



COMPUTER ARCHITECTURE

Parallel Architectures: Models and Tools

□ Performance improvements:

▣ Improvements in semiconductor technology

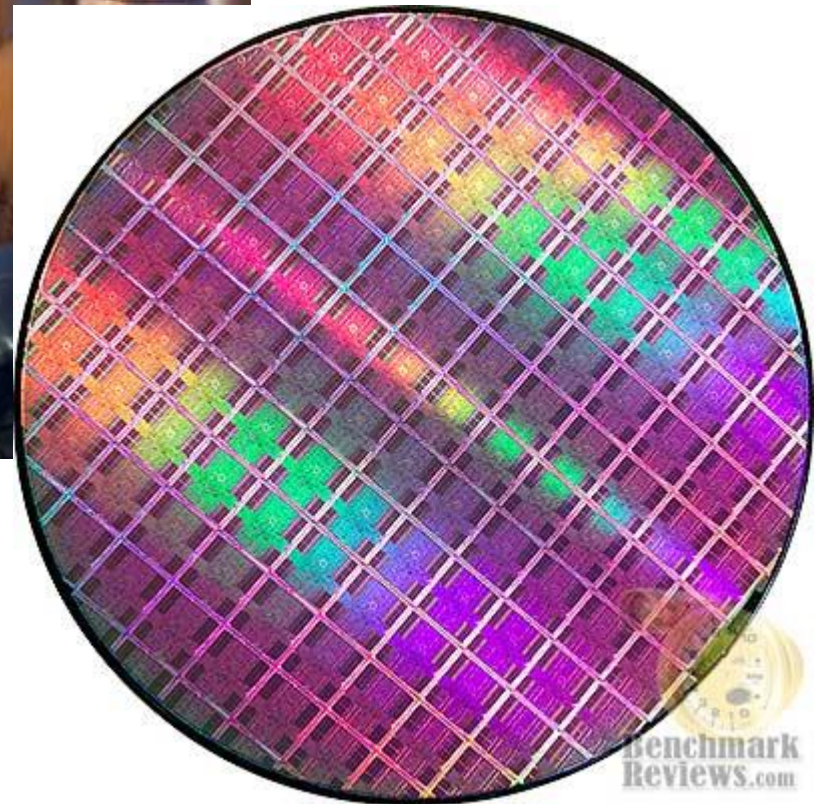
- Feature size, clock speed

▣ Improvements in computer architectures

- Enabled by HLL compilers, UNIX
- Lead to RISC architectures

▣ Together have enabled:

- Lightweight computers
- Productivity-based managed/interpreted programming languages



□ Integrated circuit

$$\text{Cost of integrated circuit} = \frac{\text{Cost of die} + \text{Cost of testing die} + \text{Cost of packaging and final test}}{\text{Final test yield}}$$

$$\text{Cost of die} = \frac{\text{Cost of wafer}}{\text{Dies per wafer} \times \text{Die yield}}$$

$$\text{Dies per wafer} = \frac{\pi \times (\text{Wafer diameter}/2)^2}{\text{Die area}} - \frac{\pi \times \text{Wafer diameter}}{\sqrt{2} \times \text{Die area}}$$

□ Bose-Einstein formula:

$$\text{Die yield} = \text{Wafer yield} \times 1 / (1 + \text{Defects per unit area} \times \text{Die area})^N$$

- Defects per unit area = 0.016-0.057 defects per square cm (2010)
- N = process-complexity factor = 11.5-15.5 (40 nm, 2010)

- Wafer with a diameter of 30 cm.
 - ▣ Dies of 1.5 cm side.
 - Dies per wafer: 269.
 - ▣ Dies of 1 cm side
 - Dies per wafer: 640.

- Integrated circuit technology
 - ▣ Transistor density: 35%/year
 - ▣ Die size: 10-20%/year
 - ▣ Integration overall: 40-55%/year

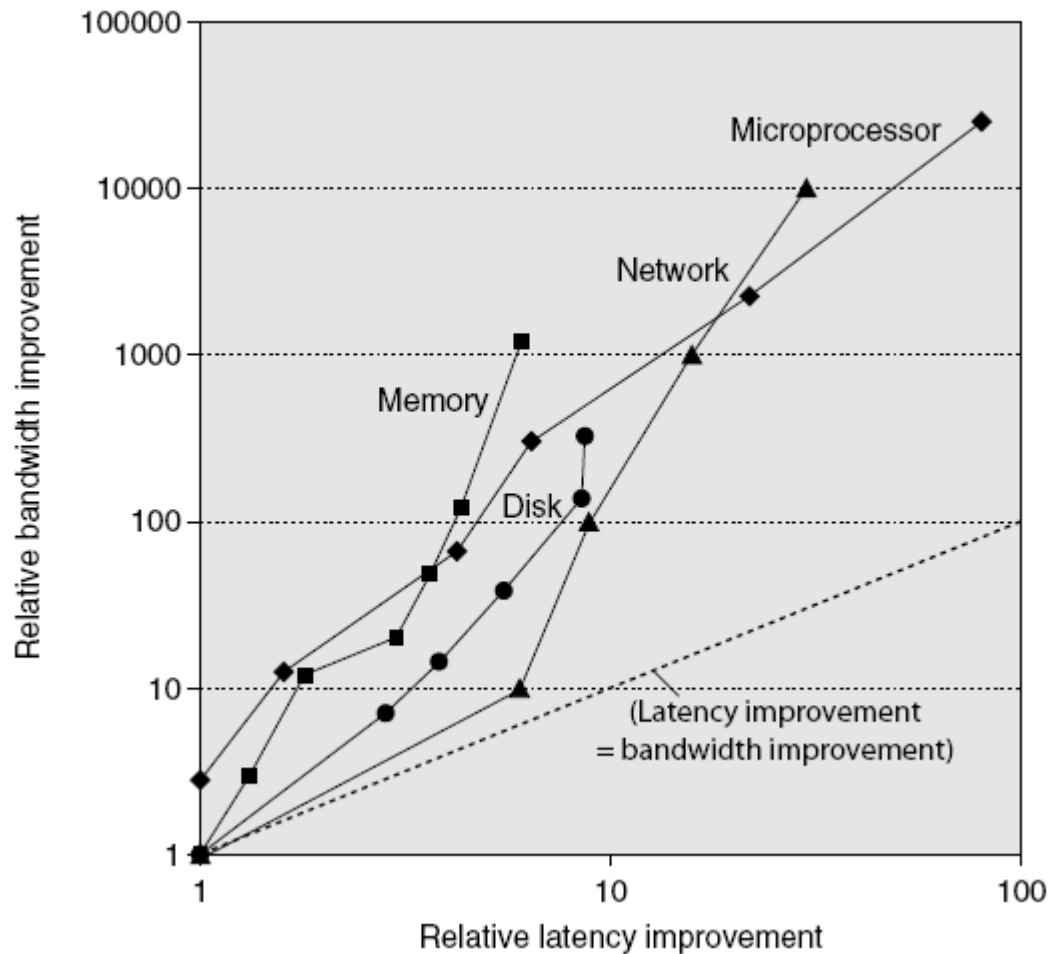
- DRAM capacity: 25-40%/year (slowing)

