

中国科学院国家天文台

National Astronomical Observatories, CAS



RECRUITMENT  
PROGRAM OF GLOBAL EXPERTS

UNIVERSITÄT  
HEIDELBERG  
Zukunft. Seit 1386.



Introduction to GPU  
Accelerated Computing:  
1. History of Computer Architecture  
Many-Core, GPU, and other ideas...

University

Rainer Spurzem

Astronomisches Rechen-Inst., ZAH, Univ. of Heidelberg, Germany  
National Astronomical Observatories (NAOC), Chinese Academy of Sciences  
Kavli Institute for Astronomy and Astrophysics (KIAA), Peking University

the SILK ROAD PROJECT at NAOC/KIAA

丝绸之路计划

[spurzem@ari.uni-heidelberg.de](mailto:spurzem@ari.uni-heidelberg.de)  
<http://silkroad.bao.ac.cn>



北京大學  
PEKING UNIVERSITY



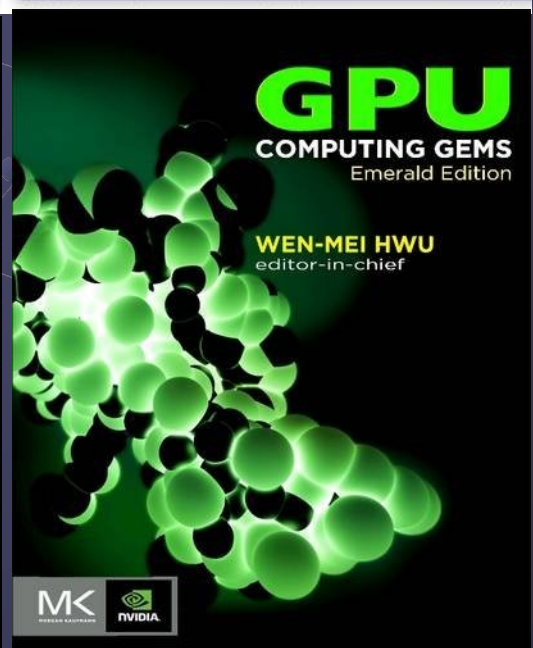
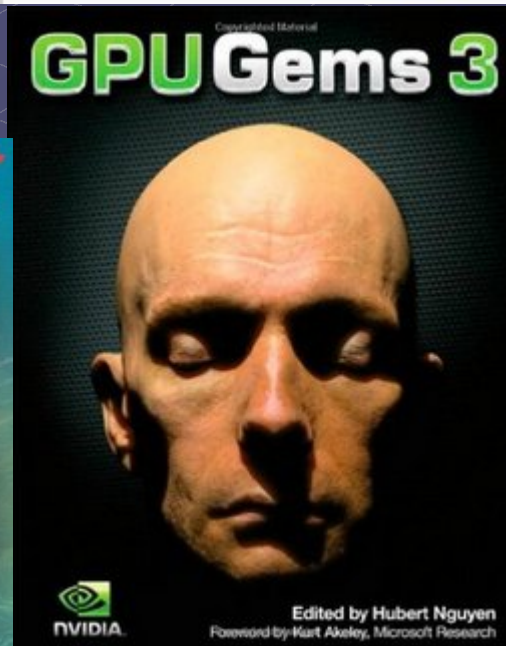
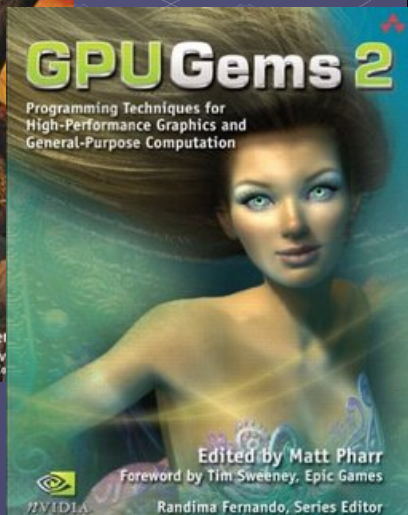
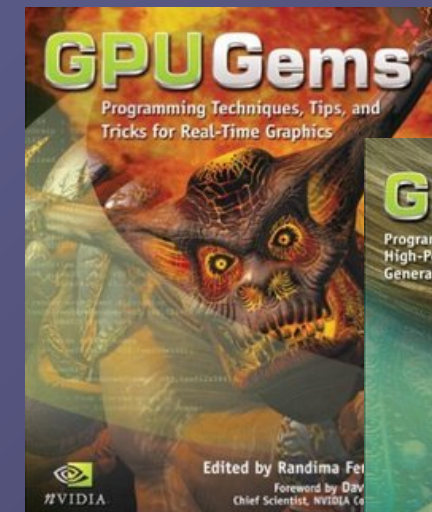
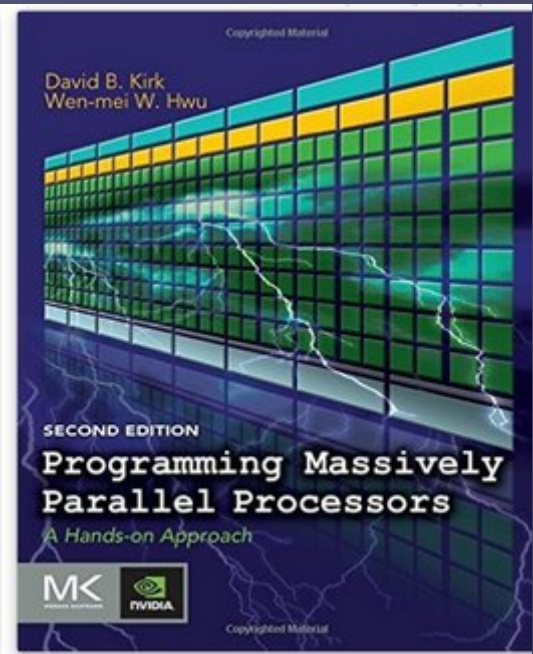
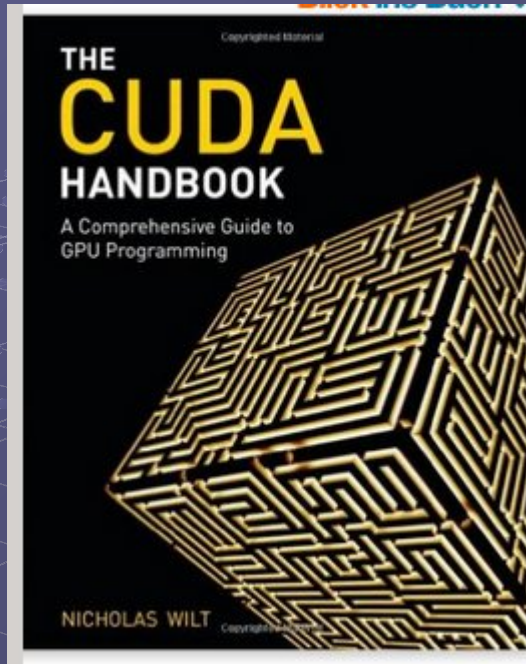
