?

12

Many-core: revolution

- Solution 2: single core -> many-core
 - Drop clock speed (i.e. by 2-3x)
 - Drop core complexity
 - Stamp out a lot more simplified cores on chip
 - Add even more cores as transistor shrinks
 - Add much more “SIMD-ness” to architecture

Added complexity: coprocessors and heterogeneous computing



- Traditionally, problems with this model:
- Two memory spaces issue (host,device)
 - CPU mediated memory access
 - Limited memory bandwidth
 - Compiler language issues (e.g. CUDA)

Slowly going away over time¹⁴

Intel Xeon Phi Knights Landing: 72 cores; 512 vectors; 3 TFLOPS



Knights Landing

Holistic Approach to Real Application Breakthroughs

Platform Memory



Up to 384 GB DDR4 (6 ch)

Over 60 Cores

Integrated Intel® Omni-Path

Processor Package

I/O

Up to 36 PCIe 3.0 lanes

Compute

- Intel® Xeon® Processor Binary-Compatible
- 3+ TFLOPS¹, 3X ST² (single-thread) perf. vs KNC
- 2D Mesh Architecture
- Out-of-Order Cores

On-Package Memory

- Over 5x STREAM vs. DDR4³
- Up to 16 GB at launch

Omni-Path
(optional)

▪ 1st Intel processor to integrate

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors.

Performance tests are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products.