- Flash capacity: 50-60%/year
 - ▣ 15-20X cheaper/bit than DRAM

- Magnetic disk technology: 40%/year
 - ▣ 15-25X cheaper/bit than Flash
 - ▣ 300-500X cheaper/bit than DRAM

- Bandwidth or throughput
 - ▣ Total work done in a given time
 - ▣ 10,000-25,000X improvement for processors
 - ▣ 300-1 200X improvement for memory and disks

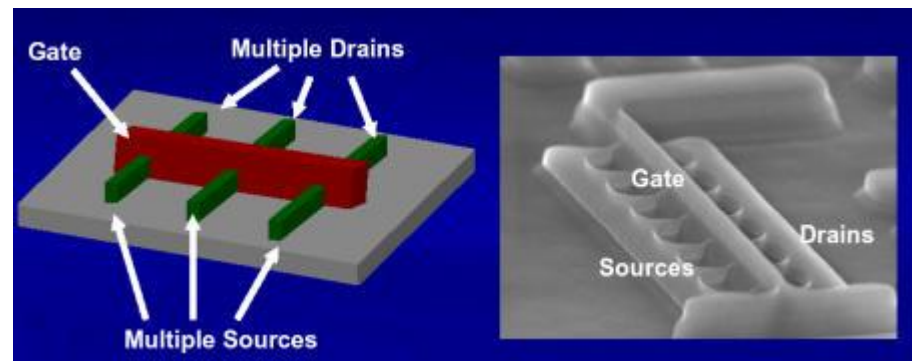
- Latency or response time
 - ▣ Time between start and completion of an event
 - ▣ 30-80X improvement for processors
 - ▣ 6-8X improvement for memory and disks



Log-log plot of bandwidth and latency milestones

□ Feature size

- ▣ Minimum size of transistor or wire in x or y dimension
- ▣ 10 microns in 1971 to .014 microns in 2014
- ▣ Transistor performance scales linearly
 - Wire delay does not improve with feature size!
- ▣ Integration density scales quadratically



- Problem: Get power in, get power out
 - ▣ Distribute power to increasingly complex circuitry
- Thermal Design Power (TDP)
 - ▣ Characterizes sustained power consumption
 - ▣ Used as target for power supply and cooling system
 - ▣ Lower than peak power, higher than average power consumption
 - ▣ Dark silicon
- Clock rate can be reduced dynamically to limit power consumption
- Energy per task is often a better measurement

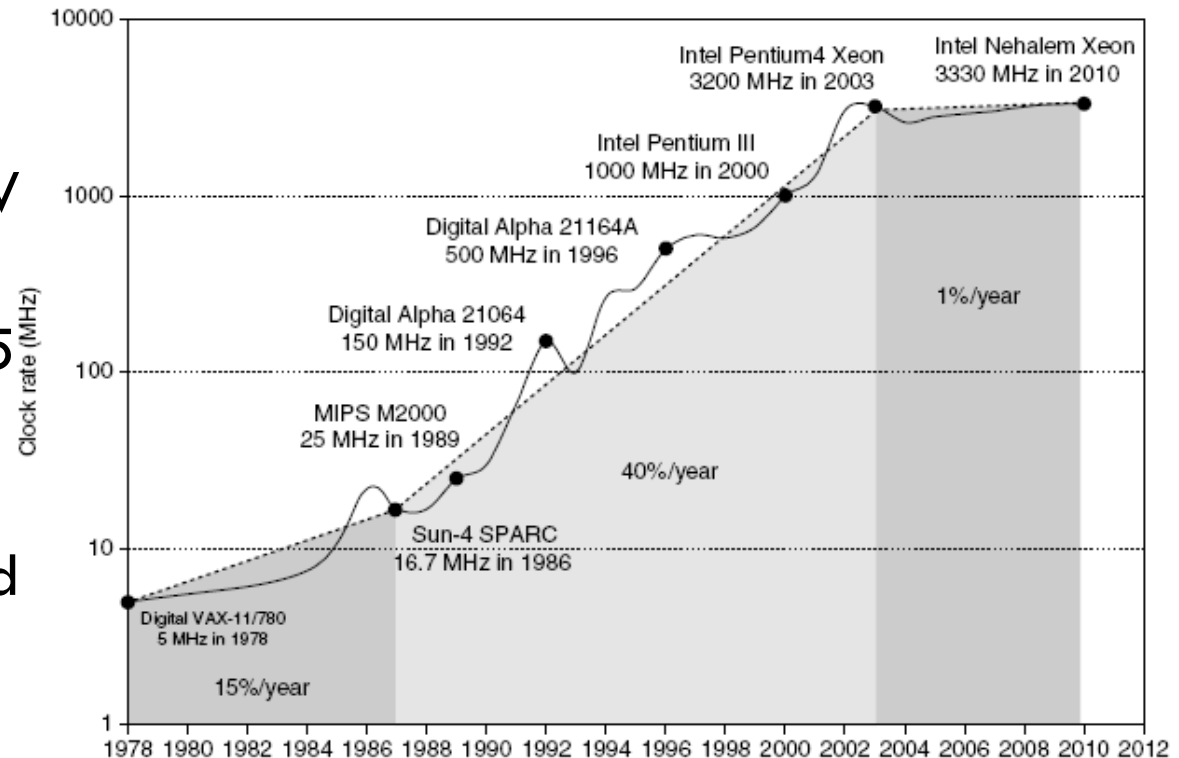
- Dynamic energy
 - ▣ Transistor switch from 0 \rightarrow 1 or 1 \rightarrow 0
 - ▣ $\frac{1}{2} \times \text{Capacitive load} \times \text{Voltage}^2$

- Dynamic power
 - ▣ $\frac{1}{2} \times \text{Capacitive load} \times \text{Voltage}^2 \times \text{Frequency switched}$

- For a fixed task reducing clock rate reduces power, not energy

- Voltage reduces both: has dropped from 5V to 1V in 20 years

- Intel 80386 consumed ~ 2 W
- 3.3 GHz Intel Core i7 consumes 130 W
- Heat must be dissipated from 1.5×1.5 cm chip
- This is the limit of what can be cooled by air



□ Static power consumption

- ▣ Due to leakage current flow

$$\text{Power}_{\text{static}} = \text{Current}_{\text{static}} \times \text{Voltage}$$

- ▣ Scales with number of transistors
- ▣ To reduce: power gating even to inactive modules
- ▣ Goal 2006 for leakage: 25% of total power consumption

- Techniques for reducing power:
 - ▣ Do nothing well
 - ▣ Dynamic Voltage-Frequency Scaling
 - ▣ Low power state for DRAM, disks
 - ▣ Overclocking, turning off cores

- Cost driven down by learning curve
 - ▣ Yield

- DRAM: price closely tracks cost

- Microprocessors: price depends on volume
 - ▣ Volume decrease the time needed to get down the learning curve.
 - ▣ Volume decreases cost, since it increases purchasing and manufacturing efficiency.
 - ▣ 10% less for each doubling of volume.

