

# Introduction to High Performance Computing

Lecture, WS 2016

Prof. Dr.rer.nat. Matthias S. Müller

- 0. Organization**
- 1. Why supercomputers?**
- 2. Modern processors**
- 3. Basic optimization techniques for serial code**
- 4. Data access optimization**
- 5. Parallel computers**
- 6. Parallelization and optimization strategies**
- 7. Parallel algorithms**
- 8. Shared-memory programming with OpenMP**
- 9. Distributed-memory programming with MPI**
- 10. Hybrid programming (MPI + OpenMP)**
- 11. Heterogeneous architectures (GPUs, Xeon Phis)**
- 12. Energy efficiency**

## 0. Organization

- About myself
- Lecturer
- Administration
- Scope of the course

1. Why supercomputers?
2. Modern processors
3. Basic optimization techniques for serial code
4. Data access optimization
5. Parallel computers
6. Parallelization and optimization strategies
7. Parallel algorithms
8. Shared-memory programming with OpenMP
9. Distributed-memory programming with MPI
10. Hybrid programming (MPI + OpenMP)
11. Heterogeneous architectures (GPUs, Xeon Phi)
12. Energy efficiency

→ First computer 1983

■ **1989: Study of physics at University Stuttgart**

→ Unix, Linux, Systemadministrator in der Mathematik

■ **1996: Promotion at Institute for Computer Applications**

→ System administrator at university

→ User of supercomputers in Stuttgart and Jülich

■ **1999: Höchstleistungsrechenzentrum Stuttgart, left as Deputy Director**

■ **2005: Deputy Director and CTO at ZIH, Dresden**

■ **2013: RWTH Aachen: Chair for High Performance Computing, Director of IT Center (*former Center for Computing and Communication*)**

## ■ Research

- Talks and publications
- Research proposals
- Conferences and workshops

## ■ Teaching

- Lecture
- Seminars
- Examinations

## ■ IT Center

- Management
- Project proposals
- Workshops and conferences
- Talks and presentations



**Julian Miller**

Accelerators, Numerics



**Jannis Klinkenberg**

Big Data, Databases



**Joachim Protze**

Correctness Checking,  
Message Passing (MPI)



**Sandra Wienke**

Accelerators, TCO  
OpenACC, OpenMP,  
CUDA



**Christian Terboven**

OpenMP, ARB



**Hristo Iliev**

Message Passing  
Interface (MPI)

## ■ Lecture: 3 SWS

- Mondays 10:15 – 11:45 in AH VI
- Tuesdays 14:15 – 15:45 in AH II

## ■ Exercises: 1 SWS

- During the dates above
- Frontal exercise (**exercises should be prepared beforehand at home**)
- In preparation, **create an HPC account**:  
<https://www.rwth-aachen.de/selfservice>  
**send TIM to contact@hpc.rwth-aachen.de**  
**with “[16ws-45645] TIM” as subject**

## ■ Prerequisites for exercises

- Basic programming knowledge in C/C++ or FORTRAN
- Using Linux/UNIX environment

## Lectures & Exercises as of 17th October 2016

*Updates will be announced on L2P*

KW	Date	Type	Date	Type
42	17.10.2016	V	18.10.2016	V
43	24.10.2016	V	25.10.2016	V
44	31.10.2016	V	01.11.2016	-
45	07.11.2016	Ü	08.11.2016	V
46	14.11.2016	Ü	15.11.2016	V
47	21.11.2016	V	22.11.2016	MF
48	28.11.2016	Ü	29.11.2016	V
49	05.12.2016	V	06.12.2016	V
50	12.12.2016	V	13.12.2016	V
51	19.12.2016	Ü	20.12.2016	V
02	09.01.2017	V	10.01.2017	Ü
03	16.01.2017	V	17.01.2017	V
04	23.01.2017	V	24.01.2017	Ü
05	30.01.2017	V	31.01.2017	V
06	06.02.2017	V	07.02.2017	F

KW: week, V: lecture, Ü: exercise

MF: Maschinenhallen/AIXCave-führung

F: Fragenstunde

- **Slides will be in English**
- **Lecture will be in German or English → Decide!**
- **Exam: written (~ 120 min)**
  - Language: German or English
- **Material**
  - Slides will be in L2P just before or right after the lecture
  - Exercises will be in L2P one week before the exercise
    - Solve them alone at home, discussion/ questions during frontal exercise
  - <https://www3.elearning.rwth-aachen.de/ws16/16ws-45645/Dashboard.aspx>
    - **Please enroll & follow updates!**



## ■ Post into discussion area in L2P

→ <https://www3.elearning.rwth-aachen.de/ws16/16ws-4564>

## ■ RWTH App: Questions, comments, surveys, etc. during the lectures

→ Chat channel “Introduction to High Performance Computing”

→ Password: “HPC”

→ More information on the app and download links:

<http://www.itc.rwth-aachen.de/cms/IT-Center/Forschung-Projekte/~fxfk/RWTH-App>

## ■ Or if you need individual help: [contact@hpc.rwth-aachen.de](mailto:contact@hpc.rwth-aachen.de)

## ■ Audience

- Computer Scientists, (Computational) Engineers, Software Systems Engineers,...
- BSCES, BSInf, MSETITTI, MSInf, MSSSE, MSTKI: 6 ECTS
- MSCES, MSVT: 4 ECTS

■ What are your expectations of the course?

■ What don't you expect (because you already know it)?

- 1. Why supercomputers?**
- 2. Modern processors**
- 3. Basic optimization techniques for serial code**
- 4. Data access optimization**
- 5. Parallel computers**
- 6. Parallelization and optimization strategies**
- 7. Parallel algorithms**
- 8. Shared-memory programming with OpenMP**
- 9. Distributed-memory programming with MPI**
- 10. Hybrid programming (MPI + OpenMP)**
- 11. Heterogeneous architectures (GPUs, Xeon Phis)**
- 12. Energy efficiency**
- 13. Tour of the cluster machine hall/ VR aixCAVE**

Motivation & recent issues

Understanding of possible architectural impacts on performance & serial performance tuning

Parallel computer architectures & how to optimize code for them

Parallel programming for different architectures

## 1. Why supercomputers?

- Examples, fundamental terms, standard benchmarks

## 2. Modern processors

- Microarchitecture
- Pipelining, superscalarity, SIMD, Prefetching, ...
- Memory hierarchies, caches, cache mapping

## 3. Basic optimization techniques for serial code

- Profiling, tracing, instrumentation
- Event-driven / sample-driven triggers
- Interval timers

## 4. Data access optimization

- Balanced metrics / roofline performance model
- Algorithm classification & access optimizations

## 5. Parallel computers

- Multicore/ multithreaded processors
- Shared-memory computers
- Distributed-memory computers
- Networking architectures for distributed memory systems

## 6. Parallelization and optimization strategies

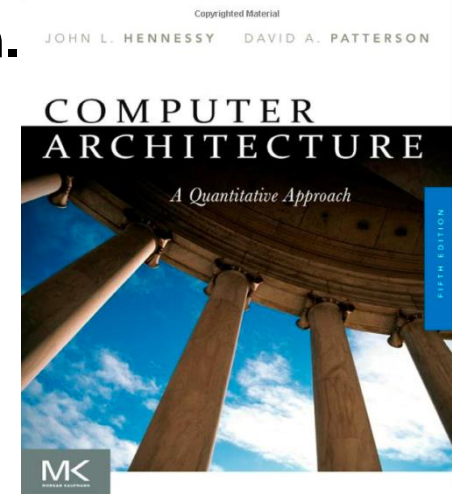
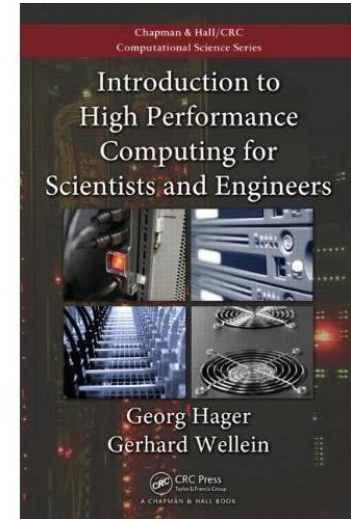
- Understanding parallel architectures & Finding concurrency
- Patterns of parallel algorithms
- Performance impacts & solutions (load imbalances, domain decomposition,...)

- 7. Parallel algorithms**
- 8. Distributed-memory programming with MPI**
- 9. Shared-memory programming with OpenMP**
- 10. Hybrid programming (MPI + OpenMP)**
- 11. Heterogeneous architectures (GPUs, Xeon Phis)**
- 12. Energy efficiency**

**Parallel  
programming**

- TOP500 trends
- Benchmarking
- Energy efficiency of microarchitectures

- G. Hager and G. Wellein:  
**Introduction to High Performance Computing for Scientists and Engineers**  
CRC Computation Science Series, 2010  
ISBN: 978-1-4398-1192-4.  
54,90 €.
- J. Hennessy and D. Patterson:  
**Computer Architecture. A Quantitative Approach.**  
Morgan Kaufmann Publishers, Elsevier, 2011,  
ISBN: 978-0123838728.  
69,50 €.



## 1. Why supercomputers?

- A few examples of supercomputers
- Top 500 list
- Why is parallelism important?

## 2. Modern processors

## 3. Basic optimization techniques for serial code

## 4. Data access optimization

## 5. Parallel computers

## 6. Parallelization and optimization strategies

## 7. Parallel algorithms

## 8. Shared-memory programming with OpenMP

## 9. Distributed-memory programming with MPI

## 10. Hybrid programming (MPI + OpenMP)

## 11. Heterogeneous architectures (GPUs, Xeon Phi)

## 12. Energy efficiency



# Why Supercomputers?

**VIDEO**

## ■ From Wikipedia:

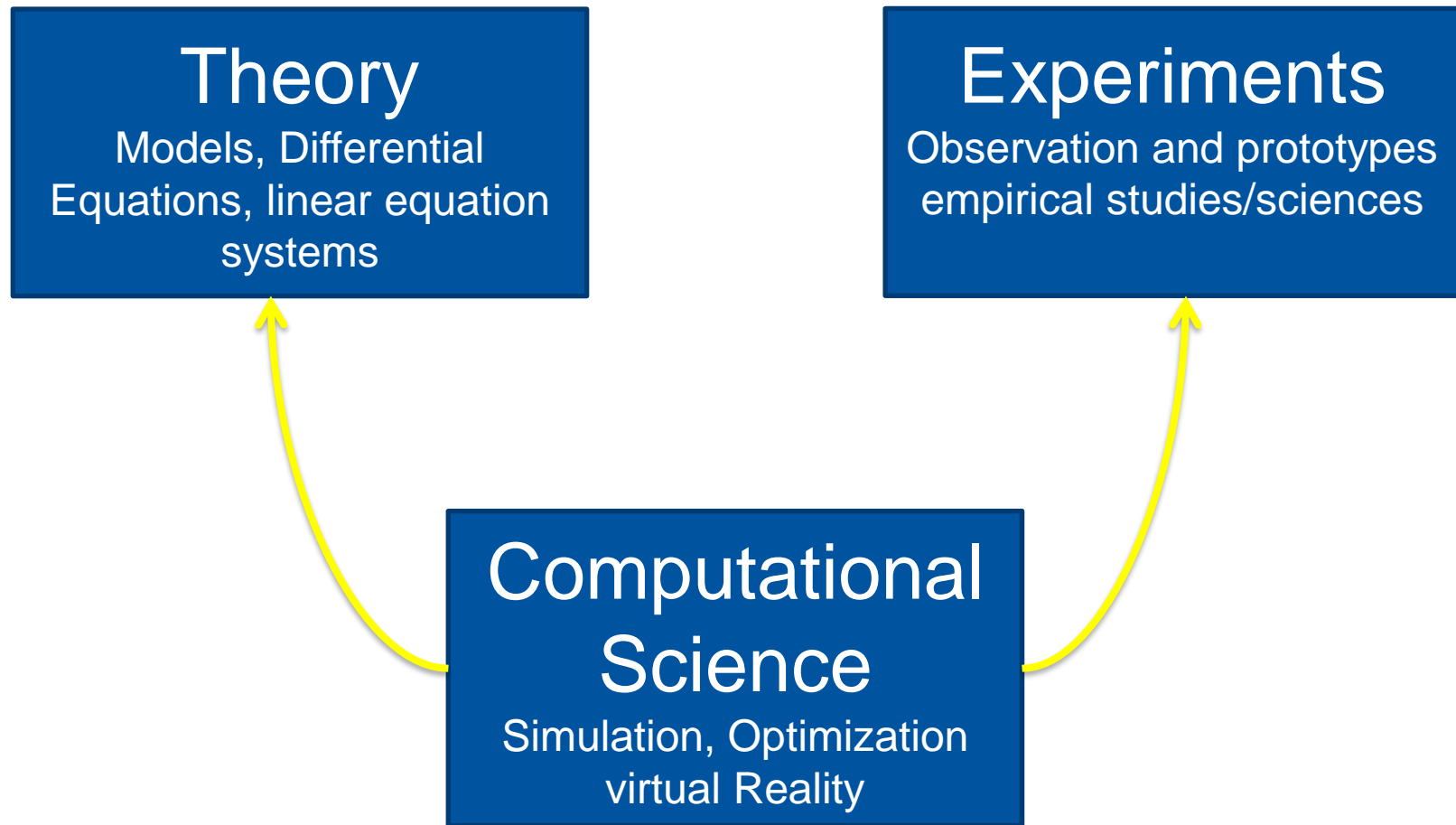
**“A supercomputer is a computer at the frontline of current processing capacity, particularly speed of calculation.**

**The term ‘Super Computing’ was first used in the *New York World* in 1929 to refer to large custom-built tabulators that IBM had made for Columbia University.”**