Many-core architecture example: NVIDIA Kepler GPU (K20X) circa 2012



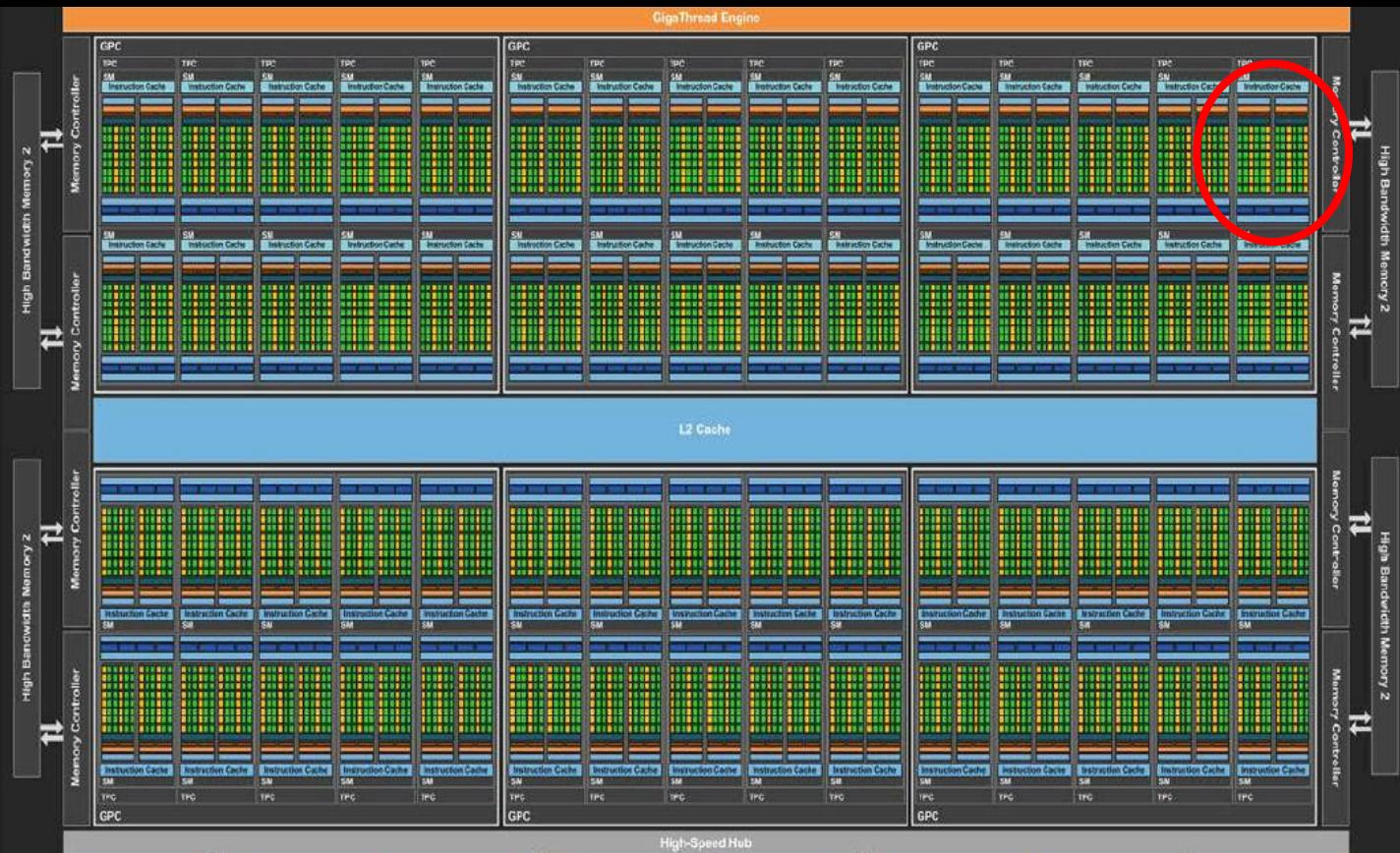
SM = Streaming
Multiprocessor

192 SIMD
CUDA
processors
per SM
X

15 SMX pipes
= 2688 cores

K20X memory
bandwidth:
= 250 GB/s

Many-core architecture example: NVIDIA Pascal GPU circa 2016

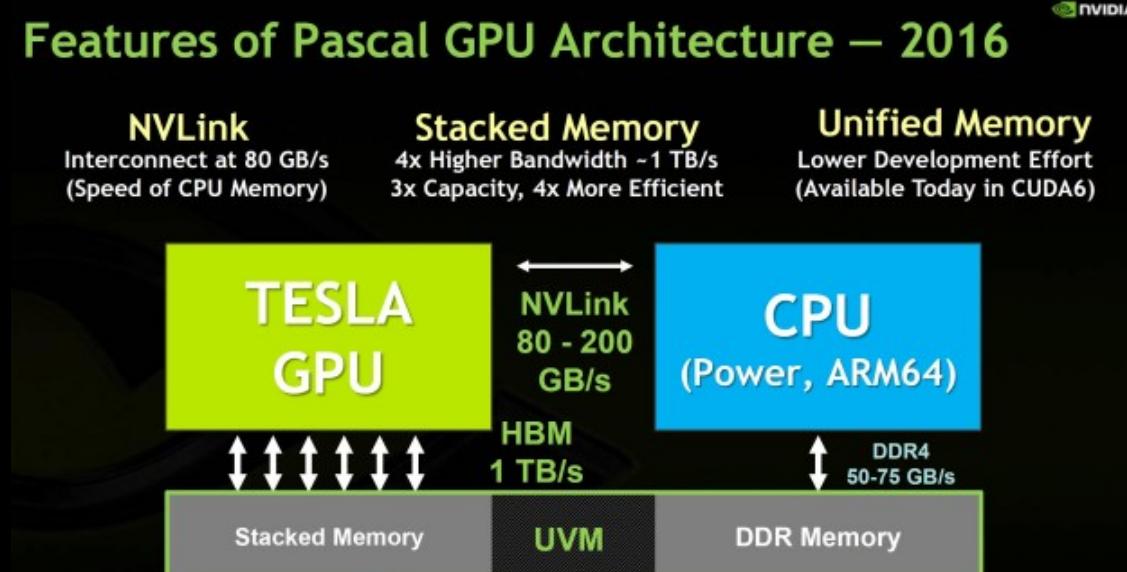


64 SIMD
CUDA cores
per SM
X
56 SMs pipes
= 3584 cores

stacked memory
bandwidth:
= 1000 GB/s

8-GPU DGX-1 architecture with 8 NVIDIA Pascal GPUs and NVLink

- $3584 \times 8 = 28,672$, 32-bit CUDA cores
- 4 dp teraflops $\times 8 = 32$ dp teraflops/node
- One dp petaflops peak in 32 nodes



Pascal: double, single, *half* precision?

- Interestingly, pascal supports fp16 for neural net calculations (>170 TF in a box)
- Question: is fp16 useful for *anything* in weather and climate?
- In a larger context this is an attempt to address the **incredible inefficiency** of conventional computers relative to biological neural networks (brains!)

Computers vs Brains



>35,000 kg
1,900,000 watts



1.4 kg (3 lbs)
20 watts

Why are supercomputers so bad and being “smart”?

- Supercomputers are designed to be:
 - Very good at **mathematical calculations**.
 - Remember everything (data integrity)
 - Predictable (reproducibility)
- Comparatively, brains are:
 - Are terrible at math
 - Forget stuff
 - Are generally unpredictable

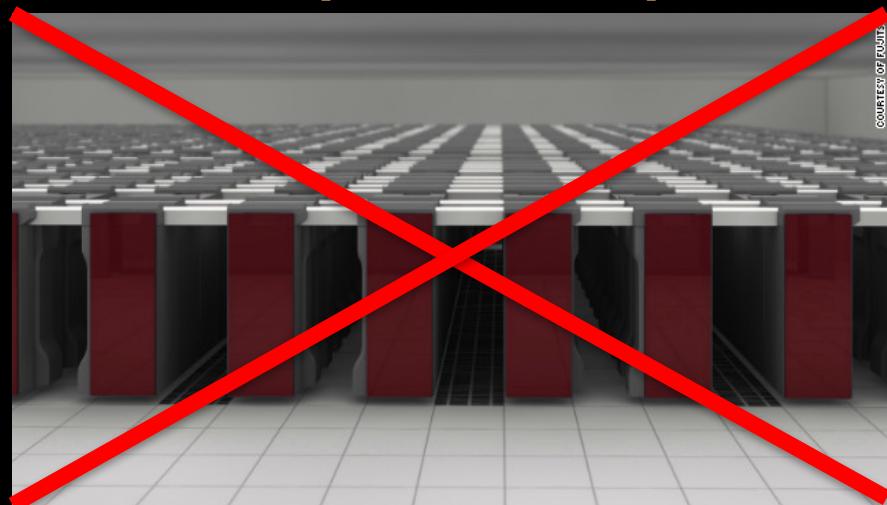
How “smart” is a supercomputer?

House Cat



OR

K computer in Japan



760 million neurons
10 trillion synapses

Simulated 1.7 billion neurons
with 10 trillion synapses
But 2400 times slower!

trillion = 1,000,000,000,000 (a million million)

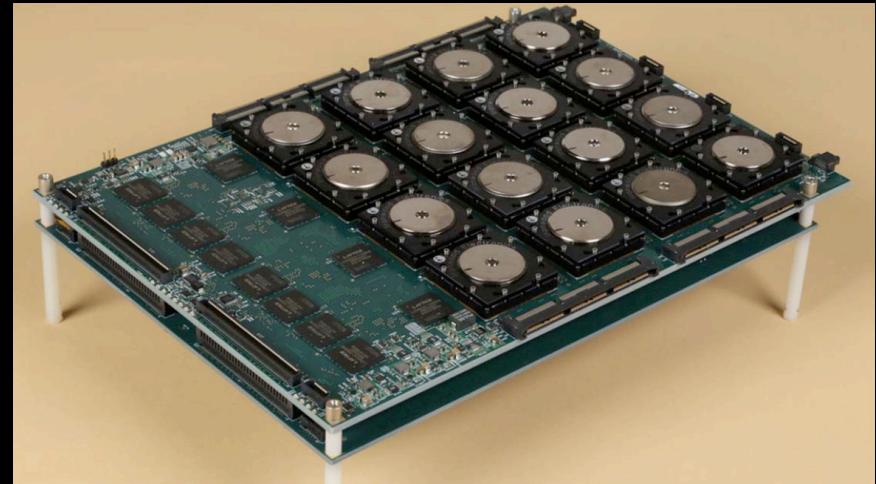
Neuromorphic computers: work like the brain!

Honey Bee



=

IBM TrueNorth Processor
(2015)