□ Module reliability

- ▣ Mean time to failure (MTTF)
- ▣ Mean time to repair (MTTR)
- ▣ Mean time between failures (MTBF) = $MTTF + MTTR$
- ▣ Availability = $MTTF / MTBF$

- Typical performance metrics:
 - ▣ Response time
 - ▣ Throughput

- Speedup of X relative to Y
 - ▣ $\text{Execution time}_Y / \text{Execution time}_X$

- Execution time
 - ▣ Wall clock time: includes all system overheads
 - ▣ CPU time: only computation time in the CPU

- Benchmarks
 - ▣ Kernels (e.g. matrix multiply)
 - ▣ Toy programs (e.g. sorting)
 - ▣ Synthetic benchmarks (e.g. Dhrystone)
 - ▣ Benchmark suites (e.g. SPEC06fp, TPC-C)

- Embedded
 - ▣ Dhrystone .
 - ▣ EEMBC (kernels).

- Desktop:
 - ▣ SPEC2006 (interger and floating point programs).

- Servers:
 - ▣ SPECWeb, SPECSFS, SPECjbb, SPECvirt_Sc2010.
 - ▣ TPC

- The only valid performance metric is the execution of real programs.
 - ▣ Any other metric is prone to errors.
 - ▣ Any other alternative to real programs is prone to errors.

SPEC2006 benchmark description	Benchmark name by SPEC generation				
	SPEC2006	SPEC2000	SPEC95	SPEC92	SPEC89
GNU C compiler					gcc
Interpreted string processing			perl		espresso
Combinatorial optimization		mcf			li
Block-sorting compression		bzip2		compress	eqntott
Go game (AI)	go	vortex	go	sc	
Video compression	h264avc	gzip	jpeg		
Games/path finding	astar	eon	m88ksim		
Search gene sequence	hmmer	twolf			
Quantum computer simulation	libquantum	vortex			
Discrete event simulation library	omnetpp	vpr			
Chess game (AI)	sjeng	crafty			
XML parsing	xalancbmk	parser			
CFD/blast waves	bwaves				fpppp
Numerical relativity	cactusADM				tomcatv
Finite element code	calculix				doduc
Differential equation solver framework	dealll				nasa7
Quantum chemistry	gamess				spice
EM solver (freq/time domain)	GemsFDTD			swim	matrix300
Scalable molecular dynamics (~NAMD)	gromacs		apsi	hydro2d	
Lattice Boltzman method (fluid/air flow)	lbm		mgrid	su2cor	
Large eddy simulation/turbulent CFD	LESlie3d	wupwise	applu	wave5	
Lattice quantum chromodynamics	milc	apply	turb3d		
Molecular dynamics	namd	galgel			
Image ray tracing	povray	mesa			
Sparse linear algebra	soplex	art			
Speech recognition	sphinx3	equake			
Quantum chemistry/object oriented	tonto	facerec			
Weather research and forecasting	wrf	ammp			
Magneto hydrodynamics (astrophysics)	zeusmp	lucas			
		fma3d			
		sixtrack			

- Speedup (plus low prog. effort and resource needs)

$$\text{Speedup}(p) = \frac{\text{Performance}(p)}{\text{Performance}(1)}$$

- For a fixed problem:

$$\text{Speedup}(p) = \frac{\text{Time}(1)}{\text{Time}(p)}$$

- Take Advantage of Parallelism
 - ▣ e.g. multiple processors, disks, memory banks, pipelining, multiple functional units

- Principle of Locality
 - ▣ Reuse of data and instructions

- Focus on the Common Case
 - ▣ Amdahl's Law

$$\text{Execution time}_{\text{new}} = \text{Execution time}_{\text{old}} \times \left((1 - \text{Fraction}_{\text{enhanced}}) + \frac{\text{Fraction}_{\text{enhanced}}}{\text{Speedup}_{\text{enhanced}}} \right)$$

$$\text{Speedup}_{\text{overall}} = \frac{\text{Execution time}_{\text{old}}}{\text{Execution time}_{\text{new}}} = \frac{1}{(1 - \text{Fraction}_{\text{enhanced}}) + \frac{\text{Fraction}_{\text{enhanced}}}{\text{Speedup}_{\text{enhanced}}}}$$

□ Suppose a fraction f of your application is not parallelizable

□ $1-f$: parallelizable on p processors

$$\text{Speedup}(P) = T_1 / T_p$$

$$\leq T_1 / (f T_1 + (1-f) T_1 / p) = 1 / (f + (1-f)/p)$$

$$\leq 1/f$$

- A web server has the following ratio of the execution time:
 - ▣ Computation: 40%
 - ▣ I/O: 60%

- If we replace this computer with another that is 10 times faster in computation, what is the overall speedup?