## ■ Why should we care about technical and scientific computing?

- Computational research complements experimental methods in all facets of engineering and science
- Certain cases: Computational simulations are the only option
  - Experiment may be cost prohibitive (e.g. Flight testing)
  - Experiment may be impossible (e.g. interaction effects between space station and space ship during docking)
- Scope of simulation depends heavily on the available computational power like amount of memory, amount of processors, performance of single processors, networking capabilities
  - super-computers
  - A comparison of pure performance is the TOP500 list

- Historically two principles of science – computational science extends them as a third



- **Simulation has become the third pillar of natural sciences along with theory and experimentation**

- **Applications vary from engineering...**

- Crash simulations

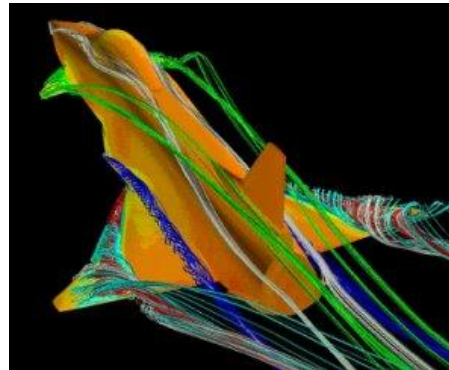
- Aerodynamics

- **over medicine...**

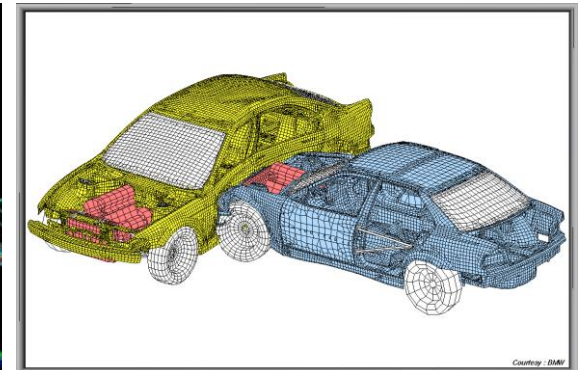
- **to meteorology**

- Weather forecast

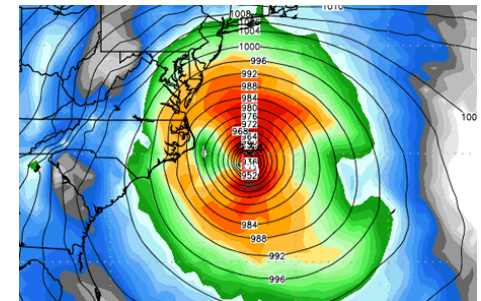
- Hurricane warnings



Courtesy of DLR



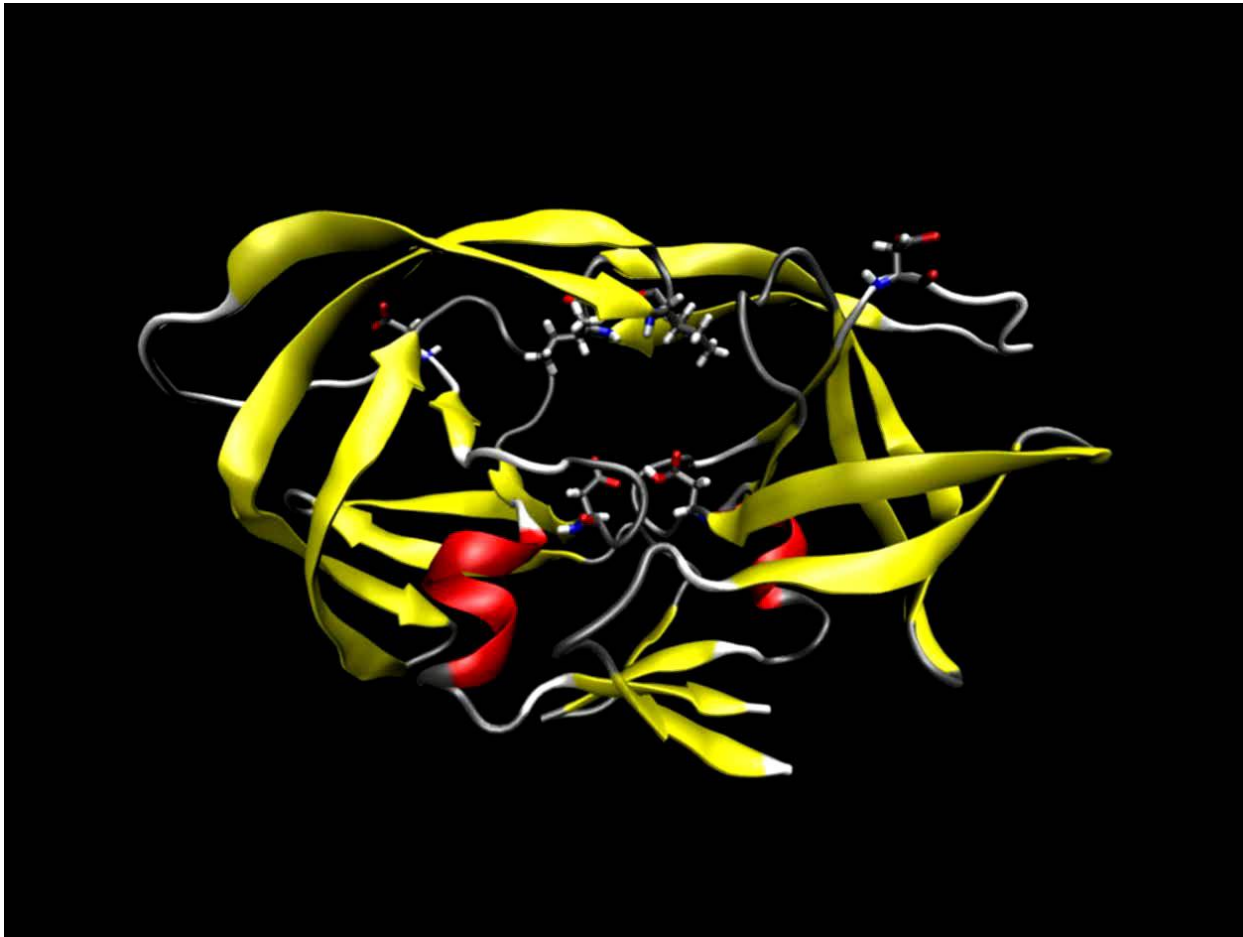
Courtesy of BMW



Courtesy of weatherbell

# Why Supercomputing?

## MD Simulation of HIV protease dynamics

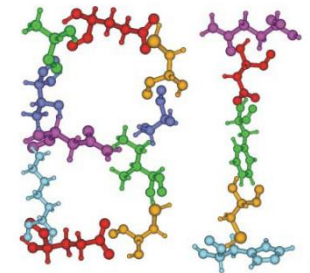


Time step:  
1fs ( $10^{-15}$  s)

Real time:  
10 ns ( $10^{-8}$  s)

Compute time:  
18.000 CPU-hrs

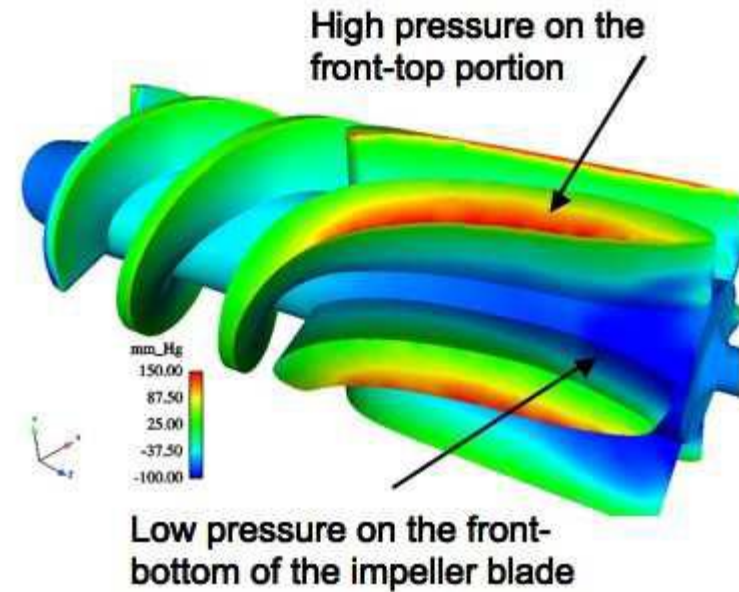
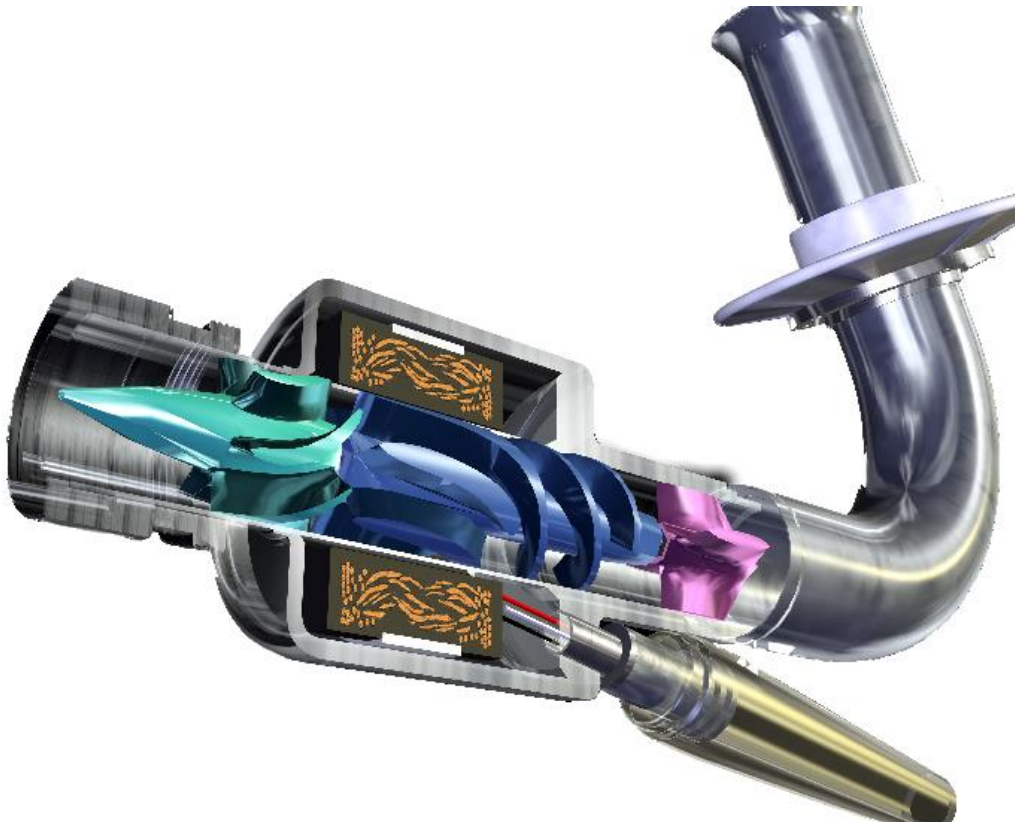
8 CPUs – 90 days



■ Courtesy: Prof Sticht, Bio-Informatics,  
Emil-Fischer Center, FAU

# Why Supercomputing?

## CFD Analysis – blood pump



**Navier stokes equations:**

$$\rho \left( \frac{\partial \mathbf{u}}{\partial t} + \mathbf{u} \cdot \nabla \mathbf{u} - \mathbf{f} \right) - \nabla \cdot \boldsymbol{\sigma} = 0 \quad \text{on } \Omega_t \quad \forall t \in (0, T)$$
$$\nabla \cdot \mathbf{u} = 0 \quad \text{on } \Omega_t \quad \forall t \in (0, T)$$

RWTH Aachen, CATS

## 1. Why supercomputers?

- A few examples of supercomputers
- Top 500 list
- Why is parallelism important?