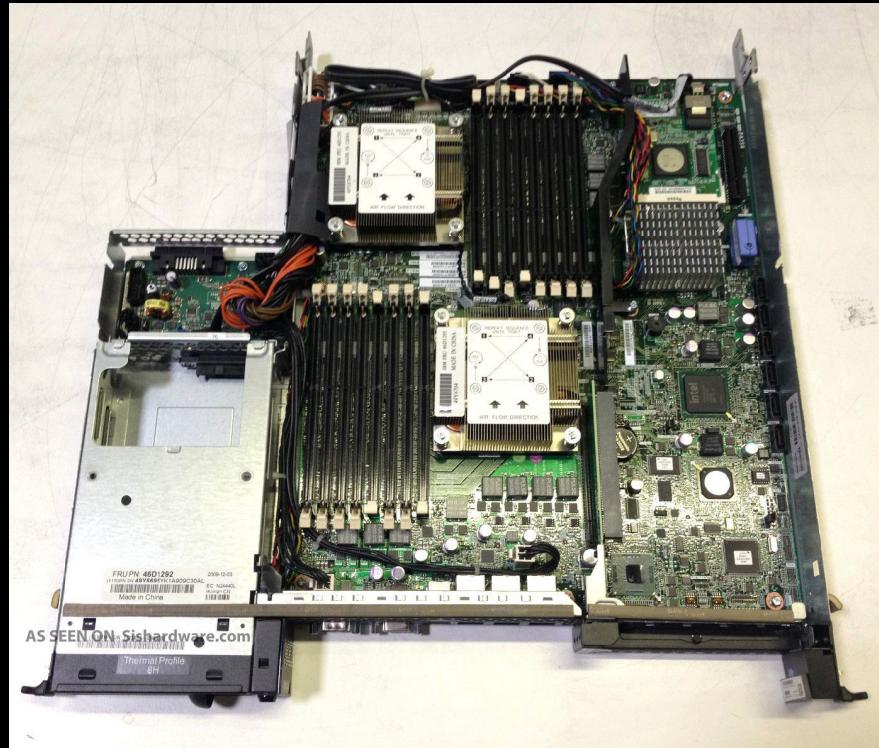
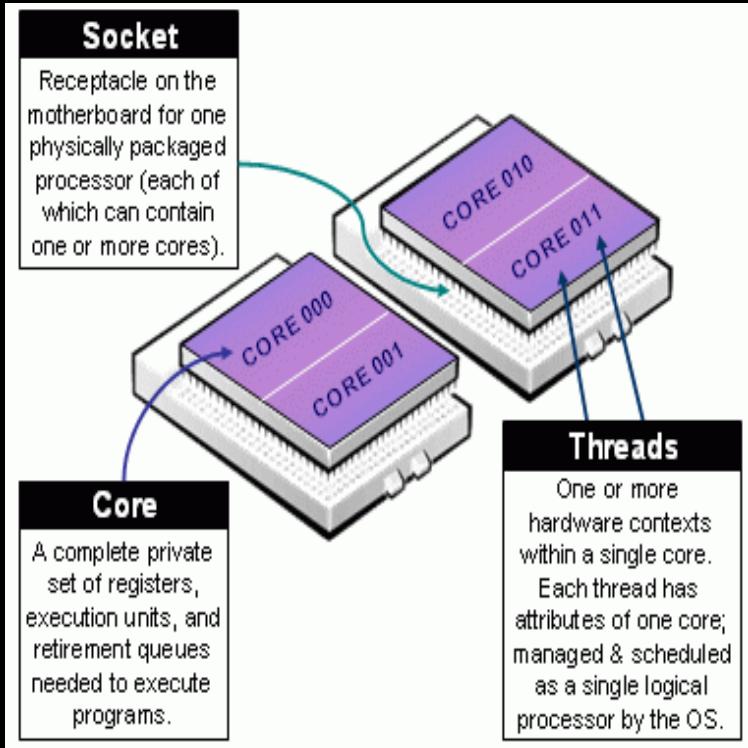
1 million neurons,
1 billion synapses
2 mW

Which is faster?

1 million neurons,
268 million synapses
@ 80 mW TrueNorth is
10,000 x more efficient
Than regular computers!

Memory and Bandwidth Trends

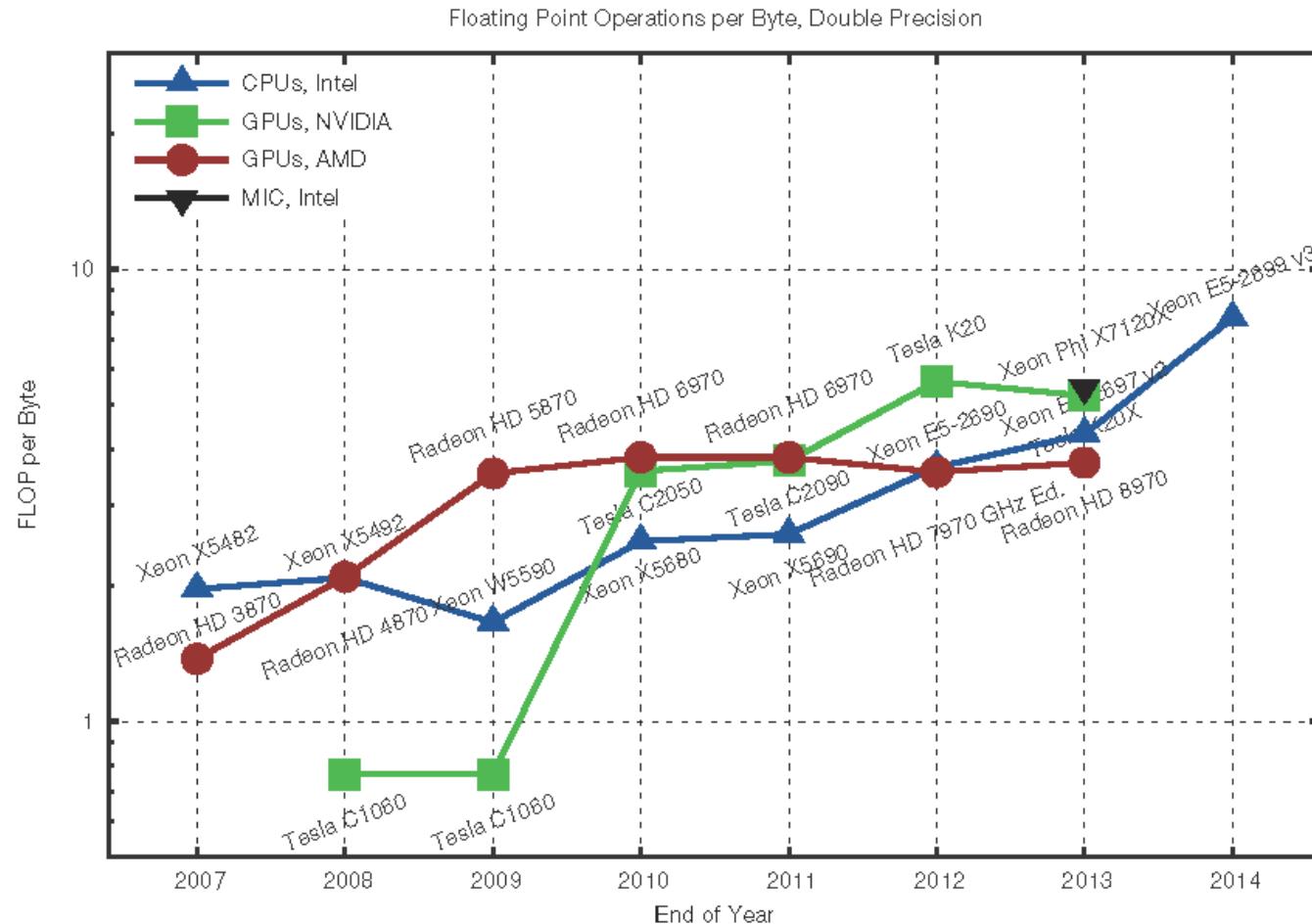
Traditional (Current) multi-core node architecture