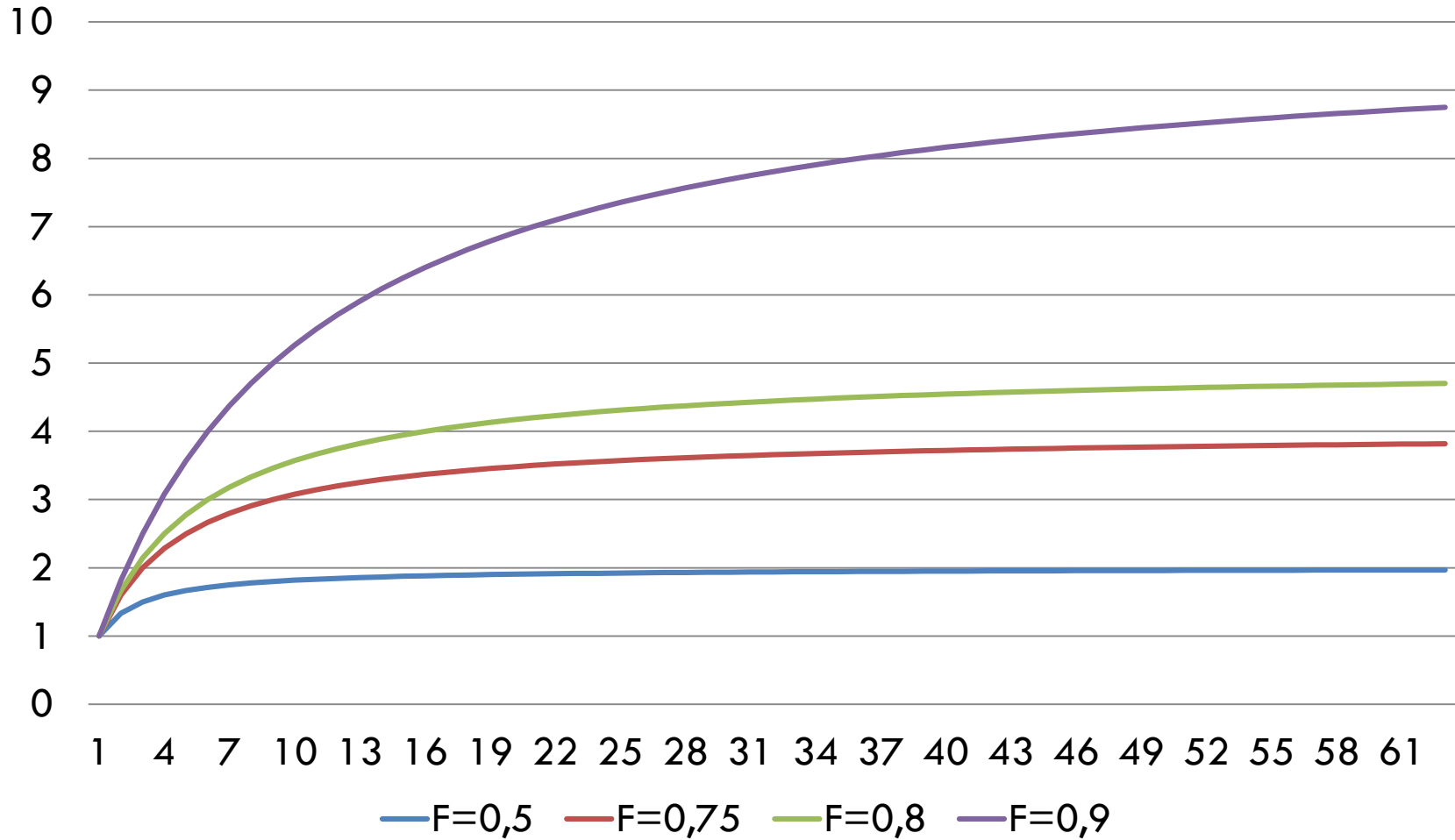
$$S = \frac{1}{0.6 + \frac{0.4}{10}} = \frac{1}{0.64} = 1.5625 < 1.666 = 1/0.6$$

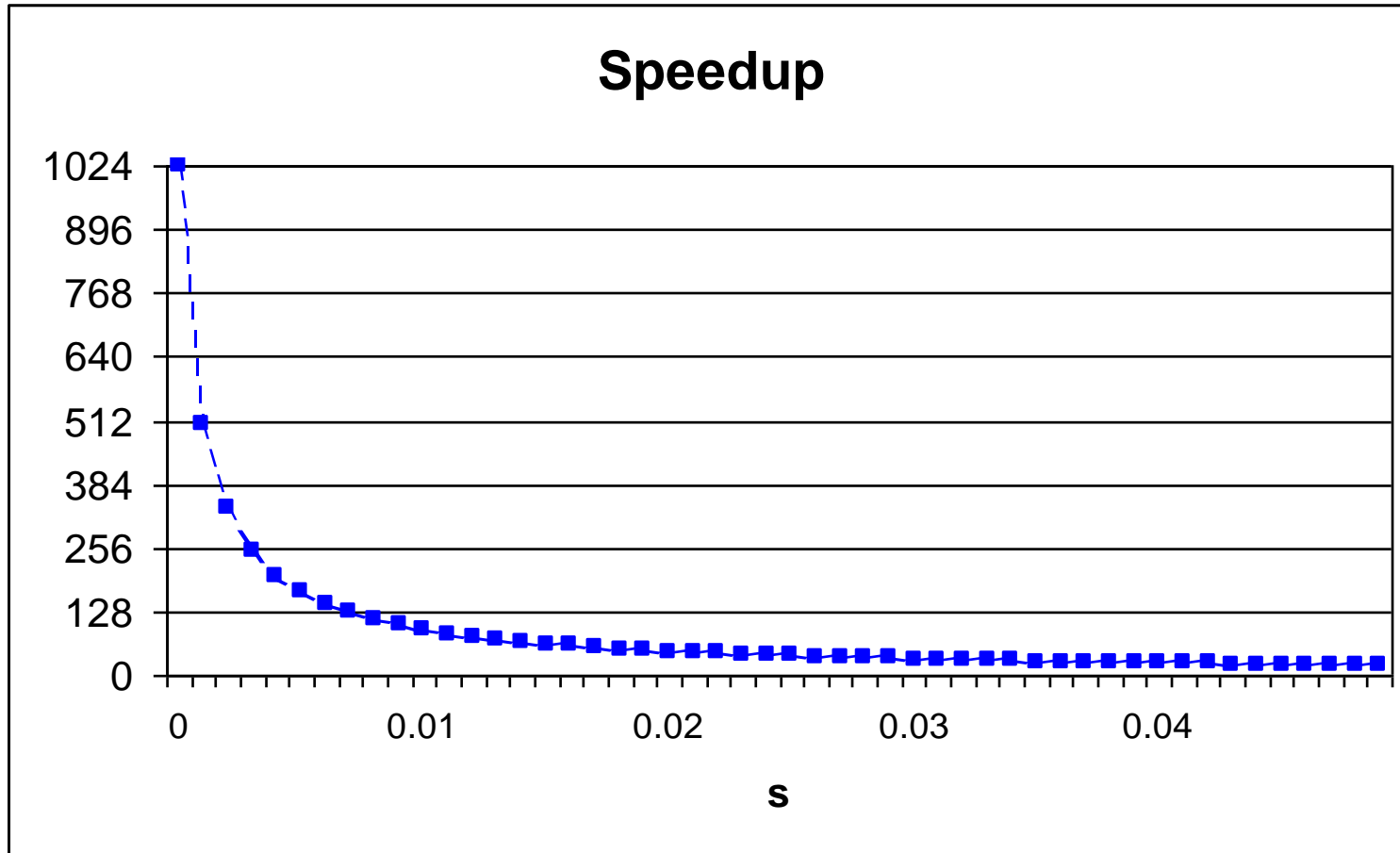
- An application has a parallel portion that takes 50% of the execution time.