2. Modern processors
3. Basic optimization techniques for serial code
4. Data access optimization
5. Parallel computers
6. Parallelization and optimization strategies
7. Parallel algorithms
8. Shared-memory programming with OpenMP
9. Distributed-memory programming with MPI
10. Hybrid programming (MPI + OpenMP)
11. Heterogeneous architectures (GPUs, Xeon Phi)
12. Energy efficiency





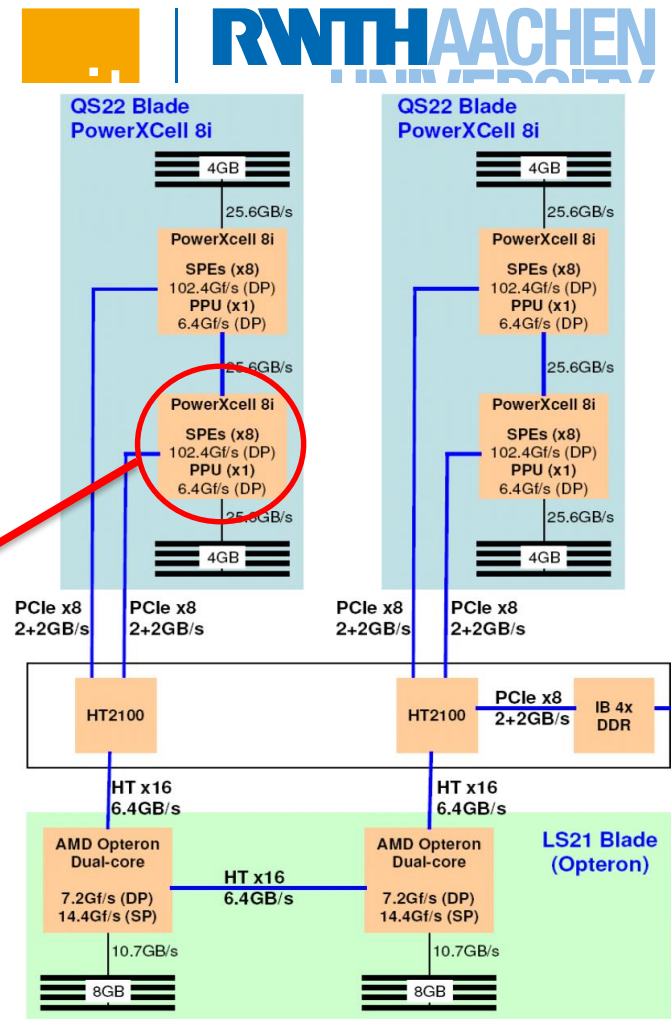
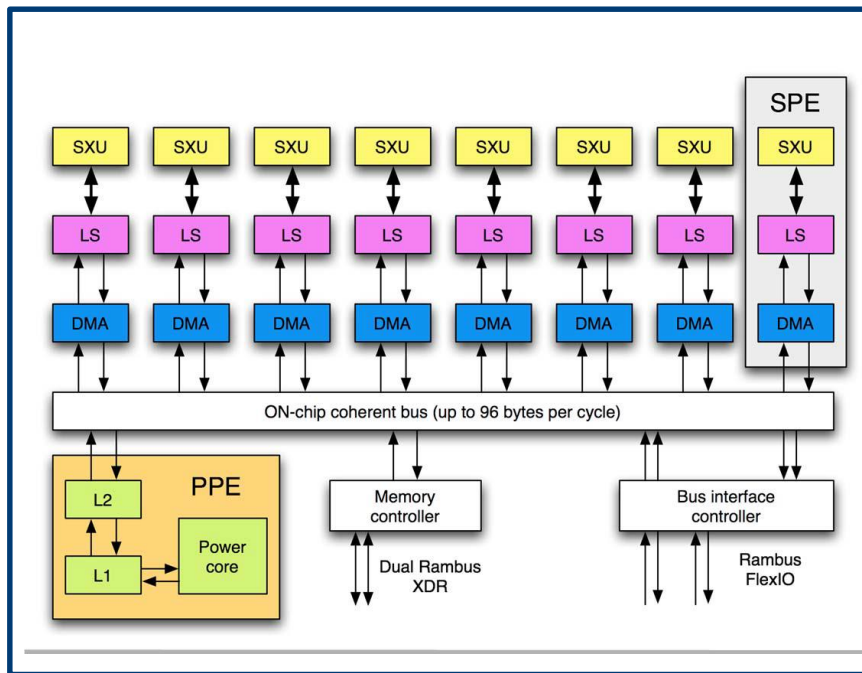
- **MPP System**
- **64-2048 CPUs**
  - DEC Alpha processor (300 MHz)
- **3D-Torus network**
- **UNICOS/mk operating system**
- **No host system needed**
- **Introduced in 1995**
- **Very popular**
  - 50% of the TOP10 fastest supercomputers in June 1997 were 512 CPU T3Es
  - Systems in Jülich and Stuttgart

- **1st computer to surpass the 1 Petaflop barrier**
- **Installed at Los Alamos National Laboratories**
- **Hybrid Architecture**
  - 13,824 AMD Opteron 1.8Ghz cores (6,912 dual-cores)
  - 116,640 IBM 3.2Ghz PowerXCell 8i cores (103,680 SPE, 12,960 PPE) (→Playstation 3 )
- **Shut down in March 2013**



# IBM Roadrunner node (,triblade') structure

- One LS21 Opteron blade + 2x QS22 Cell blades (all in all 400 GF/s)
- LS21: 2x 1.8 GHz dual-core Opterons
- QS22: 2x PowerXCell 8i CPU @ 3.2 GHz
- Single 4x DDR InfiniBand Adapter



Data needs to be transferred explicitly between SPE local store (256KB) and main memory via DMA transfers

## Top 5: K-Computer at RIKEN



- **K** off the Japanese word *kei*, meaning 10 quadrillion ( $10^{16}$ )

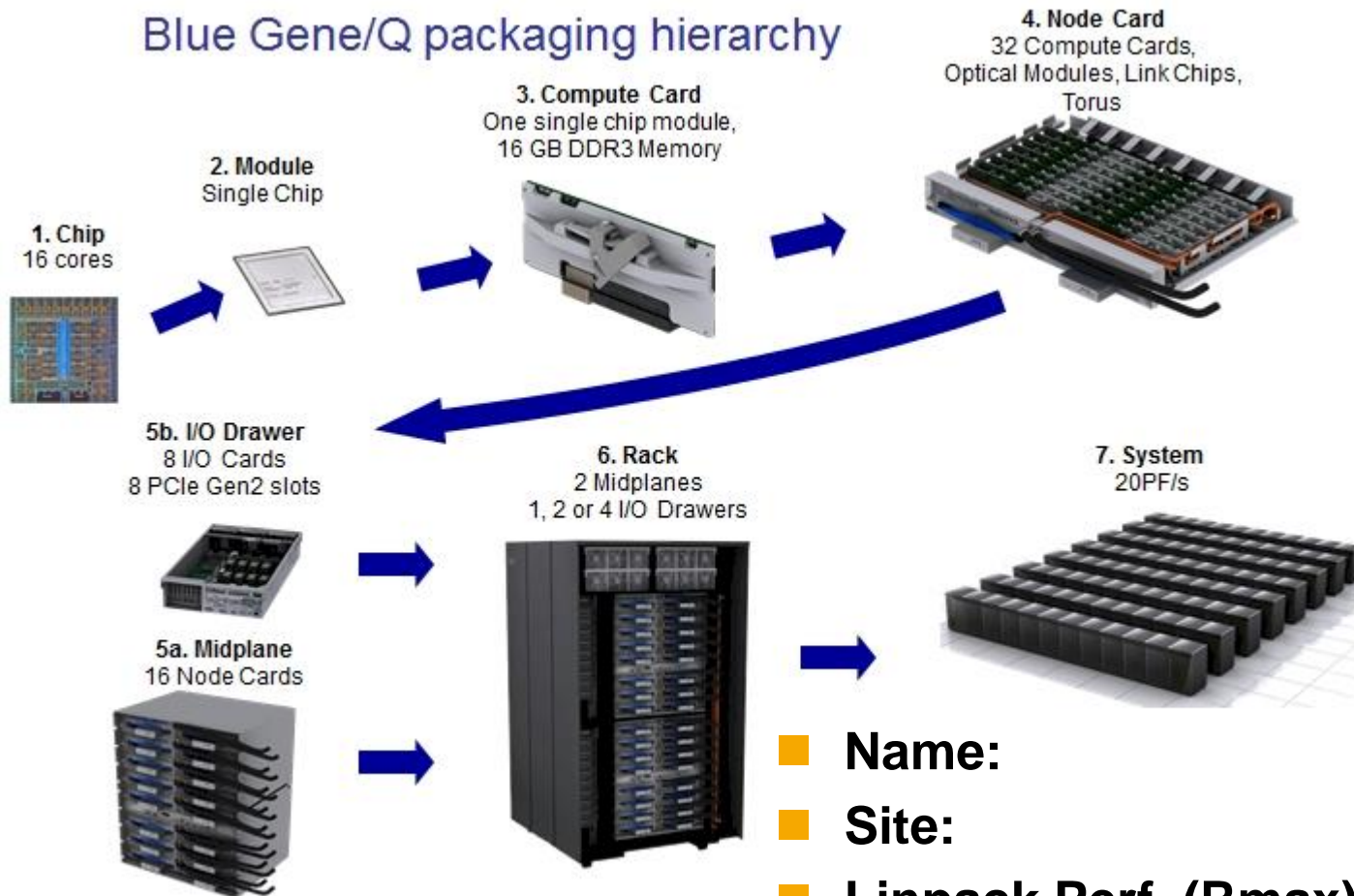
→ Goal that was met 11/2011

■ <b>Name:</b>	<b>K-Computer</b>
■ <b>Vendor:</b>	<b>Fujitsu</b>
■ <b>Site:</b>	<b>RIKEN</b>
■ <b>Linpack Perf. (Rmax):</b>	<b>10,510.0 TFlop/s</b>
■ <b>Power:</b>	<b>12,659.89 kW</b>
■ <b>Cores:</b>	<b>705,024</b>
■ <b>Memory:</b>	<b>1,410,048 GB</b>

# Top 4: IBM Blue Gene/Q at LLNL



## Blue Gene/Q packaging hierarchy



Source: top500.org

■ Name:	IBM Sequoia
■ Site:	DOE/NNSA/LLNL
■ Linpack Perf. (Rmax):	17,173.2 TFlop/s
■ Power:	7890 kW
■ Cores:	1,572,864
■ Memory:	1,572,864 GB



# Top 4: IBM Blue Gene/Q at LLNL

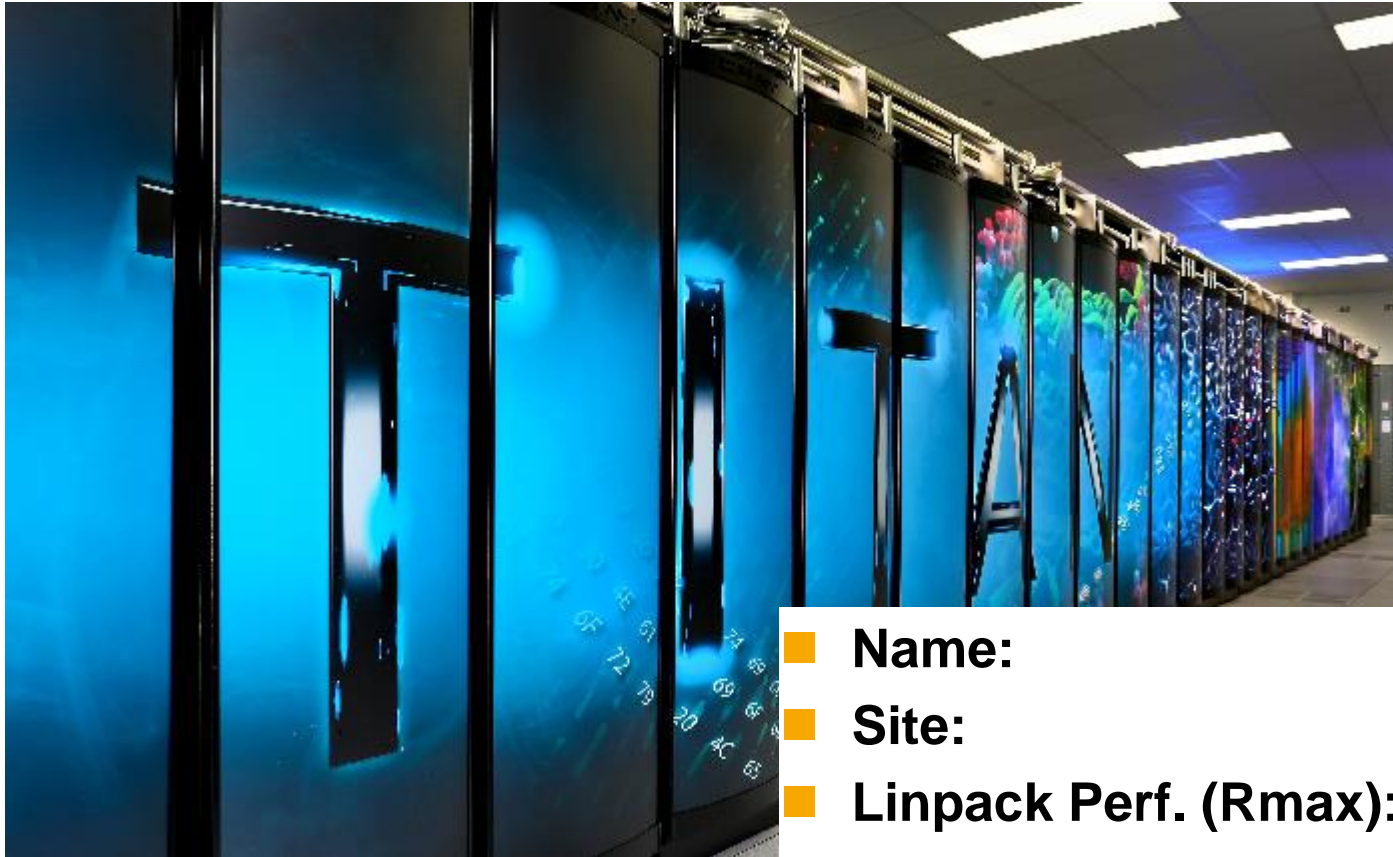
## Node architecture



- 18 IBM PowerPC A2 @ 1.6 GHz, 4-way SMP
- 16 cores are used for computations, 1 runs the OS, IO, MPI, ... the 18. is for redundancy in case one of the other cores is permanently damaged
- 16 GB SDRAM-DDR3
- At least 2 threads on each processor to outbid the machines potential
- 204.8 GFLOps/s / 44W
- 1.47 Billion transistors
- 5D-torus interconnect with 2 GB/s between nodes