Sockets, cores, threads

Danse hall: memory segregated from processors

Peak flops is driving top 10 system design

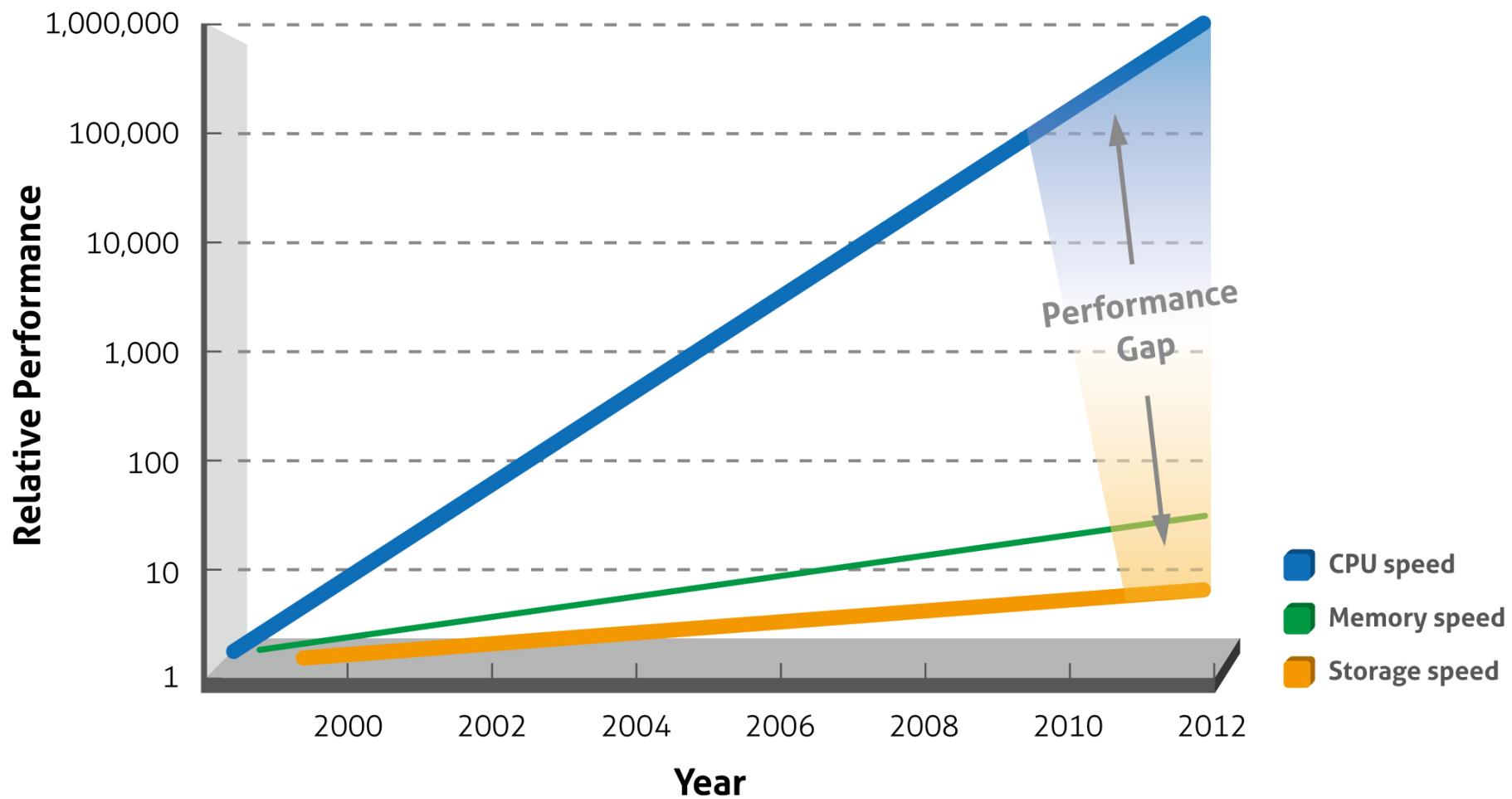


[c/o Karl Rupp]

Slide courtesy of Jed Brown, CU

the performance gap

Technology Trend



Power and Memory Bandwidth

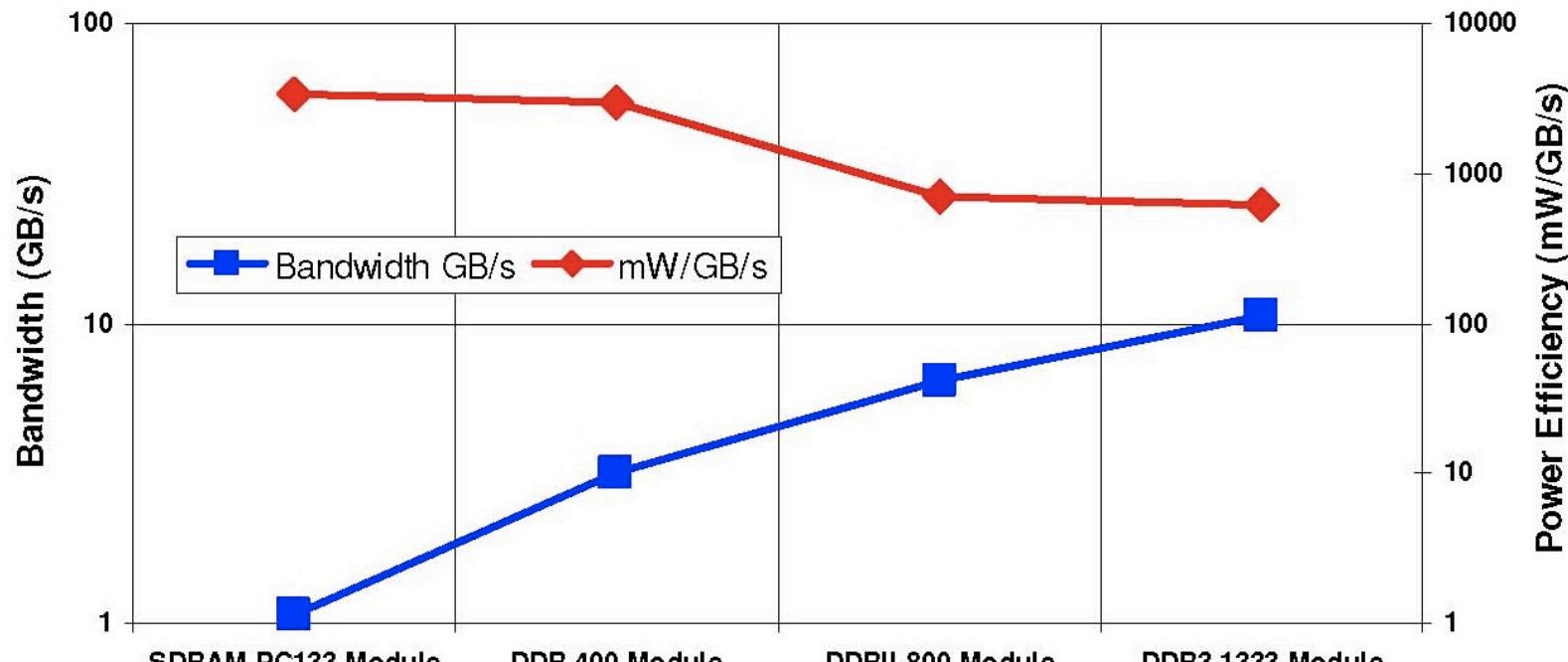


Figure 6.22: Commodity DRAM module power efficiency as a function of bandwidth.