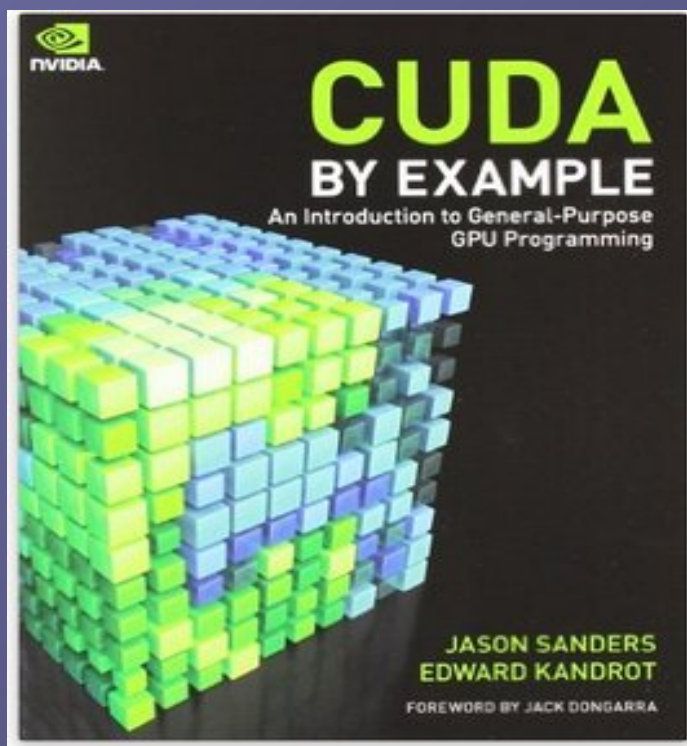
# Introduction to GPU Accelerated Computing

February 11 – 14, 2019

## Table of Contents (subject to adjustment/change):

1. Monday morning: General Introduction Computer Architecture, Many-Core, GPU and others..., Access...
2. Monday afternoon: Access to kepler, CUDA Hello, GPU Properties, Simple Add, Vector Add
3. Tuesday morning: More on GPU Software and Hardware
4. Tuesday afternoon: CUDA More Vector Add, Scalar Products, Using Blocks and Threads
5. Wednesday morning: Parallelization and Amdahl's Law, GPU Acceleration, Future Architecture
6. Wednesday Afternoon: Events, Histograms, Matrix Multiplication
7. Thursday Morning: Astrophysical N-Body Code
8. Thursday Afternoon: Astrophysical Parallel N-Body Code Using MPI and GPU
9. Access: Use **ssh-keygen -t rsa** (give passphrase)  
Send **id\_rsa.pub** to **spurzem@ari.uni-heidelberg.de**