## Top 3: Cray TITAN at ORNL (Cray XK7)



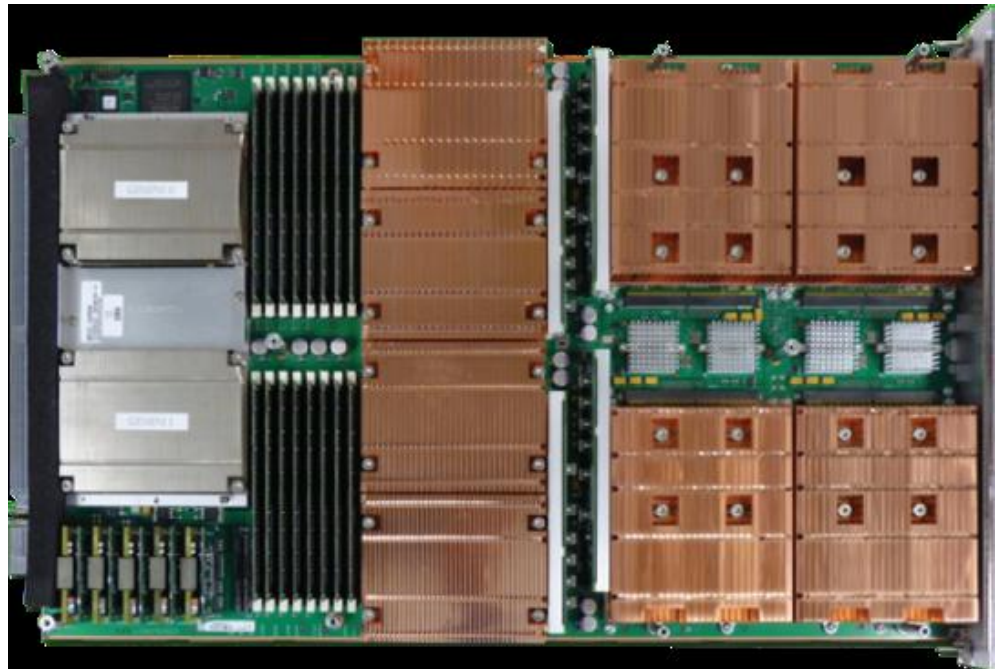
■ <b>Name:</b>	<b>Titan</b>
■ <b>Site:</b>	<b>DOE/SC/ORNL</b>
■ <b>Linpack Perf. (Rmax):</b>	<b>17,590.0 TFlop/s</b>
■ <b>Power:</b>	<b>8,209 kW</b>
■ <b>Cores:</b>	<b>560,640</b>
■ <b>Memory:</b>	<b>710,144 GB</b>
■ <b>Accelerator:</b>	<b>NVIDIA K20x</b>

# Top 3: Cray TITAN at ORNL

## Cray XK7 blade architecture



- **4 nodes per blade**
- **1 node contains**
  - 1 AMD Opteron 6274, 16 cores @ 2.2 GHz, 141 GFLOps/s, 32 GB DDR3
  - 1 NVIDIA Tesla K20x, 14 MP, 1.31 TFLOp/s peak, 6 GB GDDR5
- **1 Cray Gemini 3D-torus interconnect per 2 nodes with 160 GB/s**
- **45 – 54.1 kW**
- **5.8 TFLOp/s per blade**





# Top 2: Tianhe-2/ MilkyWay-2 in China



Source: top500.org

■ <b>Name:</b>	<b>Tianhe-2 (MilkyWay-2)</b>
■ <b>Site:</b>	<b>National Super Computer Center in Guangzhou</b>
■ <b>Linpack Perf. (Rmax):</b>	<b>33,862.7 TFlop/s</b>
■ <b>Power:</b>	<b>17,808 kW</b>
■ <b>Cores:</b>	<b>3,120,000</b>
■ <b>Memory:</b>	<b>1,024,000 GB</b>
■ <b>Accelerator:</b>	<b>Intel Xeon Phi 31S1P</b>

# Top 2: Tianhe-2/ MilkyWay-2 in China

## Node architecture



- **16,000 compute nodes**

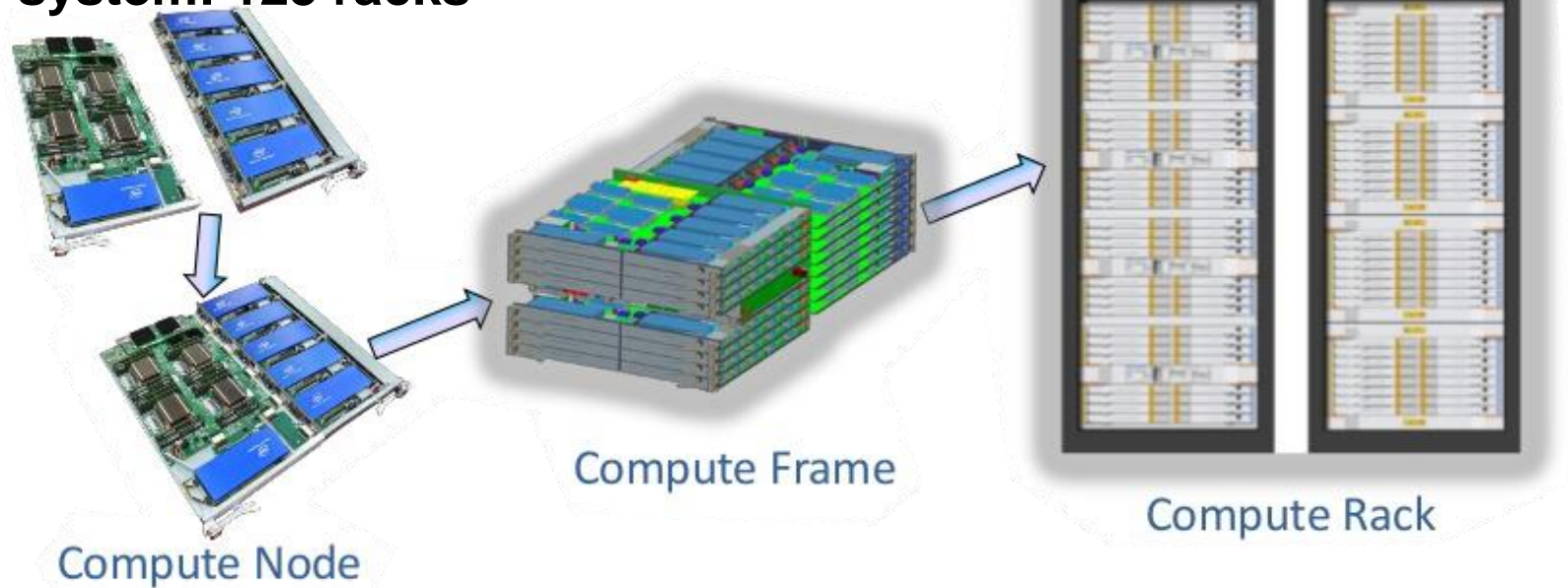
- 2 Intel Ivy Bridge 12-core CPUs @ 2.2 GHz (→ total of 32,000)

- 3 Intel Xeon Phis (→ total of 48,000)

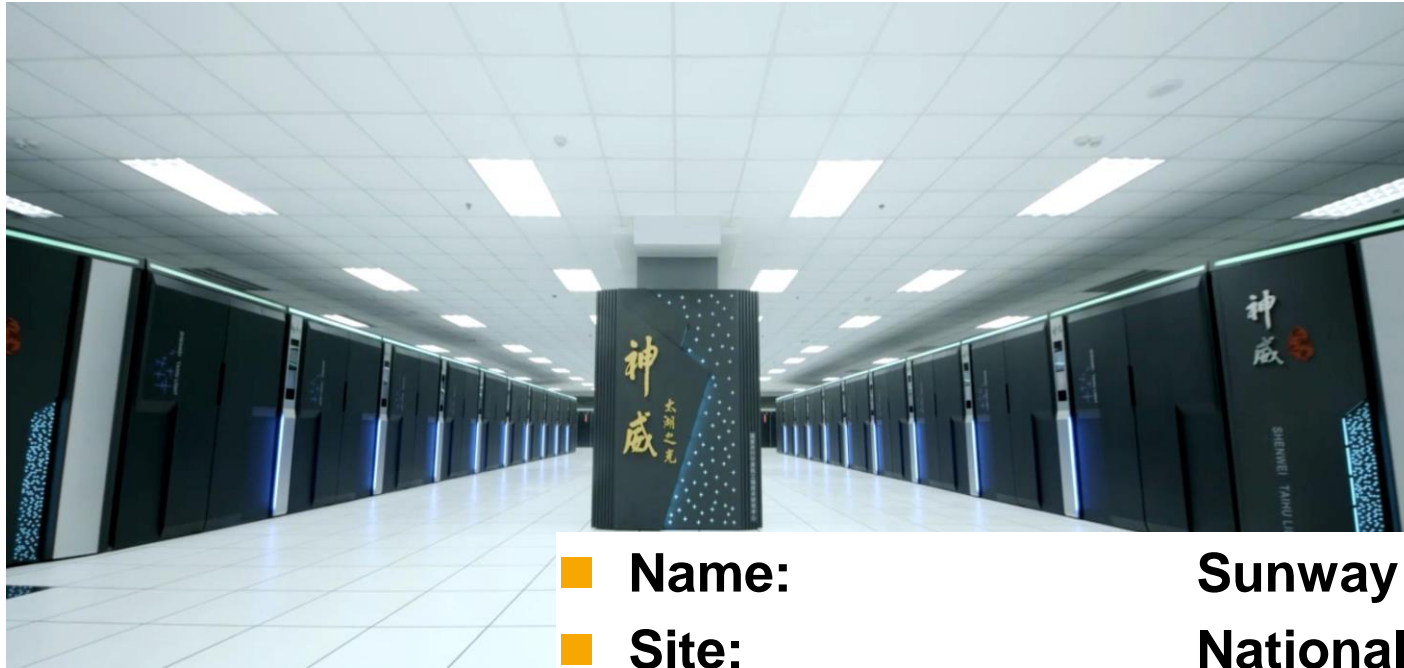
- **Frame: 32 compute nodes**

- **Rack: 4 compute frames**

- **Whole system: 125 racks**



# Top 1: Sunway TaihuLight in China



Source: top500.org

■ <b>Name:</b>	<b>Sunway TaihuLight</b>
■ <b>Site:</b>	<b>National Super Computer Center in Wuxi</b>
■ <b>Linpack Perf. (Rmax):</b>	<b>93,649.6 TFlop/s</b>
■ <b>Power:</b>	<b>15,371 kW</b>
■ <b>Cores:</b>	<b>10,649,600</b>
■ <b>Memory:</b>	<b>1,310,720 GB</b>

Source: top500.org

# Top 1: Sunway TaihuLight in China

## Node architecture

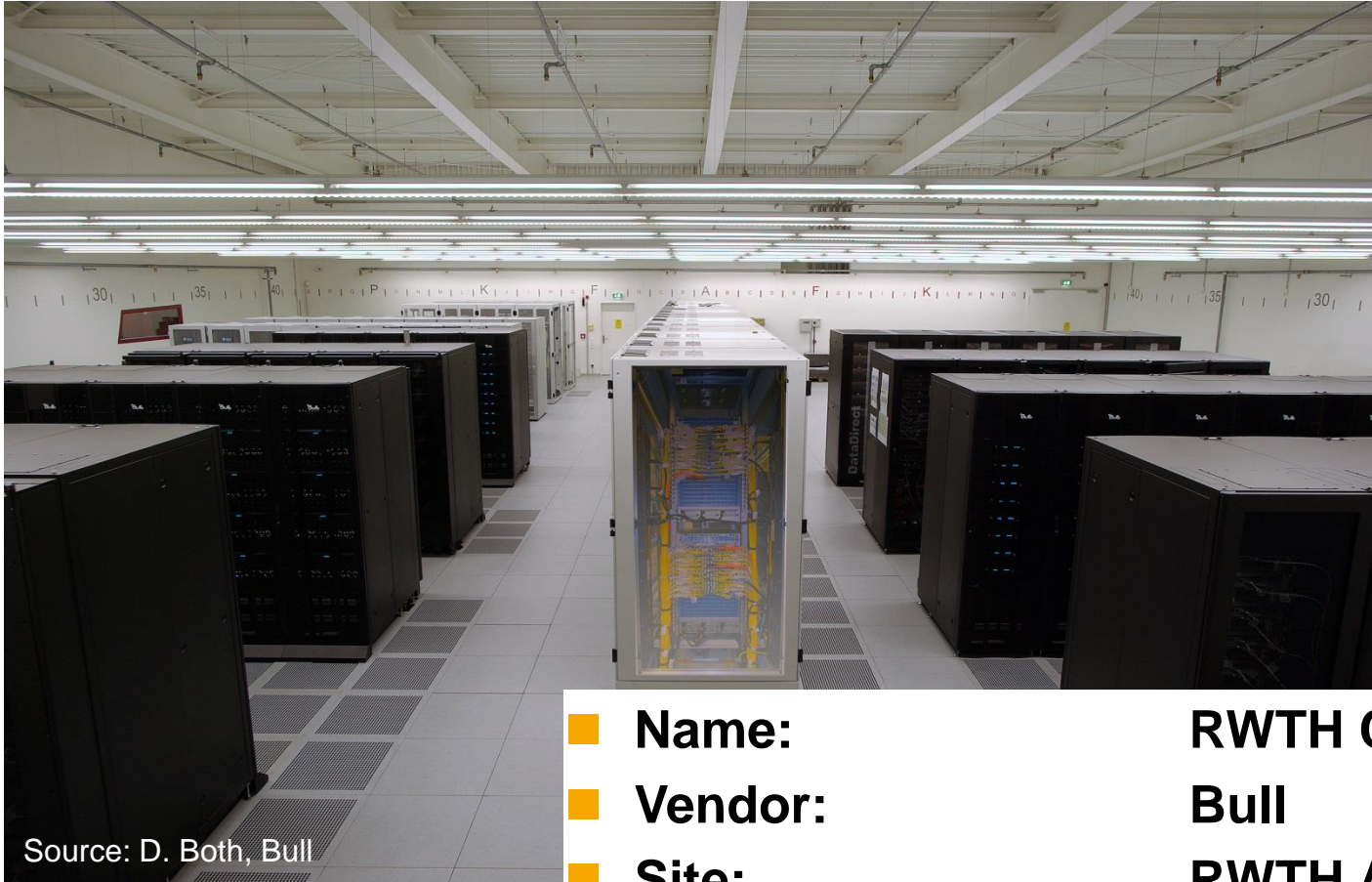


- Chinese “homegrown” hardware and software
- **SW26010 processor**

- Developed by Shanghai IC Design Center
- 260 cores @ 1.45 GHz
- 3.06 TFlop/s double precision peak performance