O(1) watt/GB/s is relatively constant over time

Memory Bandwidth means power: consider an exaflops system

- Assume one, 8-byte word loaded from memory per flop
- To sustain **1% of peak**, i.e. $O(10)$ PFLOPS, we will need **80 PB/sec**
- Using the DDR figure of merit, we the memory we need is **~80 Megawatts**.
- This suggests a ***new technology*** is required to achieve a practical exascale capability.

Off chip bandwidth limitations

- Off chip bandwidth = pin count x signaling rate
 - Pin count might increase at 5% per year
- Upping electrical signaling rate increases power consumption.
 - Intel considers a 10x improvement a grand challenge



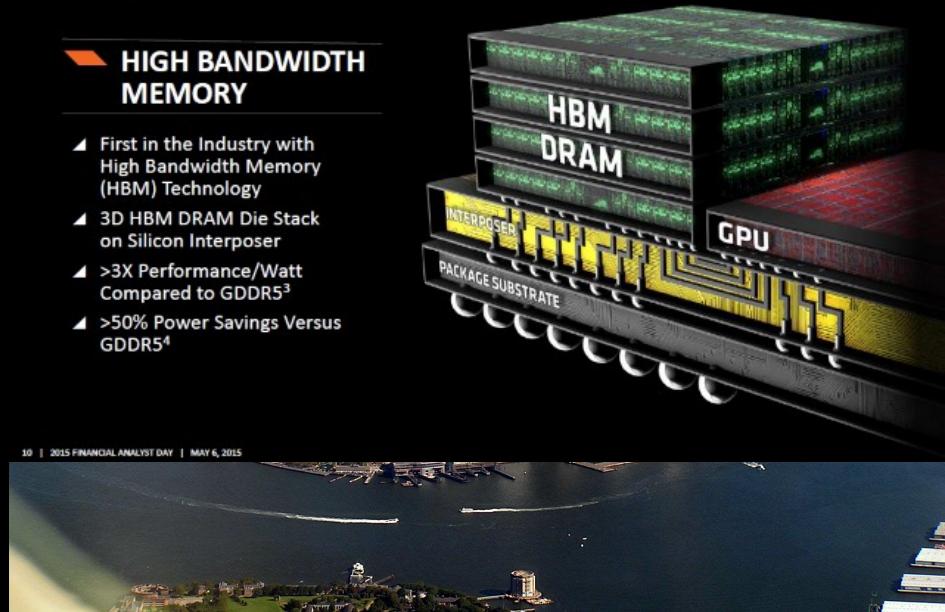
Going vertical: fuel-efficient memory bandwidth through 3-D Chips

- Advantages
 - Increased density
 - Shorter wires
 - lower capacitance
 - lower power
 - higher GB/sec/Watt
- Disadvantages
 - Tighter thermal constraints on memory
 - Limits to capacity?

GRAPHICS TECHNOLOGY LEADERSHIP

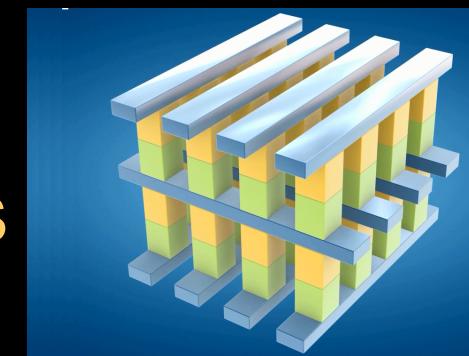
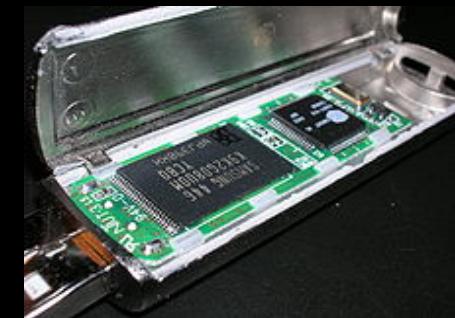
HIGH BANDWIDTH MEMORY

- ▲ First in the Industry with High Bandwidth Memory (HBM) Technology
- ▲ 3D HBM DRAM Die Stack on Silicon Interposer
- ▲ >3X Performance/Watt Compared to GDDR5³
- ▲ >50% Power Savings Versus GDDR5⁴



So-long disk? Non-volatile memory and NAND technologies

- Flash memory-based storage
 - Floating gate or NAND transistors
 - Limitations: density and endurance
 - Great for high IOPS/disk latency sensitive applications
- Nonvolatile storage technologies
 - 1000x faster than NAND
 - 1000x the endurance of NAND
 - 10x denser than DRAM



3D XPoint

Hybrid electro-optical systems

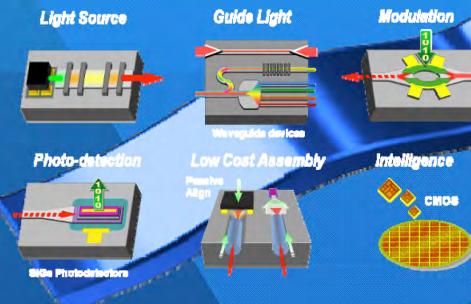
- Lots of bandwidth available in an optical fiber
- Parallelism: data can be carried on different wavelengths.
- Coupling optics to silicon is undergoing rapid development
 - Optically active silicon
 - Miniaturization of optical devices

Trends in micron-scale opto-electronics

Integration Vision

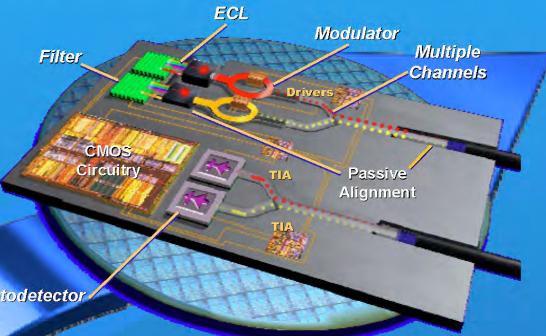
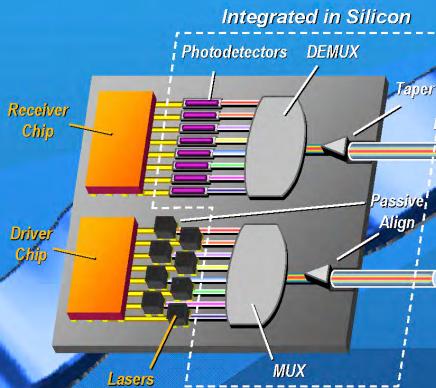
Time →