- ▣ We execute the application in a 32-processor computer, what is the maximum speedup?

$$S = \frac{1}{0.5 + \frac{0.5}{32}} = \frac{1}{0.515625} = 1.9393$$

Speedup

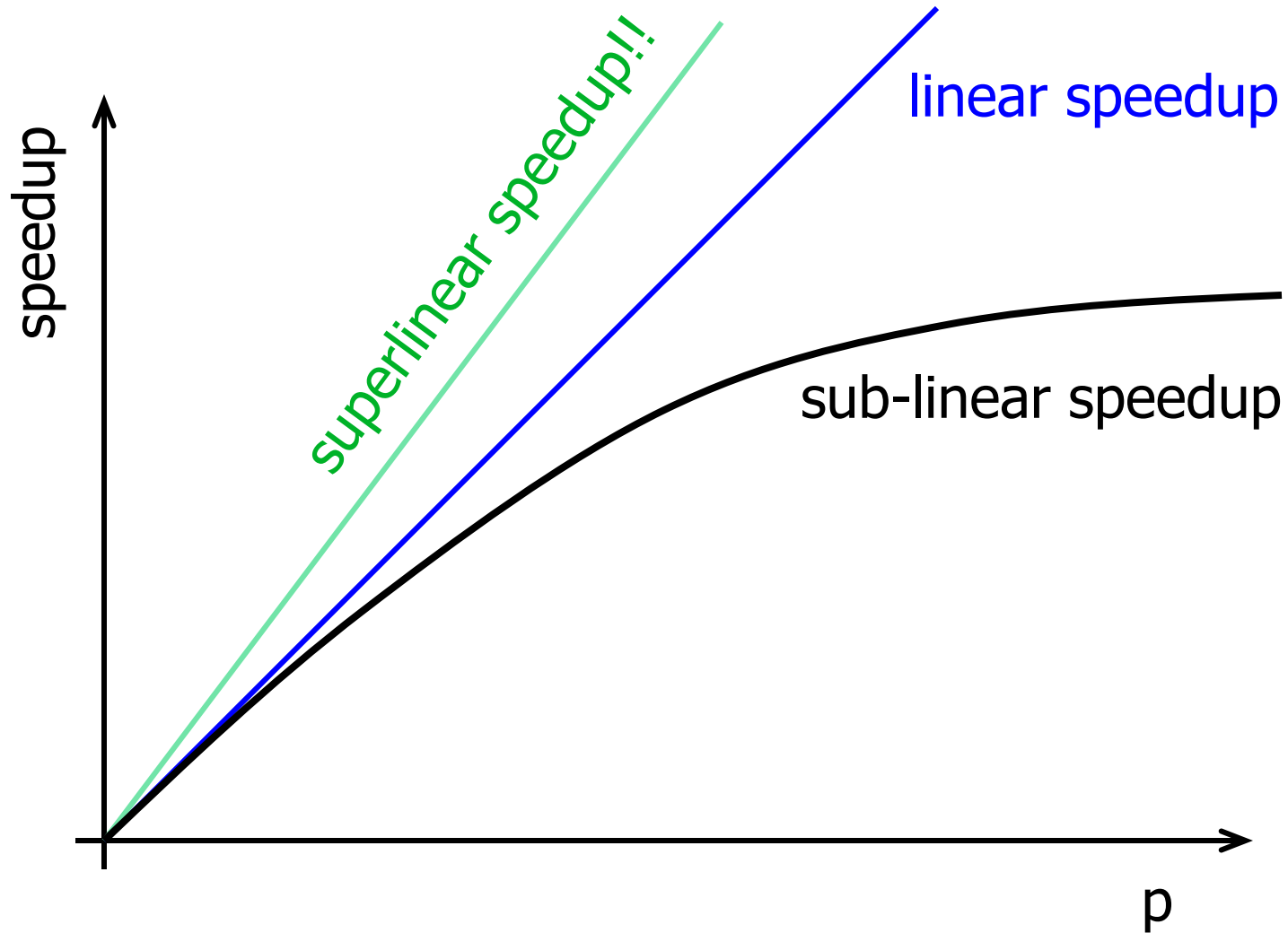




See: Gustafson, Montry, Benner, "Development of Parallel Methods for a 1024 Processor Hypercube", SIAM J. Sci. Stat. Comp. 9, No. 4, 1988, pp.609.

□ But:

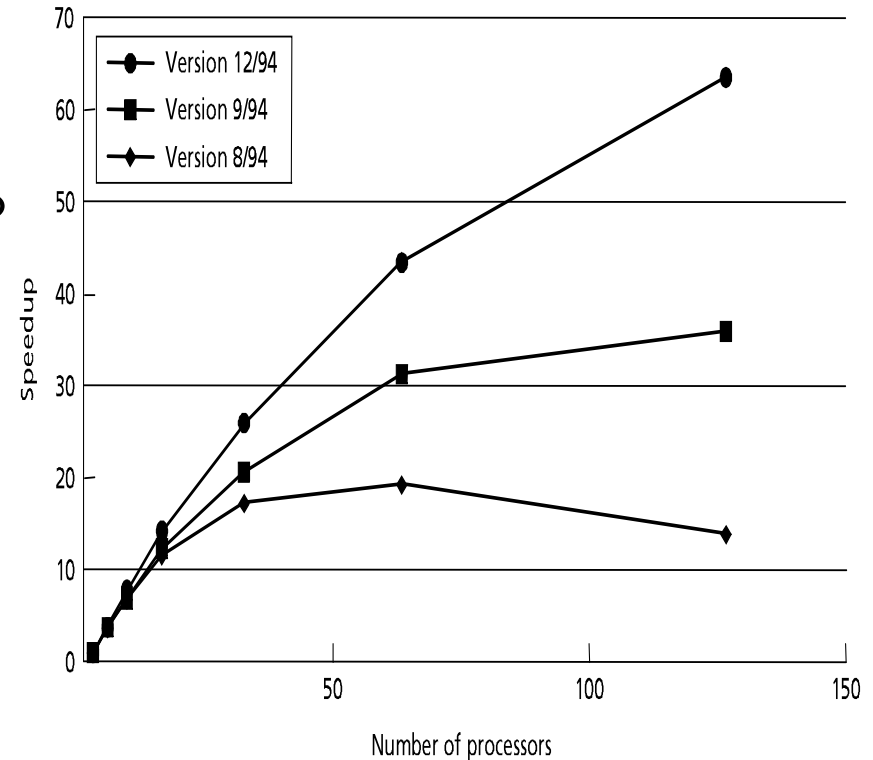
- There are many problems can be “embarrassingly” parallelized
 - Ex: image processing, differential equation solver
- “In some cases the serial fraction does not increase with the problem size
- Additional speedup can be achieved from additional resources (super-linear speedup due to more memory)