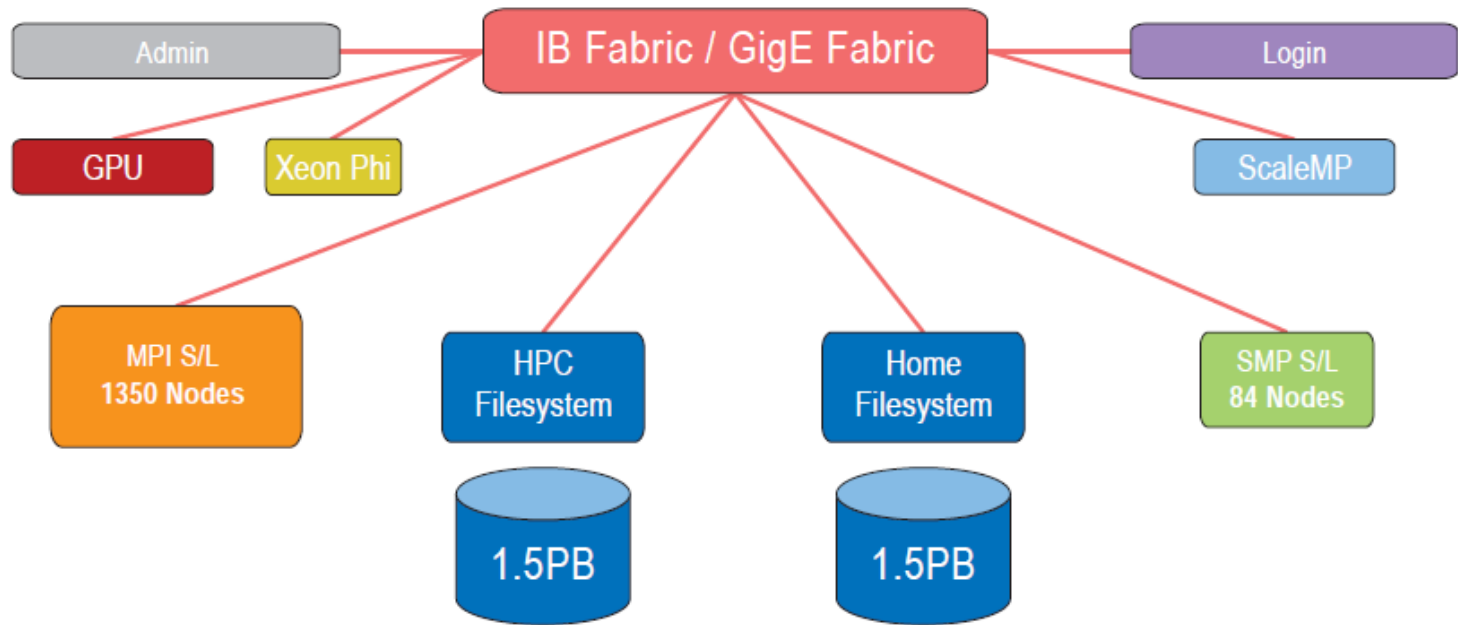
- **40 cabinets à 4 supernodes à 256 nodes → 40,960 nodes**
- **Nodes are connected by PCI-E 3.0 with 16 GB/s**
- **15.37 MW during LINPACK Benchmark (→ 6 Gflops/W)**
- **Software**
  - Sunway Raise OS 2.05 based on Linux
  - Own software stack: compilers, vectorization tools, math libraries, etc.



- **Name:** RWTH Compute Cluster
- **Vendor:** Bull
- **Site:** RWTH Aachen University
- **Linpack Perf. (Rmax):** 219.84 TFlop/s
- **Cores:** 25,448

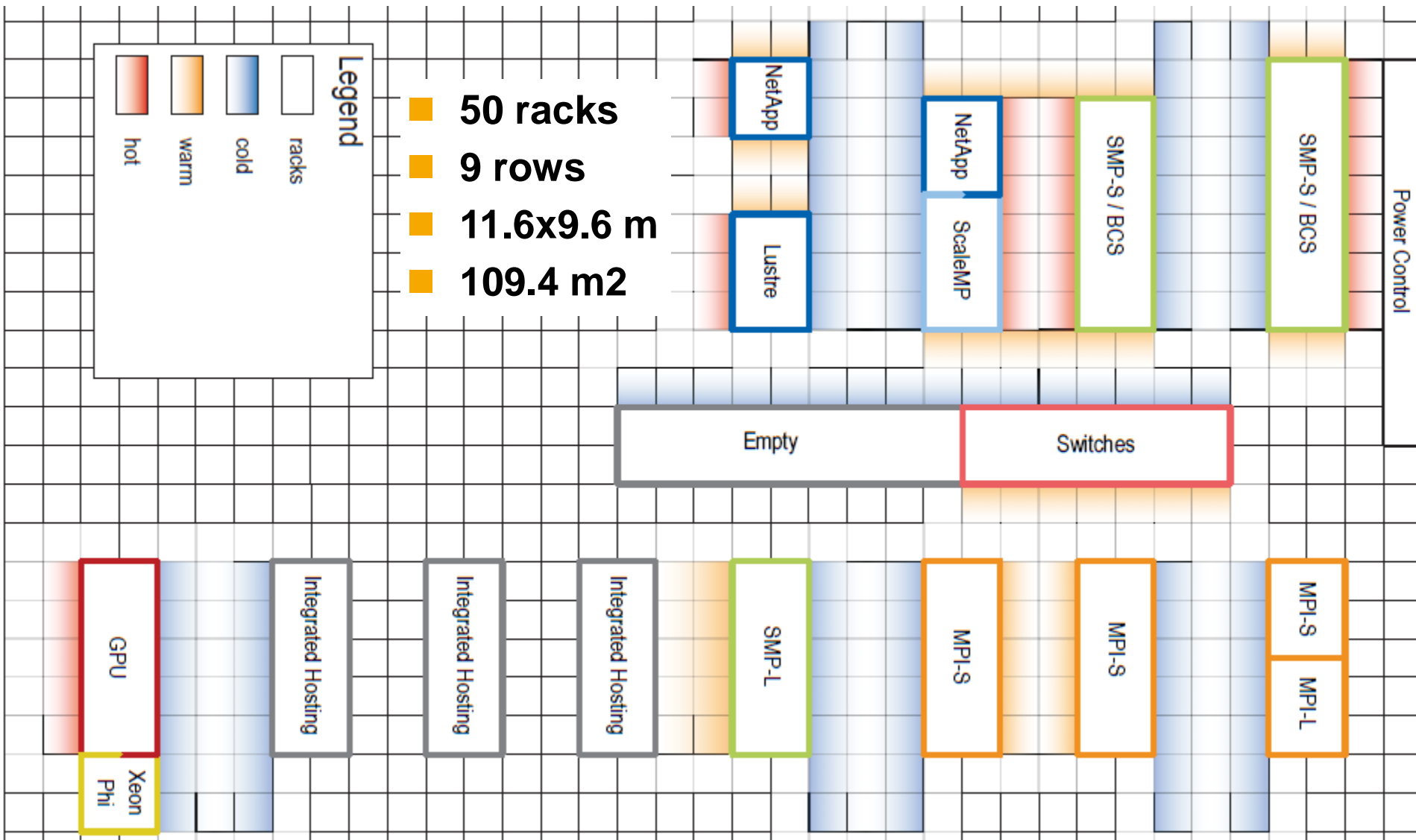


# RWTH Compute Cluster: Overview



Group	# Nodes	Sum TFLOP	Sum Memory	# Procs per Node	Memory per Node
<b>MPI-Small</b>	1098	161	26 TB	2 x Westmere-EP (3.06 GHz, 6 Cores, 12 Threads)	24 GB
<b>MPI-Large</b>	252	37	24 TB		96 GB
<b>SMP-Small</b>	135	69	17 TB	4x Nehalem-EX (2.0 GHz, 8 Cores, 16 Threads) BCS: 4, vSMP: 16	128 GB
<b>SMP-Large</b>	36	18	18 TB		512 GB
<b>ScaleMP-vSMP</b>	8 gekoppelt	4	4 TB		4 TB gekoppelt
<b>Xeon Phi</b>	6 w/ 2 Phis each	12 + 1.5	96 + 192 GB	2x SandyBridge (2 GHz, 8 cores, 16 threads) 2x Xeon Phi (1.053 GHz, 60 cores)	32 + 16 GB
<b>GPU</b>	24 w/ 2 GPUs each	28 + 7	336 + 672 GB	2 x Westmere (2.67 GHz, 6 Cores, 12 Threads) 2x Quadro 6000 (448 cores)	24 + 12 GB

# RWTH Compute Cluster: Floorplan



# New RWTH HPC System

## Cluster Aix-la-Chapelle (CLAIX)



- **Vendor: NEC**

- **Intel Omni-Path Network**

- **>600 MPI nodes**

  - 2-socket (24 cores) Intel Broadwell

  - 128 GB main memory

- **8 SMP nodes**

  - 8-socket (144 cores) Intel Broadwell

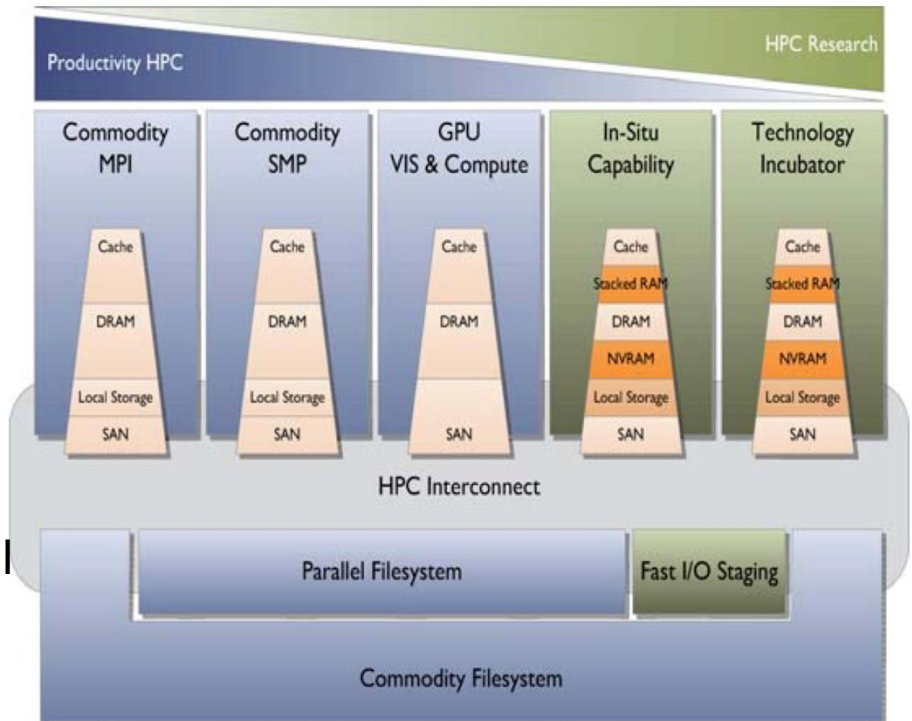
  - 1024 GB main memory

- **Some specialized nodes**

  - GPU nodes for visualization and general-purpose computing (GPGPU)

  - Many Integrated Core (MIC) architecture, Intel Xeon Phi (Knights Landing)

- **Commissioning: November 2016**





## 1. Why supercomputers?

- A few examples of supercomputers
- **Top 500 list**
- Why is parallelism important?