**First: Prove Silicon good optical material
*Many at 40Gb/s**



Next Integration: silicon devices into hybrid modules

Increasing silicon integration over time



**FUTURE
Monolithic?**

**Level of integration
Determined by
Application/cost**

all data out...
012.5 Gbps
per channel

1001110
01001110
01001110
101001110

Credit - Intel: Antuino, et al. Hot Chips 2010

Impacts on Atmospheric Models

Parallel programming models

- MPI – distributed memory programming
 - Co-Array Fortran
- OpenMP – directive-based shared memory
 - CPU moving towards GPU support
- GPU programming models:
 - CUDA (NVIDIA's SIMD language)
 - OpenCL (SIMD language)
 - In theory portable to CPUs and GPUs
 - OpenACC (directive-based)
 - In theory portable to CPUs and GPUs

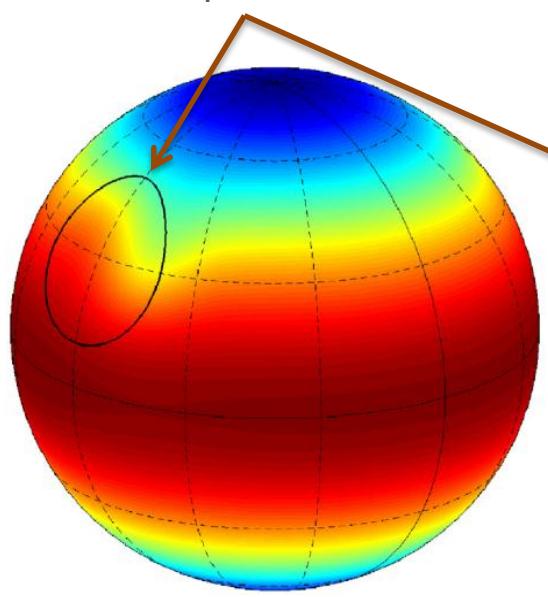
CPU-GPU intercomparison

- Relevant algorithmic characteristics
 - Available parallelism
 - “SIMD-ness”
 - Computational intensity (flops/memory accesses)
 - Workload imbalances/conditionals
 - Communication overhead
- Other metrics of performance
 - Hardware cost-performance
 - Energy-performance (watts/flops)
 - Software cost-performance (\$/SLOC optimized)

Which is Faster? Benchmark Problem

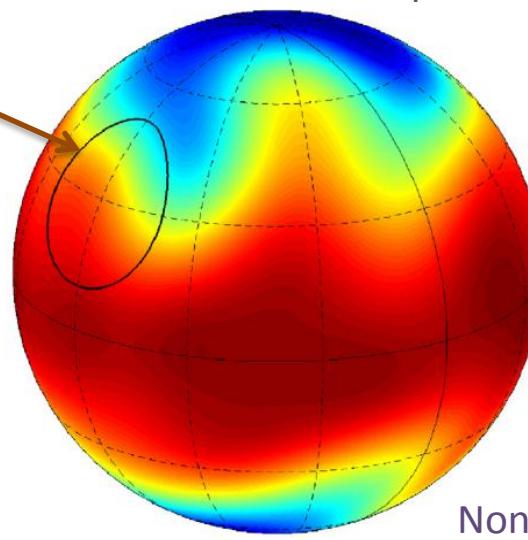
- Shallow Water Equations (SWE)
 - A set of non-linear partial differential equations (PDE)
 - Capture features of atmospheric flow around the Earth
- Radial basis function-generated finite difference (RBF-FD) methods

Cone-shaped mountain



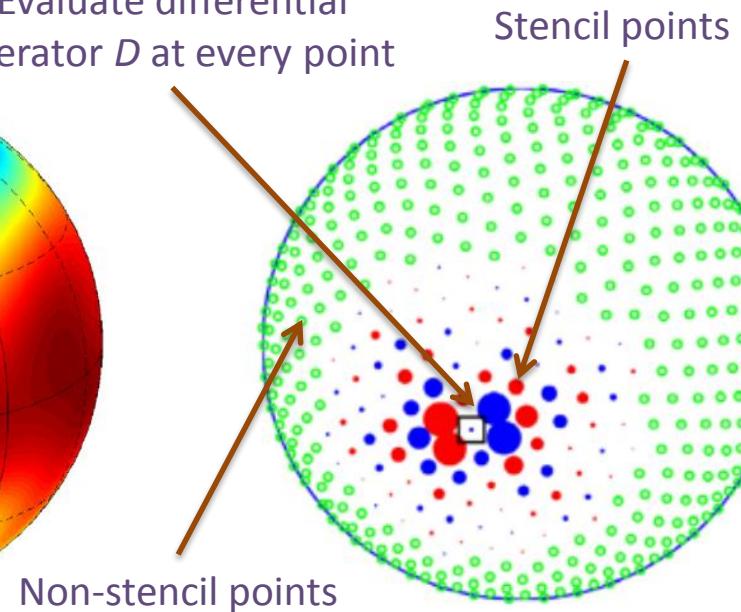
Day 1

Evaluate differential operator D at every point



Day 15

RBF-FD solution to SWE test case “Flow over an isolated mountain” using 655,532 points [1]