- Possible causes
- Algorithm
 - ▣ e.g., with optimization problems, throwing many processors at it increases the chances that one will “get lucky” and find the optimum fast
- Hardware
 - ▣ e.g., with many processors, it is possible that the entire application data resides in cache (vs. RAM) or in RAM (vs. Disk)

- $\text{Eff}_p = S_p / p$
- Typically 1, unless superlinear speedup
- Used to measure how well the processors are utilized
 - ▣ If increasing the number of process by a factor 10 increases the speedup by a factor 2, perhaps it's not worth it: efficiency drops by a factor 5

- Architect Goal
 - ▣ observe how program uses machine and improve the design to enhance performance
- Programmer Goal
 - ▣ observe how the program uses the machine and improve the implementation to enhance performance



- Amdahl's law focuses on the negative point of view of parallel processing
- However:
 - ▣ Parallel machines are used for solving large problems.
 - ▣ A sequential computer could never execute a large parallel program.
 - Memory limits.
 - Processing limits.



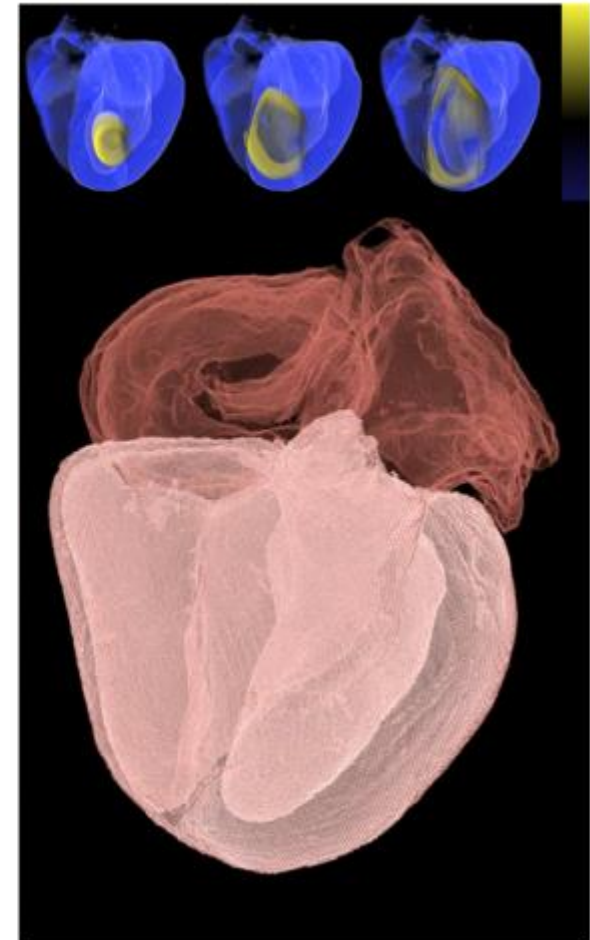
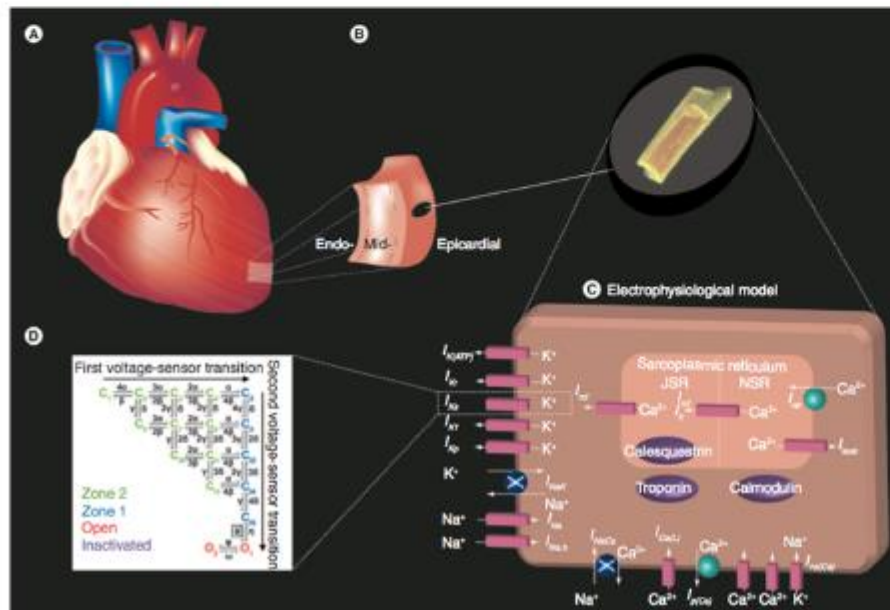
$$S = \frac{T_s}{T_p}$$

T_s = Time in a sequential machine

T_p = Time in a parallel machine

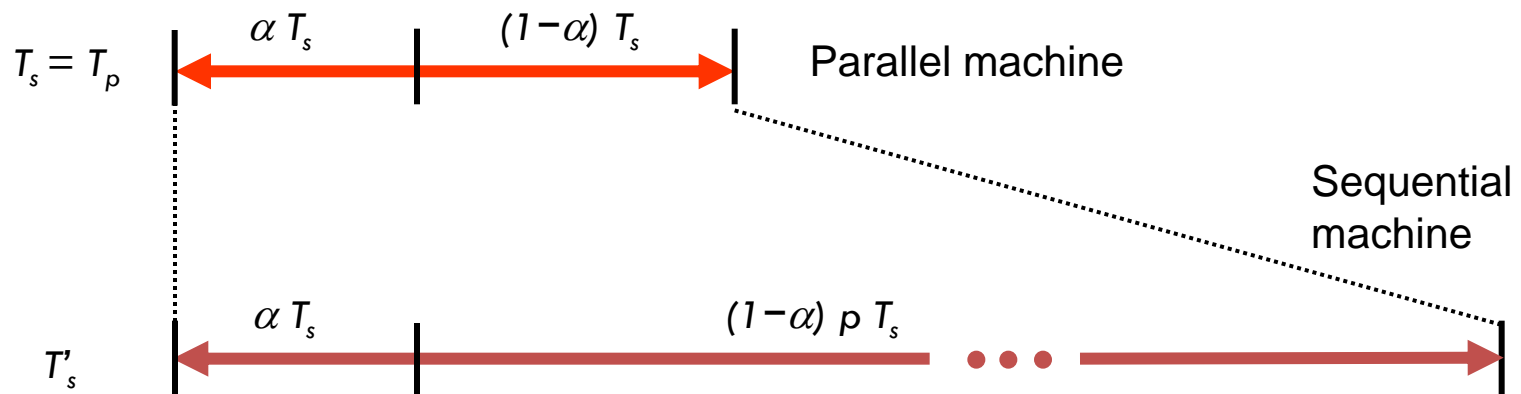
Computational Medicine: Whole Organ Simulation