2. Modern processors
3. Basic optimization techniques for serial code
4. Data access optimization
5. Parallel computers
6. Parallelization and optimization strategies
7. Parallel algorithms
8. Shared-memory programming with OpenMP
9. Distributed-memory programming with MPI
10. Hybrid programming (MPI + OpenMP)
11. Heterogeneous architectures (GPUs, Xeon Phi)
12. Energy efficiency

- List of the fastest 500 supercomputers in the world
- Ranking set up after the Rmax value result of a machine running the LINPACK Benchmark
  - Solve a large system of linear equations:  $A\vec{x} = \vec{b}$
  - Result is measured in Flop/s (Floating Point Operations per Second)
  - Floating Point operations: double-precision (64bit) add & mult. operations
  - Originally LINPACK was a library for the solution of LES, superseded by LAPACK
- Published twice a year (ISC in Germany – Jun., SC in USA – Nov.)
- Established in 1993 (CM5/1024): 60 GFlop/s
- June 2016 (Sunway TaihuLight): 93,014,600 GFlop/s
- $10^6$  MFlop/s;  $10^9$  GFlop/s;  $10^{12}$  TFlop/s;  $10^{15}$  PFlop/s;  $10^{18}$  EFlop/s

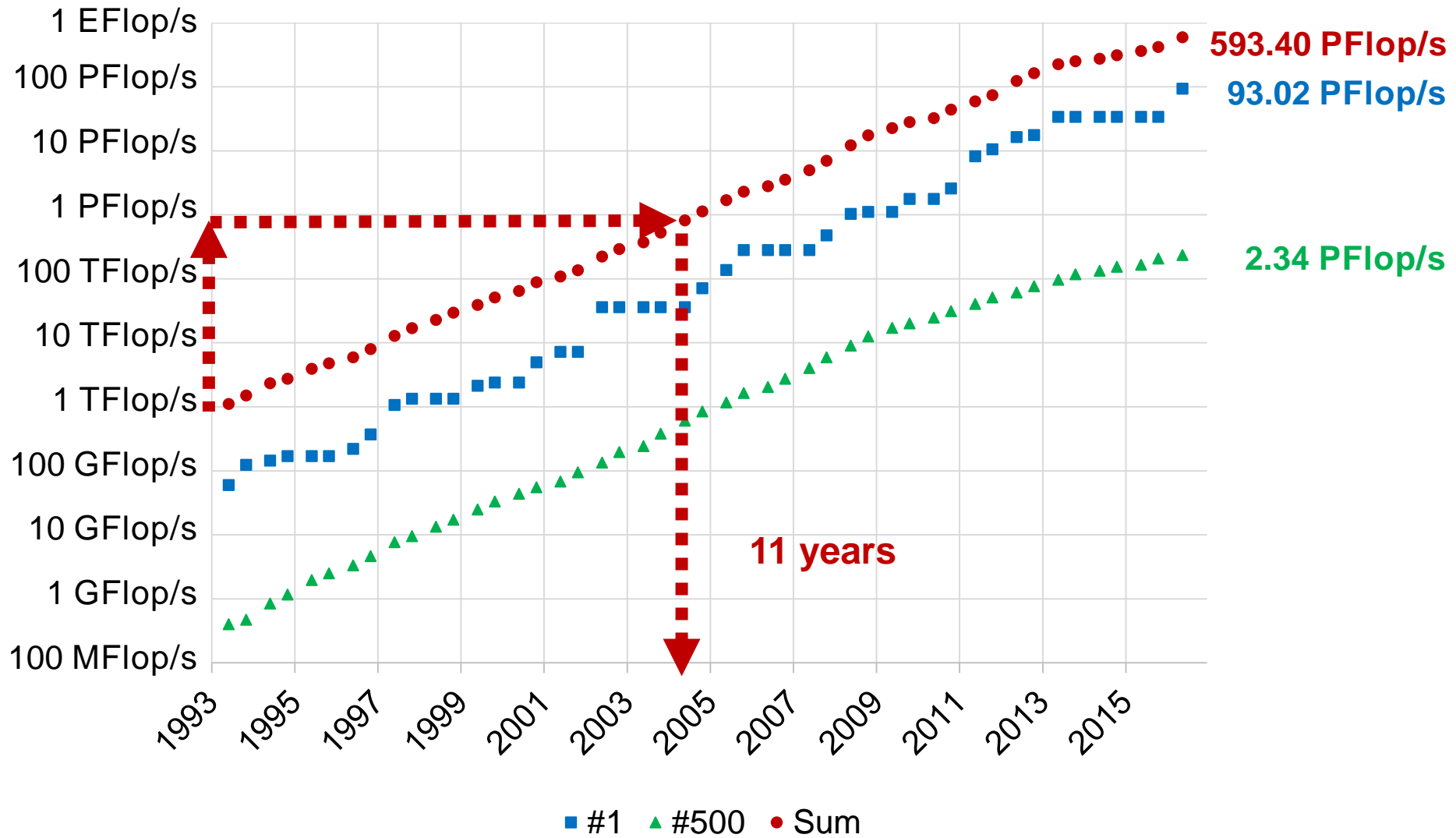
# TOP500 as of November 2013



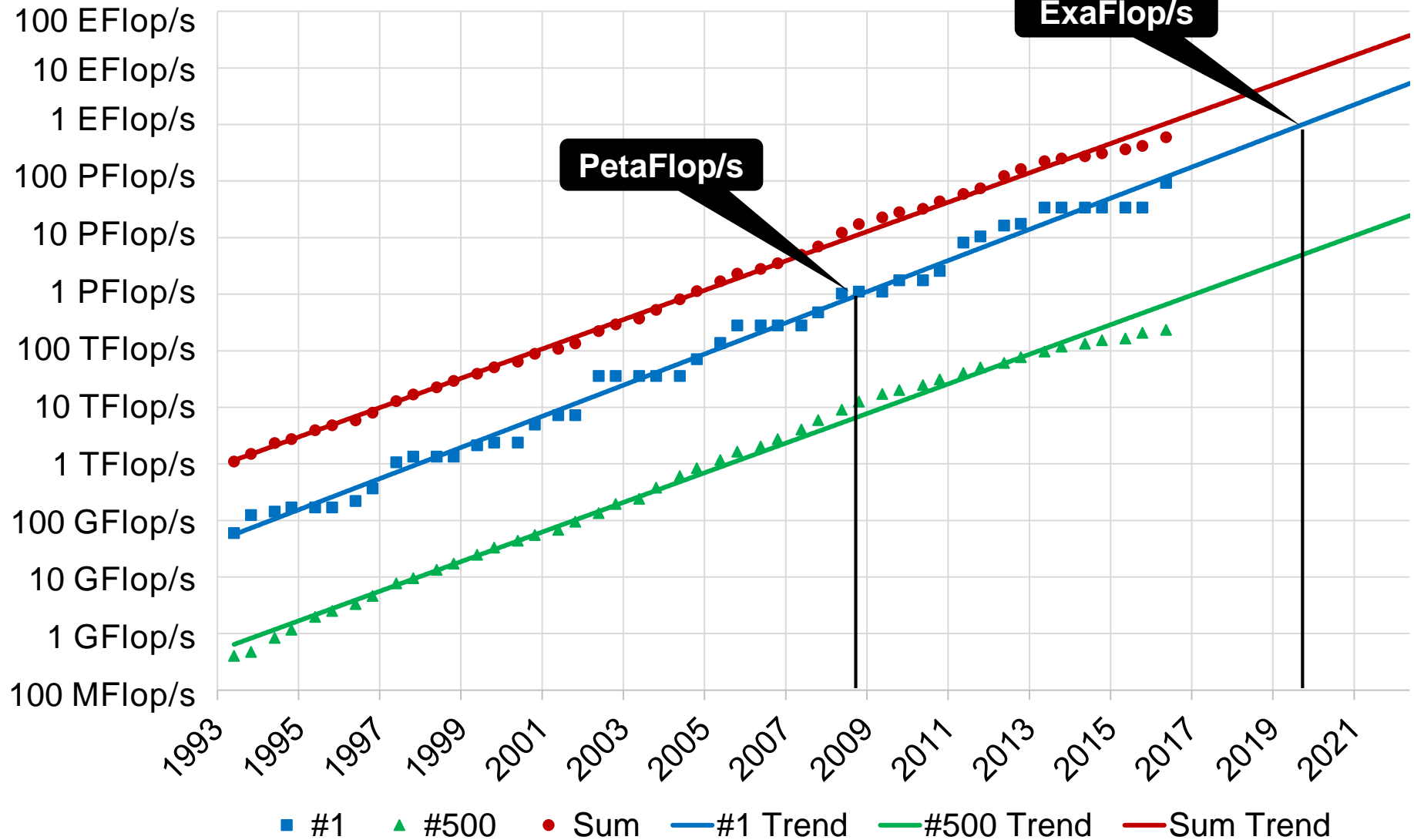
#	Name	Computer	Site	Manu- facturer	#Cores	Rmax [Tflop/s]	Rpeak [Tflop/s]	Power [kW]
1	Sunway TaihuLight	Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway	National Super Computer Center in Wuxi	NRCPC	10,649,600	93,014.6	125,435.9	15,371
2	Tianhe-2 (MilkyWay-2)	TH-IVB-FEP Cluster, Intel Xeon E5-2692 12C 2.200GHz, TH Express-2, Intel Xeon Phi 31S1P	National Super Computer Center in Guangzhou	NUDT	3,120,000	33,862.7	54,902.4	17,808
3	Titan	Cray XK7 , Opteron 6274 16C 2.200GHz, Cray Gemini interconnect, NVIDIA K20x	DOE/SC/Oak Ridge National Laboratory	Cray Inc.	560,640	17,590.0	27,112.5	8,209
4	Sequoia	BlueGene/Q, Power BQC 16C 1.60 GHz, Custom	DOE/NNSA/LLNL	IBM	1,572,864	17,173.2	20,132.7	7,890
5	K computer	SPARC64 VIIIfx 2.0GHz, Tofu interconnect	RIKEN Advanced Institute for Computational Science (AICS)	Fujitsu	705,024	10,510.0	11,280.4	12,660
6	Mira	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	DOE/SC/Argonne National Laboratory	IBM	786,432	8,586.6	10,066.3	3,945
7	Trinity	Cray XC40, Xeon E5-2698v3 16C 2.3GHz, Aries interconnect	DOE/NNSA/LANL/SNL	Cray Inc.	301,056	8,100.9	11,078.9	
8	Piz Daint	Cray XC30, Xeon E5-2670 8C 2.600GHz, Aries interconnect , NVIDIA K20x	Swiss National Supercomputing Centre (CSCS)	Cray Inc.	115,984	6,271.0	7,788.9	2,325
9	Hazel Hen	Cray XC40, Xeon E5-2680v3 12C 2.5GHz, Aries interconnect	HLRS Stuttgart	Cray Inc.	185,088	5,640.2	7,403.5	
10	Shaheen II	Cray XC40, Xeon E5-2698v3 16C 2.3GHz, Aries interconnect	King Abdullah University of Science and Technology	Cray Inc.	196,608	5,537.0	7,235.2	2,834

**1000KW = 1MW  $\approx$  1.2 million € /year**

# Performance development in TOP500



# Projected performance development in TOP500



## 1. Why supercomputers?

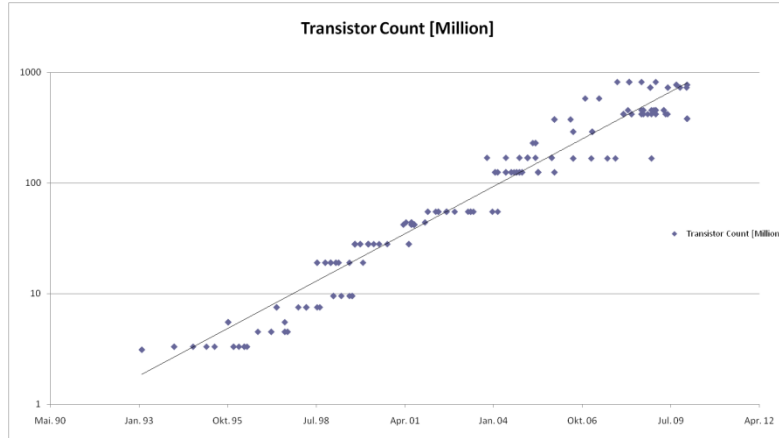
- A few examples of supercomputers
- Top 500 list
- **Why is parallelism important?**

2. Modern processors
3. Basic optimization techniques for serial code
4. Data access optimization
5. Parallelization and optimization strategies
6. Parallel algorithms
7. Shared-memory programming with OpenMP
8. Distributed-memory programming with MPI
9. Hybrid programming (MPI + OpenMP)
10. Heterogeneous architectures (GPUs, Xeon Phi)
11. Parallel computers
12. Energy efficiency

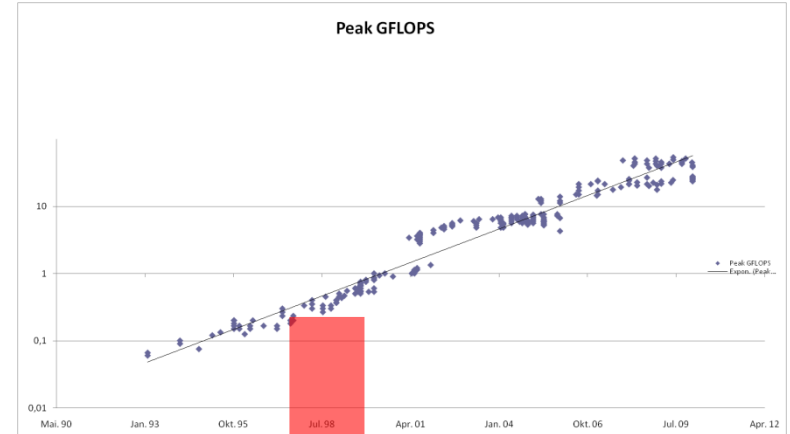
# Driving force behind TOP500 development