Non-stencil points

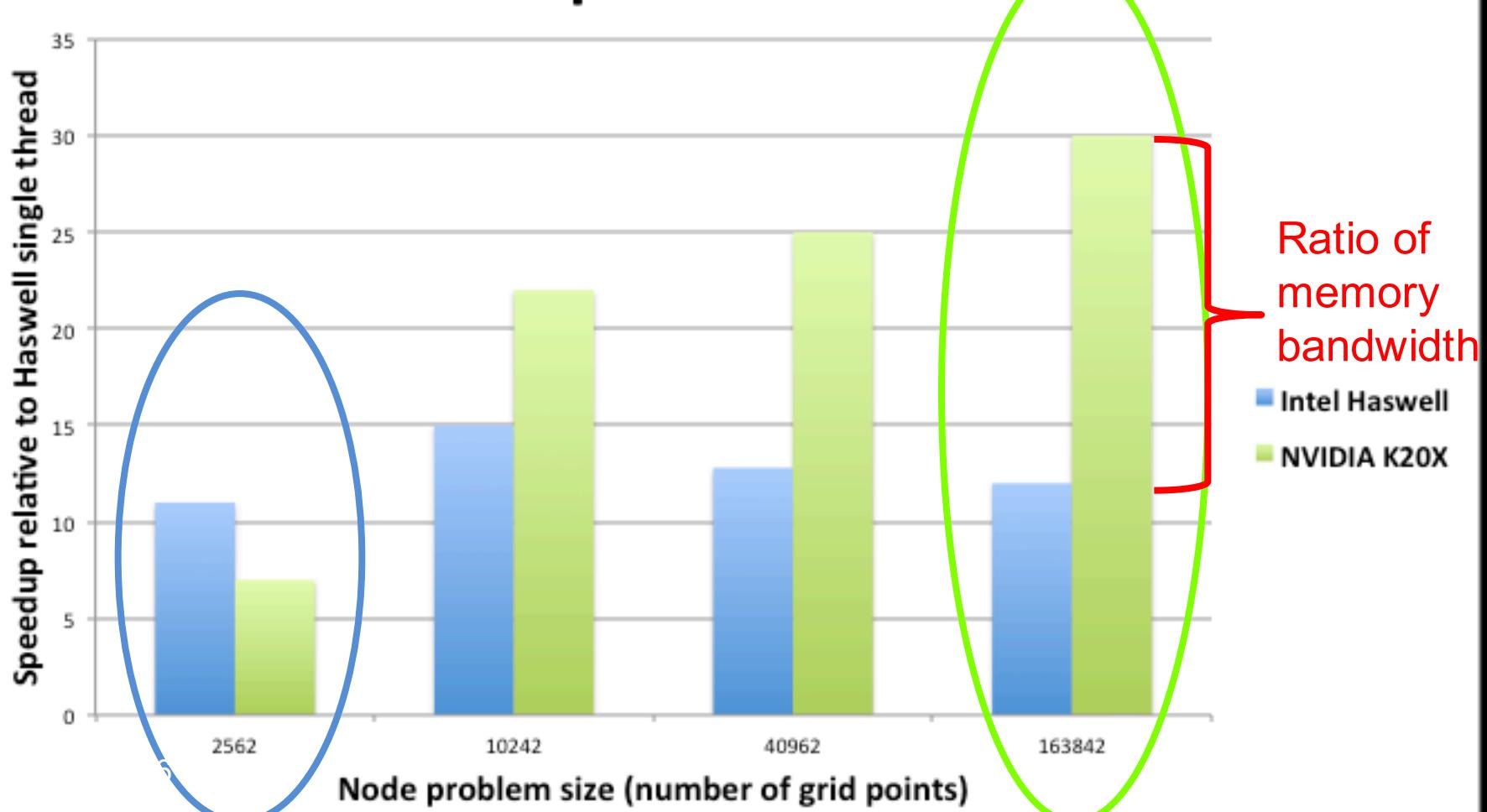
Stencil points

An example of 75-point stencil on a sphere [1]

Test case: shallow water equations, unstructured radial basis function grid, 31 point stencil

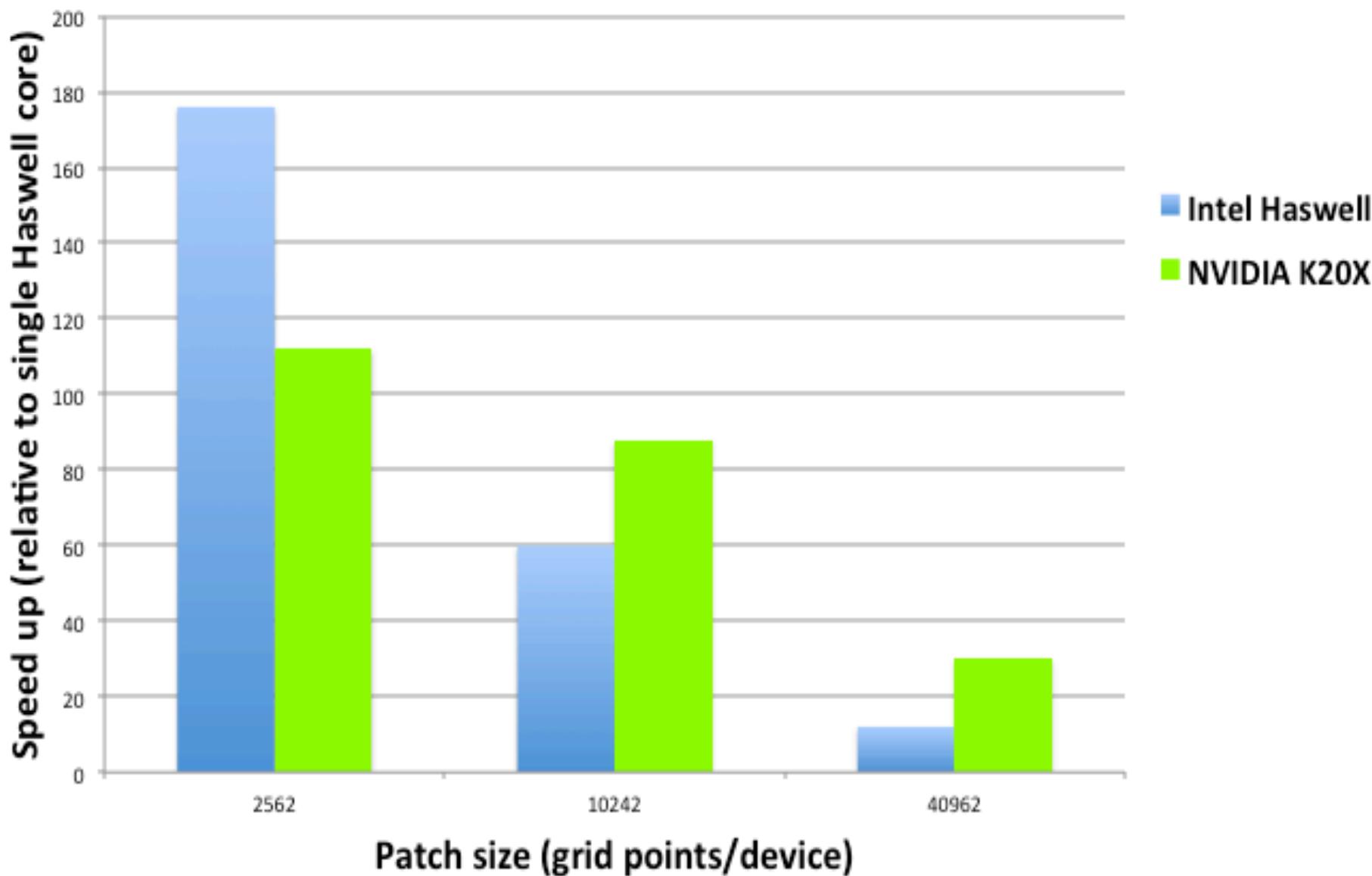
GPU's vs CPU's which is faster?

It depends...