- **Predictive Toxicology**
- **Multiscale Model of Organs**
 - from protein function through to cell function through to tissue function through to macroscale organ modeling.
- **Multiple model components and scales require Petascale to Exascale compute capability**
 - Usefulness requires “turnkey” modeling environment where many variations and scenarios can be attempted by the medical or pharmaceutical researcher quickly and accurately
 - Further increases the computational requirements



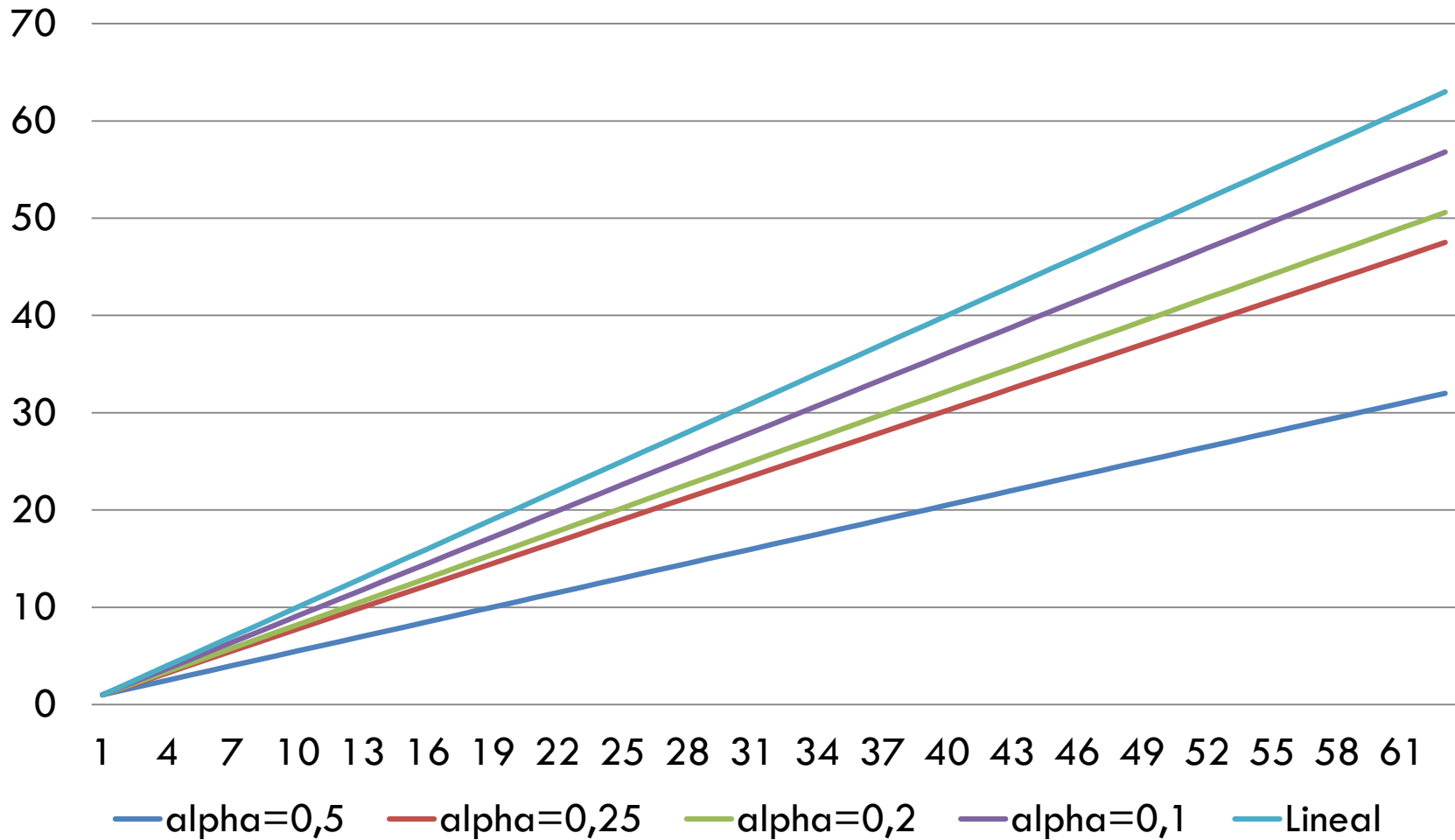
- The amount of work changes with the number of processors

$$S_p = \frac{T'_s}{T_p} = \frac{\alpha T_s + (1-\alpha)pT_s}{T_s} = p + \alpha(1-p)$$



- The sequential portion of the program decreases with program size.
 - ▣ When the problem size grows we can assume a close-to-linear speedup ($S \approx p$).
- Using parallelism, we can approach larger problems.

Speedup



□ The Processor Performance Equation

CPU time = CPU clock cycles for a program \times Clock cycle time

$$\text{CPU time} = \frac{\text{CPU clock cycles for a program}}{\text{Clock rate}}$$

$$\text{CPI} = \frac{\text{CPU clock cycles for a program}}{\text{Instruction count}}$$

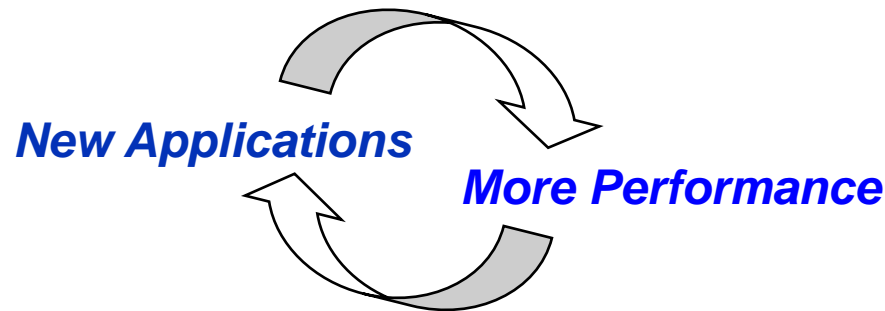
CPU time = Instruction count \times Cycles per instruction \times Clock cycle time

$$\frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Clock cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{Clock cycle}} = \frac{\text{Seconds}}{\text{Program}} = \text{CPU time}$$