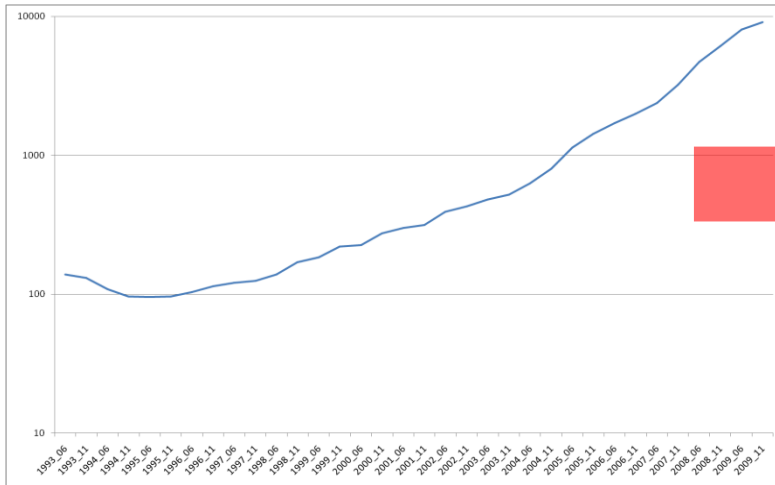
## Moore's law



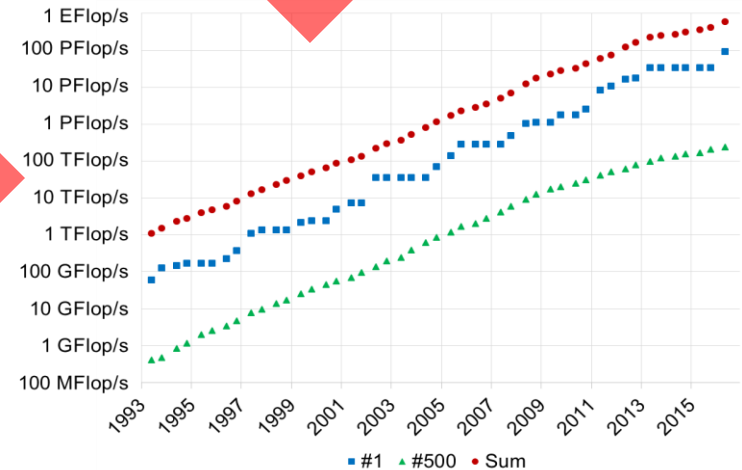
## Processor performance

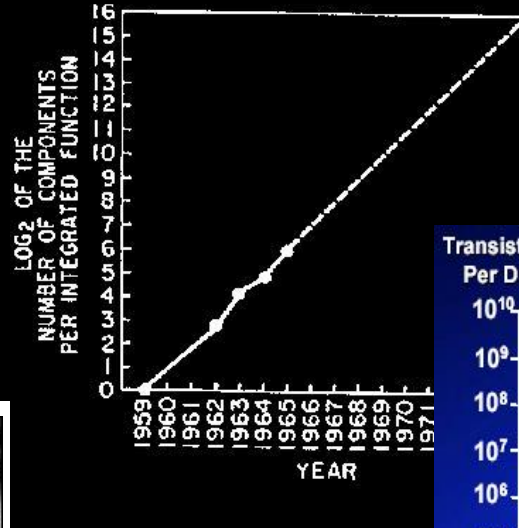
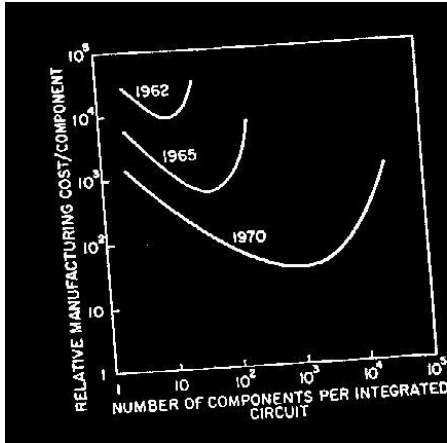


## Parallelism (#cores)



## System performance

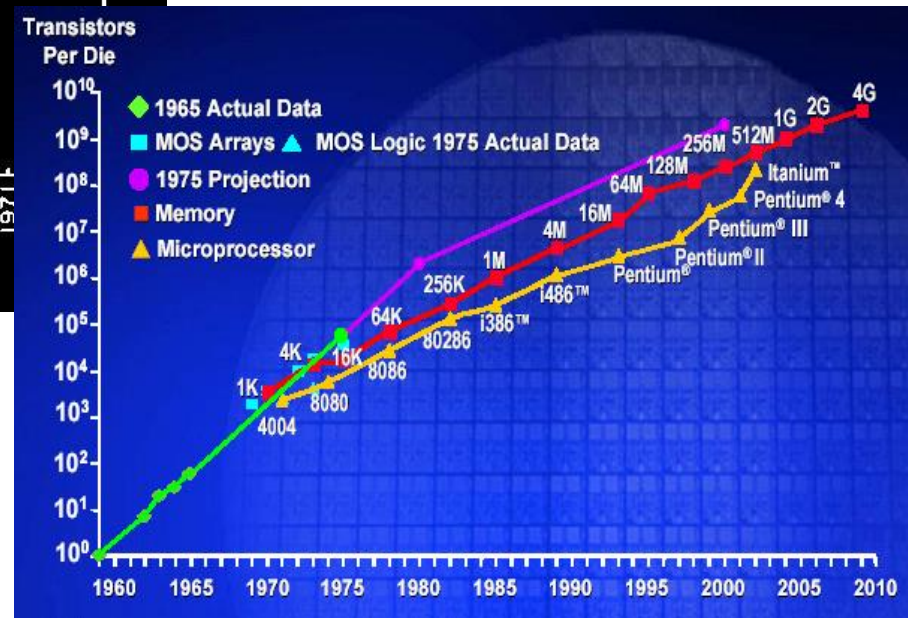
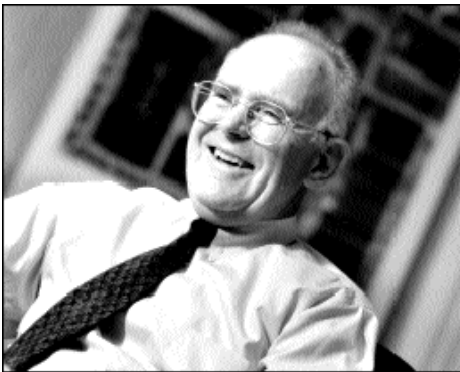




***“Cramming More Components onto Integrated Circuits”***

Gordon Moore, Electronics, 1965

([ftp://download.intel.com/museum/Moores\\_Law/Articles-Press\\_Releases/Gordon\\_Moore\\_1965\\_Article.pdf](ftp://download.intel.com/museum/Moores_Law/Articles-Press_Releases/Gordon_Moore_1965_Article.pdf))



■ # on transistors / cost-effective integrated circuit double every N months ( $12 \leq N \leq 24$ )



- **Only predicts the number of transistors!**
  - Number of transistors per chip is  $1.59^{\text{year}-1959}$
  - Now slope is less; but we should see 10 – 100x or more growth  
(65 nm – sub 10 nm)
- **But why should twice the transistors offer twice the performance?**

# Moore's Law at work: History of X86 CPUs

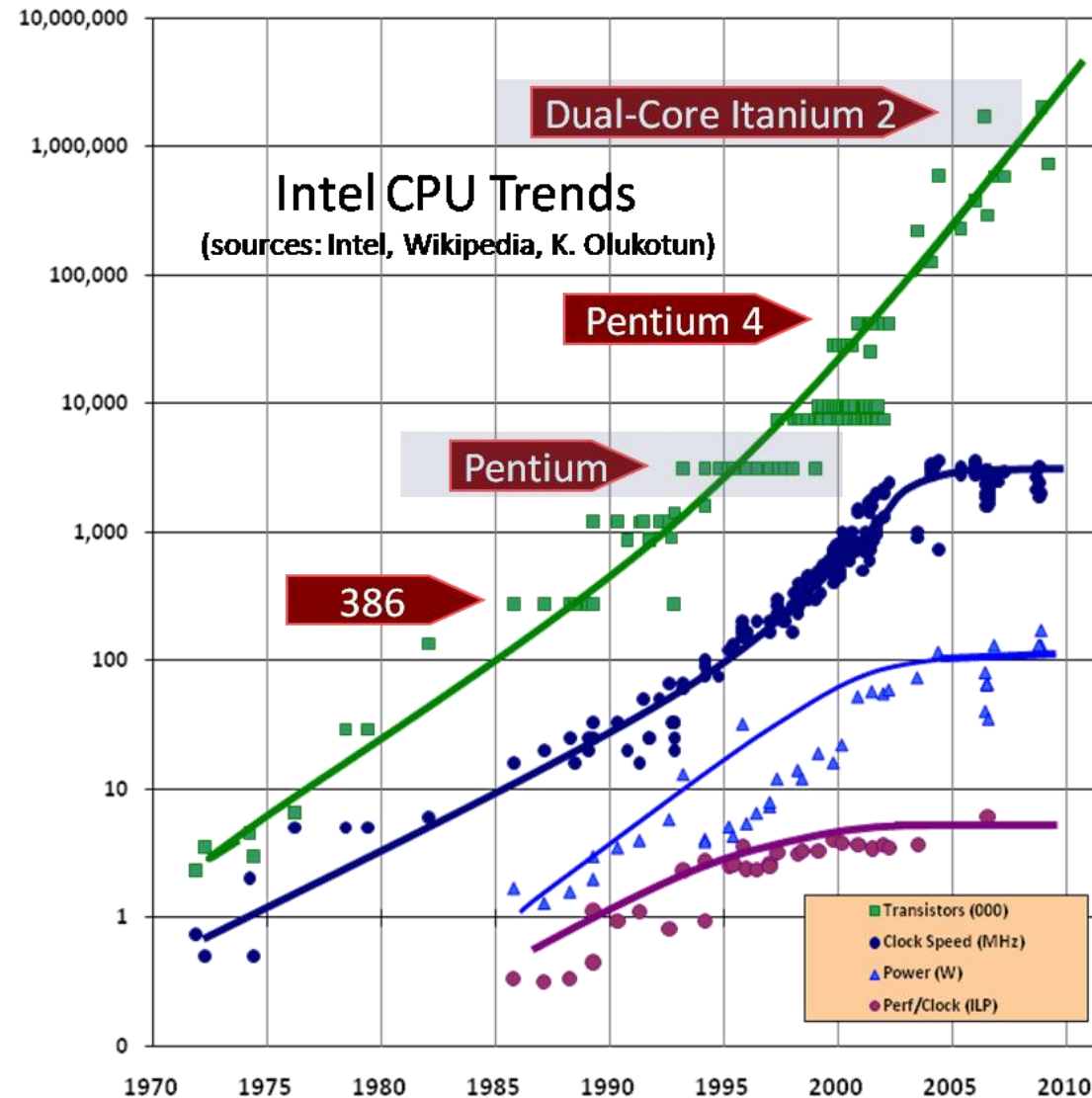
clock speed  
increases (not  
anymore)



CPU	Year	BW	#Transistors	Clock	Structure	L1 / L2 / L3
			#transistors increase			physical size decreases
4004	1971	4	2,300	740 kHz	10 micro	
8008	1972	8	3,500	500 kHz	10 micro	
8086	1978	16	29,000	10 Mhz	3 micro	
80286	1982	16	134,000	25 MHz	1.5 micro	
80386	1985	32	275,000	33 Mhz	1 micro	
80486	1989	32	1,200,000	50 MHz	0.8 micro	8K
Pentium I	1994	32	3,100,000	66 MHz	0.8 micro	8K
Pentium II	1997	32	7,500,000	300 MHz	0.35 micro	16K/512K*
Pentium III	1999	32	9,500,000	600 MHz	0.25 micro	16K/512K*
Pentium IV	2000	32	42,000,000	1.5 GHz	0.18 micro	8K/256K
P IV F	2005	64		2.8- 3.8 GHz	90 nm	16K/2MB
Core i7	2008	64	781,000,000	3.2 GHz	45 nm	32K/256K/8MB
Westmere-EP	2010	64	1,170,000,000	3.46 GHz	32 nm	32K/256K/12MB
Ivy Bridge-EP	2013	64	2,890,000,000	~3.3 GHz	22 nm	32K/256K/10-30MB