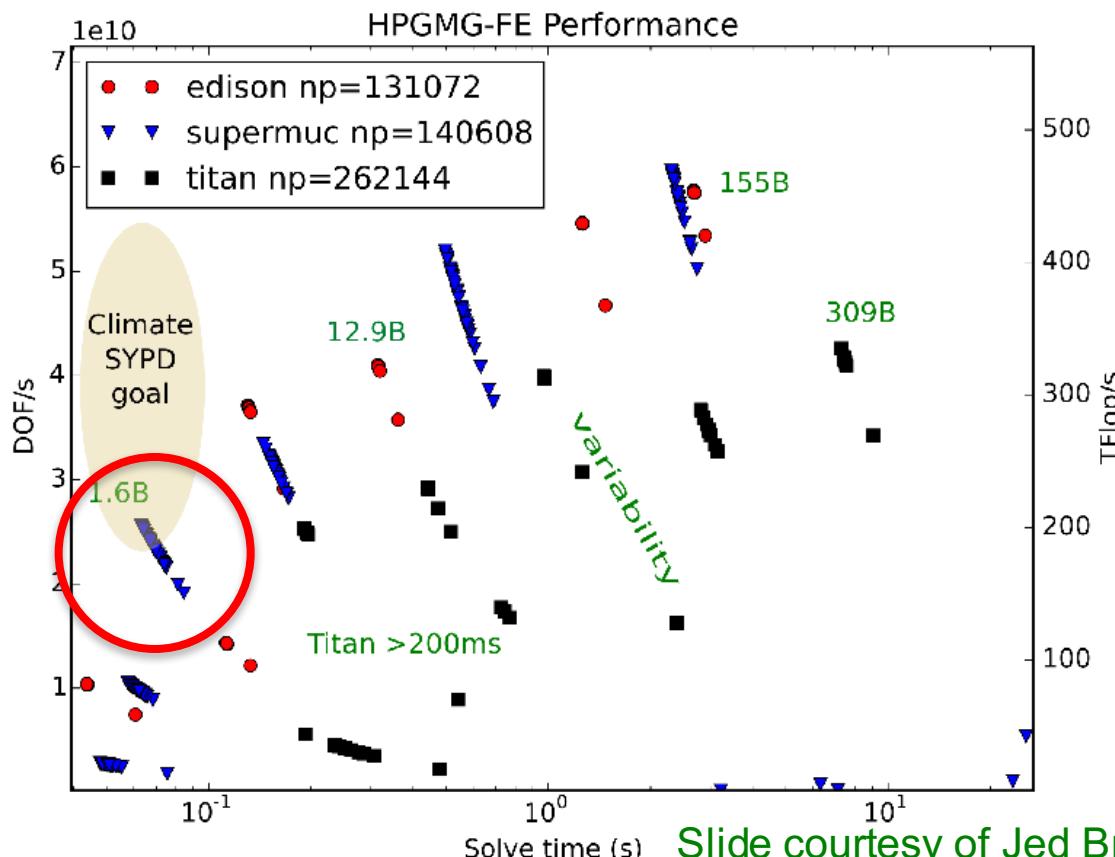
Multi-device performance: CPU vs GPU

note: projected ideal case (no comms; load imbalance)



Finite Element Benchmark Comparison of Large Systems

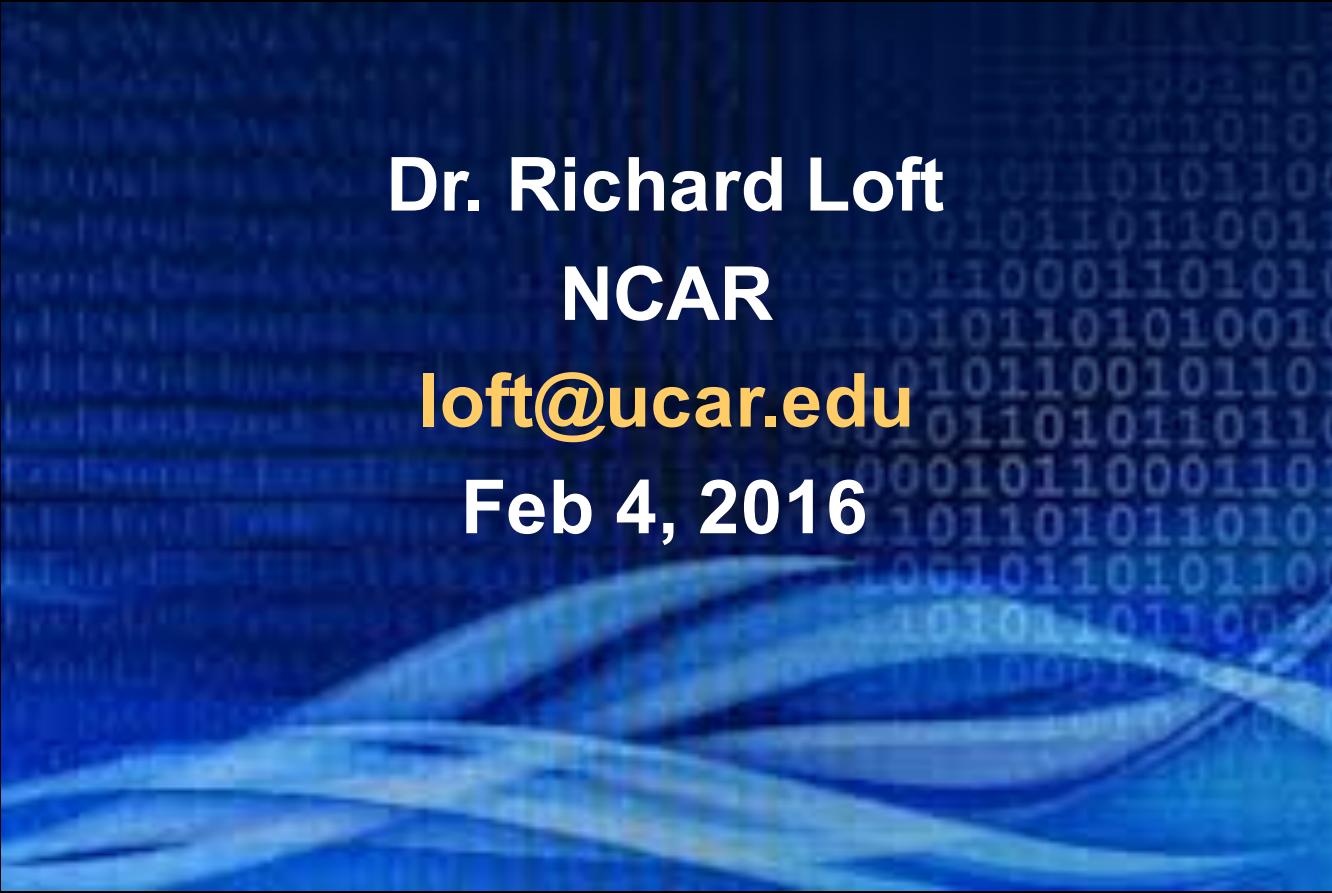
Scaling regime: HPGMG-FE on Edison, SuperMUC, Titan



Slide courtesy of Jed Brown, CU

41

Thanks! Question?



Dr. Richard Loft
NCAR
loft@ucar.edu
Feb 4, 2016

Image the size of supercomputer needed to match human brain!

Human Being



Computer X?



=

86 billion neurons
150 trillion synapses

~100,000 x honey bee
~20 x cat

1.7 million times
faster than
Yellowstone!

Cirrascale GX-8



8-way GPU cluster in a box