- Different instruction types having different CPIs

$$\text{CPU clock cycles} = \sum_{i=1}^n \text{IC}_i \times \text{CPI}_i$$

$$\text{CPU time} = \left(\sum_{i=1}^n \text{IC}_i \times \text{CPI}_i \right) \times \text{Clock cycle time}$$



- Demand for cycles fuels advances in hardware, and vice-versa
 - ▣ Cycle drives exponential increase in microprocessor performance
 - ▣ Drives parallel architecture harder: most demanding applications
- Goal of applications in using parallel machines: Speedup

$$\text{Speedup (p processors)} = \frac{\text{Performance (p processors)}}{\text{Performance (1 processor)}}$$

- For a fixed problem size (input data set), performance = 1 / time

$$\text{Speedup fixed problem (p processors)} = \frac{\text{Time (1 processor)}}{\text{Time (p processors)}}$$

- Science
 - ▣ Global climate modeling
 - ▣ Astrophysical modeling
 - ▣ Biology: genomics; protein folding; drug design
 - ▣ Computational Chemistry
 - ▣ Computational Material Sciences and Nanosciences
- Engineering
 - ▣ Crash simulation
 - ▣ Semiconductor design
 - ▣ Earthquake and structural modeling
 - ▣ Computation fluid dynamics (airplane design)
 - ▣ Combustion (engine design)
- Business
 - ▣ Financial and economic modeling
 - ▣ Transaction processing, web services and search engines
- Defense
 - ▣ Nuclear weapons -- test by simulations
 - ▣ Cryptography

- 1 PFLOP has been surpassed in 2008
- Currently:
 - ▣ 33 PFLOPS
 - ▣ 3.1M cores system
- We head toward ExaScale age
 - ▣ 1,000,000,000 cores
- Increased probabilities of failures
 - ▣ Learn to live with failures
 - ▣ Fault tolerance
 - ▣ Learn to continue in the presence of failures
- Challenges in getting a global view of the system
- New challenges for applications and algorithms
- Scale invariance targeted
 - ▣ Local versus global
 - ▣ Learn from Internet
 - ▣ Learn from nature: evolution, adaptation, swarm behaviour
- Energy efficiency target: 20MW for an Exascale system (50x improvement)

- ❑ Since 1993 twice a year: June and November
- ❑ Ranking of the most powerful computing systems in the world
- ❑ Ranking criteria: performance of the LINPACK benchmark
- ❑ Jack Dongarra alma máter
- ❑ Site web: www.top500.org
- ❑ Poster 2012:
http://www.top500.org/static/lists/2012/06/TOP500_201206_Poster.pdf



Rank	Site	Computer/Year Vendor	Cores	R_{\max}	R_{peak}	Power
1	DOE/NNSA/LLNL United States	Sequoia - BlueGene/Q, Power BQC 16C 1.60 GHz, Custom / 2011 IBM	1572864	16324.75	20132.66	7890.0
2	RIKEN Advanced Institute for Computational Science (AICS) Japan	K computer , SPARC64 VIIIfx 2.0GHz, Tofu interconnect / 2011 Fujitsu	705024	10510.00	11280.38	12659.9
3	DOE/SC/Argonne National Laboratory United States	Mira - BlueGene/Q, Power BQC 16C 1.60GHz, Custom / 2012 IBM	786432	8162.38	10066.33	3945.0
4	Leibniz Rechenzentrum Germany	SuperMUC - iDataPlex DX360M4, Xeon E5-2680 8C 2.70GHz, Infiniband FDR / 2012 IBM	147456	2897.00	3185.05	3422.7
5	National Supercomputing Center in Tianjin China	Tianhe-1A - NUDT YH MPP, Xeon X5670 6C 2.93 GHz, NVIDIA 2050 / 2010 NUDT	186368	2566.00	4701.00	4040.0
6	DOE/SC/Oak Ridge National Laboratory United States	Jaguar - Cray XK6, Opteron 6274 16C 2.200GHz, Cray Gemini interconnect, NVIDIA 2090 / 2009 Cray Inc.	298592	1941.00	2627.61	5142.0
7	CINECA Italy	Fermi - BlueGene/Q, Power BQC 16C 1.60GHz, Custom / 2012 IBM	163840	1725.49	2097.15	821.9
8	Forschungszentrum Juelich (FZJ) Germany	JuQUEEN - BlueGene/Q, Power BQC 16C 1.60GHz, Custom / 2012 IBM	131072	1380.39	1677.72	657.5
9	CEA/TGCC-GENCI France	Curie thin nodes - Bullx B510, Xeon E5- 2680 8C 2.700GHz, Infiniband QDR / 2012 Bull	77184	1359.00	1667.17	2251.0
10	National Supercomputing Centre in Shenzhen (NSCS) China	Nebulae - Dawning TC3600 Blade System, Xeon X5650 6C 2.66GHz, Infiniband QDR, NVIDIA 2050 / 2010 Dawning	120640	1271.00	2984.30	2580.0

- For a long time performance has been the only metric
 - ▣ FLOPS
 - ▣ Total cost of ownership (TCO) neglected
- Conscience about increasing costs of power, maintenance, administration, failure recovery
- Ranking of the most energy-efficient supercomputers in the world
 - ▣ MFLOPS/Watt
- First edition: November 2007
- Last release: June 2012

Green500 Rank	MFLOPS/W	Site*	Computer*	Total Power (kW)
1	2,100.88	DOE/NNSA/LLNL	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	41.10
2	2,100.88	IBM Thomas J. Watson Research Center	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	41.10
3	2,100.86	DOE/SC/Argonne National Laboratory	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	82.20
4	2,100.86	DOE/SC/Argonne National Laboratory	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	82.20
5	2,100.86	Rensselaer Polytechnic Institute	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	82.20
6	2,100.86	University of Rochester	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	82.20
7	2,100.86	IBM Thomas J. Watson Research Center	BlueGene/Q, Power BQC 16C 1.60 GHz, Custom	82.20
8	2,099.56	University of Edinburgh	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	493.10
9	2,099.50	Science and Technology Facilities Council - Daresbury Laboratory	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	575.30
10	2,099.46	Forschungszentrum Juelich (FZJ)	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	657.50