## ■ Investigating the relationship of the amount of transistors, clock speed and power consumption

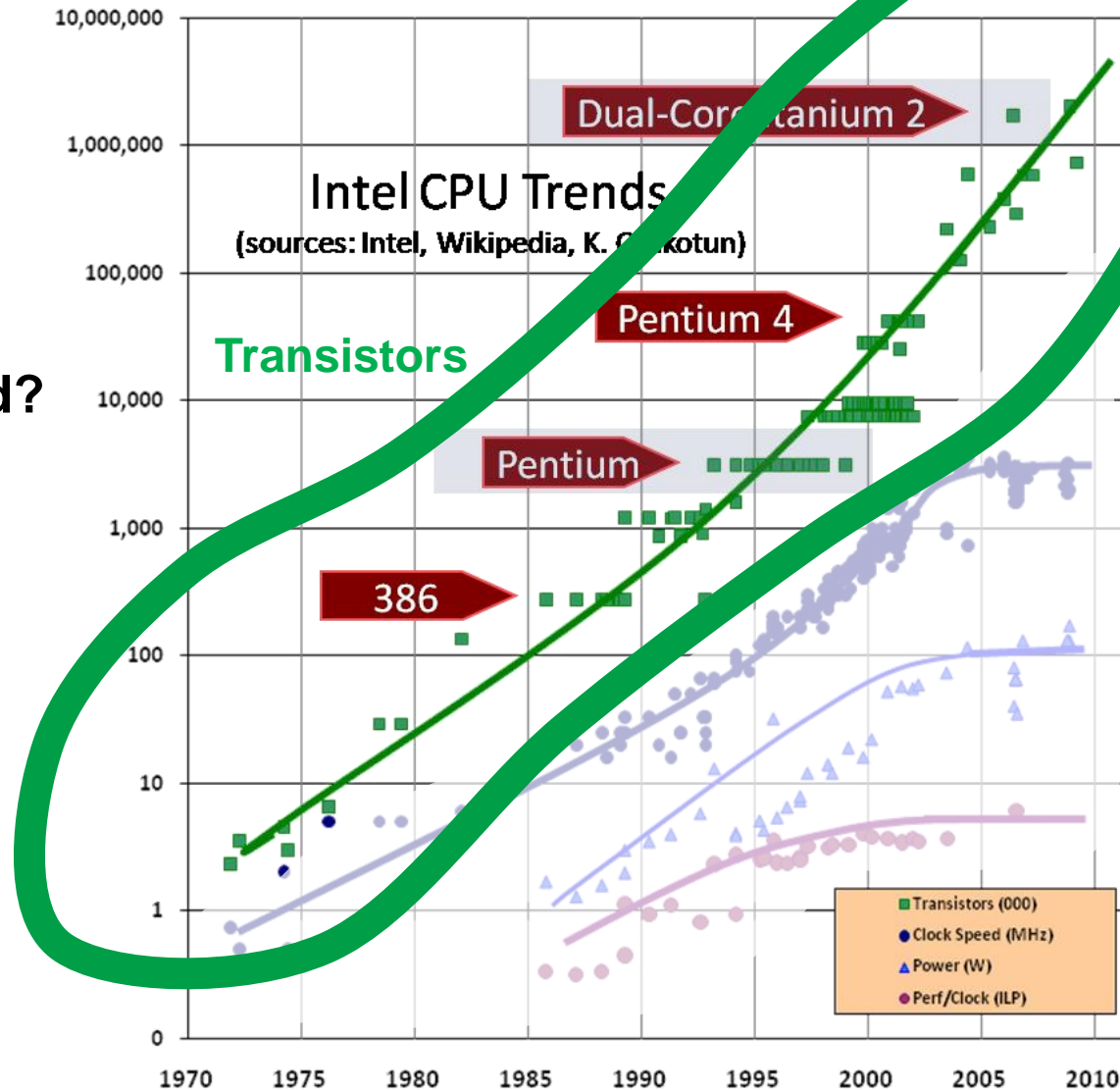
→ See following slides for details



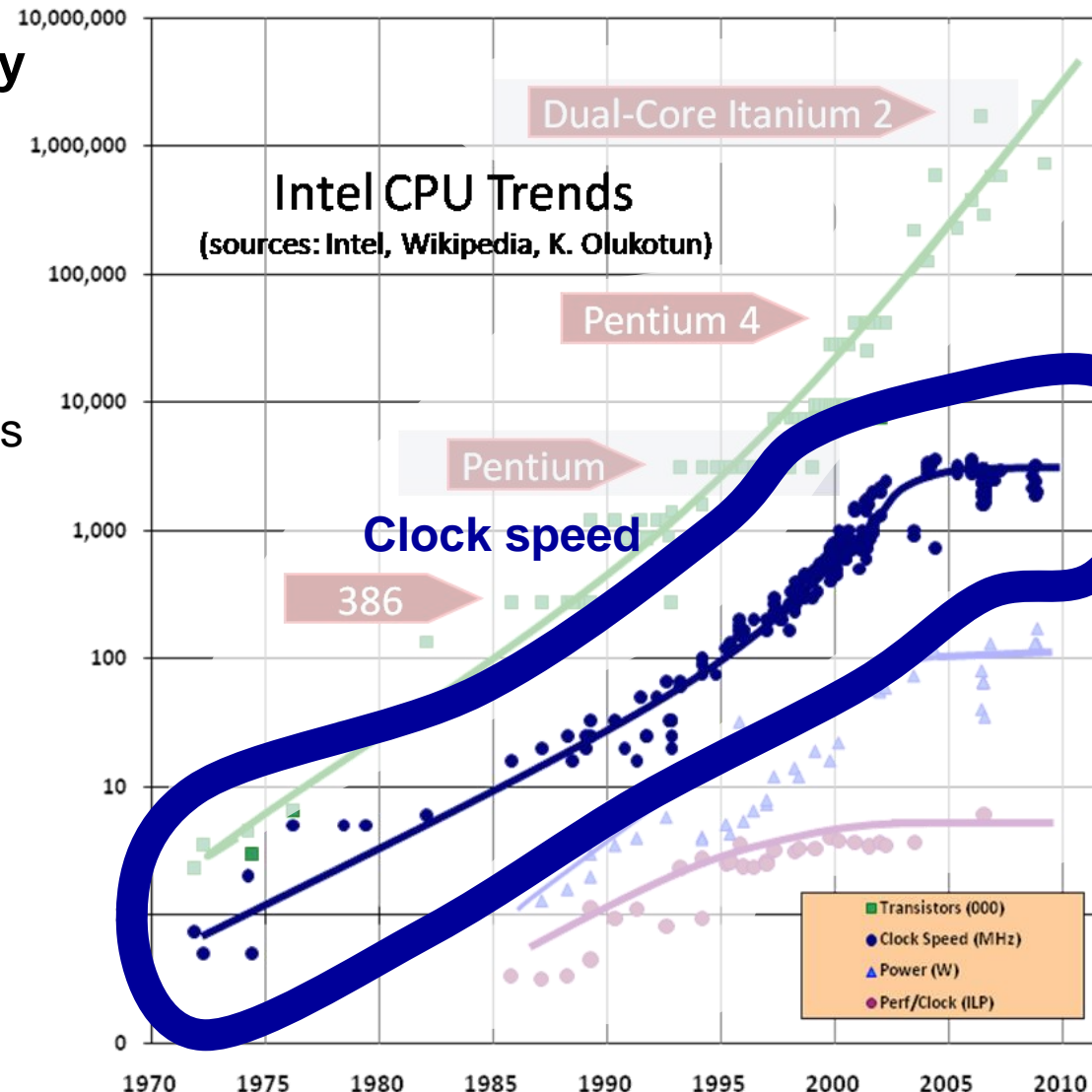
## ■ Moore's law continues

→ Transistor count is still growing exponentially

## ■ How can that be achieved?



- Increase of clock frequency is one solution, but:
- Exponential clock rate growth has ended
  - Stagnation since early 2000s
- Why?



## Reason: power consumption

$$\rightarrow P \sim f^3$$

P = power, f = frequency

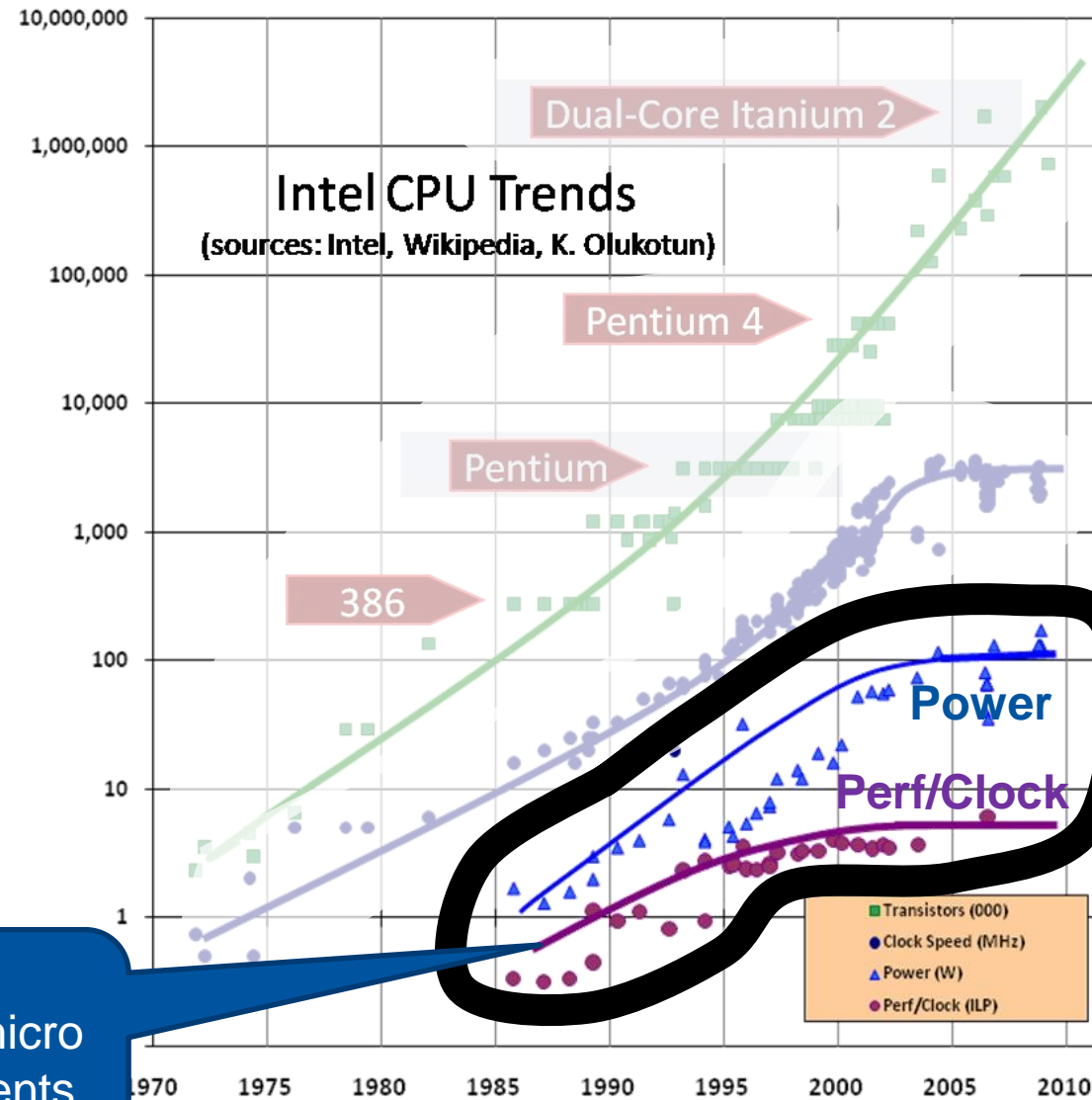
→ High heat dissipation

## Power consumption may not be increased

→ Clock frequency stagnates

→ Performance per cycle  
stagnates

Parallelism cannot be  
exploited further with micro  
architecture improvements.  
See later lecture.

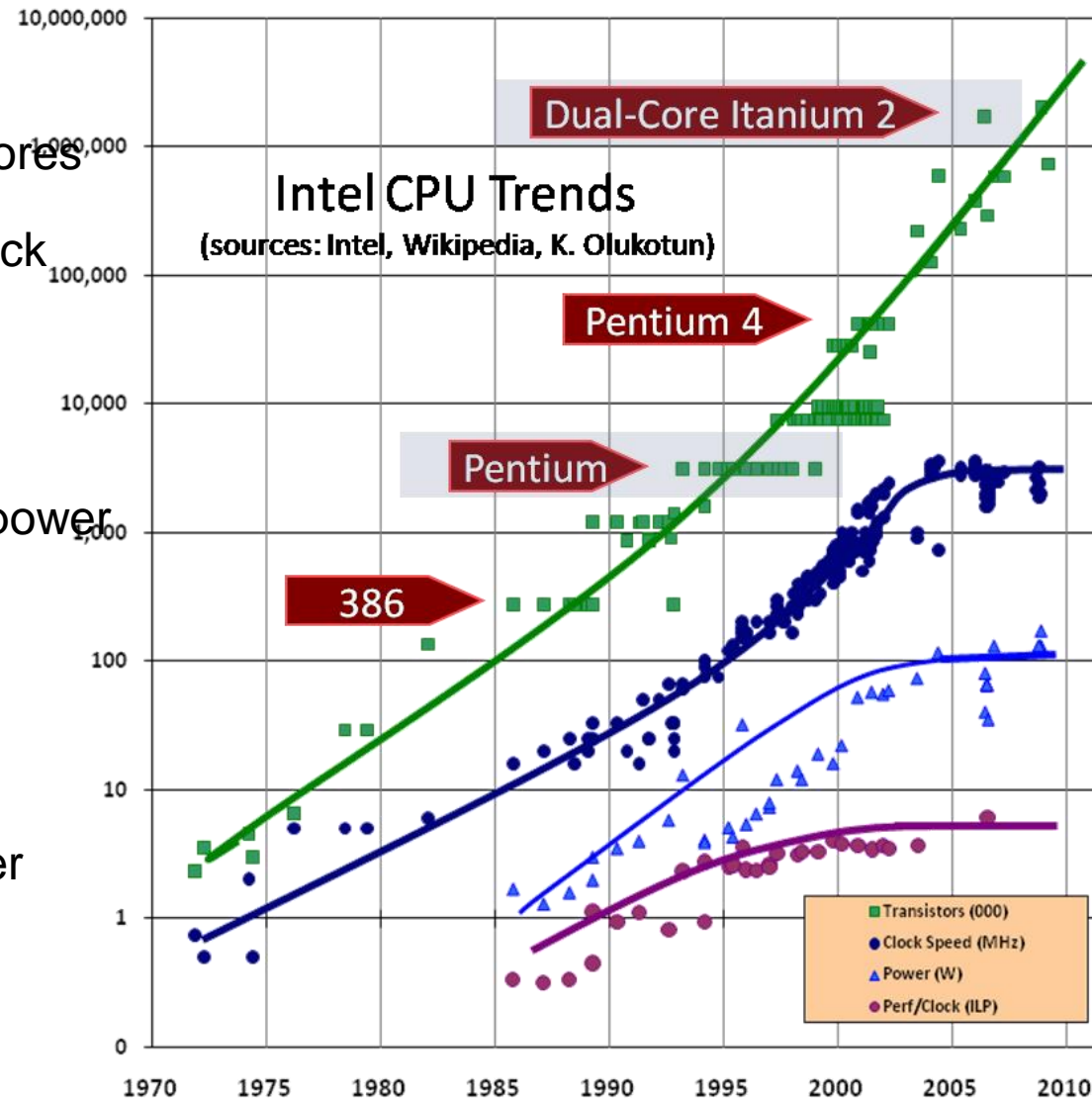


## Multi-core chips

- Transistors depict several cores
- Cores have same/ lower clock frequency
- Linear increase of power consumption instead of 3<sup>rd</sup> power

## Moore's Law continues

- Multi-core era
- Every 18 months the number of cores per chip doubles instead of clock frequency



- **What is a supercomputer?**
- **What are supercomputers used for?**
- **What are recent supercomputers?**
- **What can we read from the performance development as measured in the TOP500?**
- **What does Moore's law tell you? Is it still valid?**
- **Why do we have multi-core architectures today?**