# Literature







Observations (Experiment)



Theory



Computational Physics





GPU Computing

History



# History



**Erik Holmberg (1908-2000)**

Dissertation Univ. Lund (Schweden) (1937):

“A study of double and multiple galaxies”

Galaxies often in Groups and Pairs

Irregular Distribution of Satellite Galaxies

(Holmberg-Effect)

**Father of numerical astrophysics?**

» **...with 200 light bulbs**



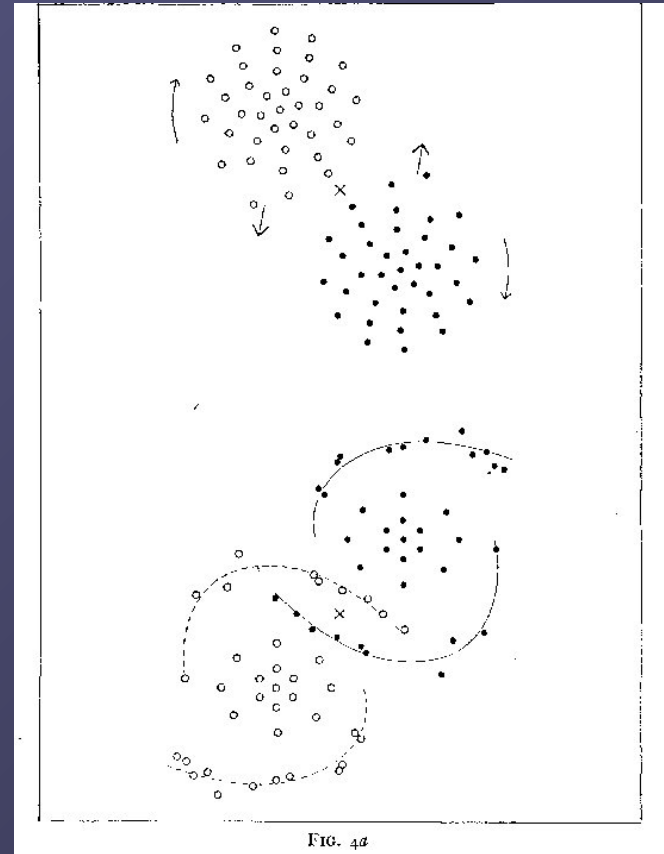
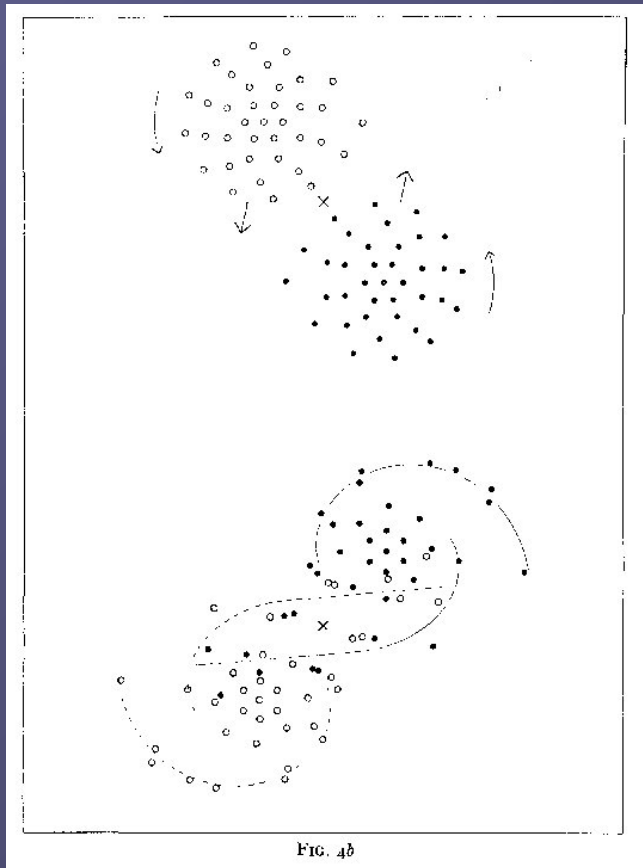
# History

<http://cdsads.u-strasbg.fr/abs/1941ApJ...94..385H>

## The Astrophysical Journal, Nov. 1941



**LUMA METALL**





# HARDWARE

...before von Neumann...

● Konrad Zuse (1910-1995) Berlin



Invented freely programmable Computer



Z1 in parental flat 1936

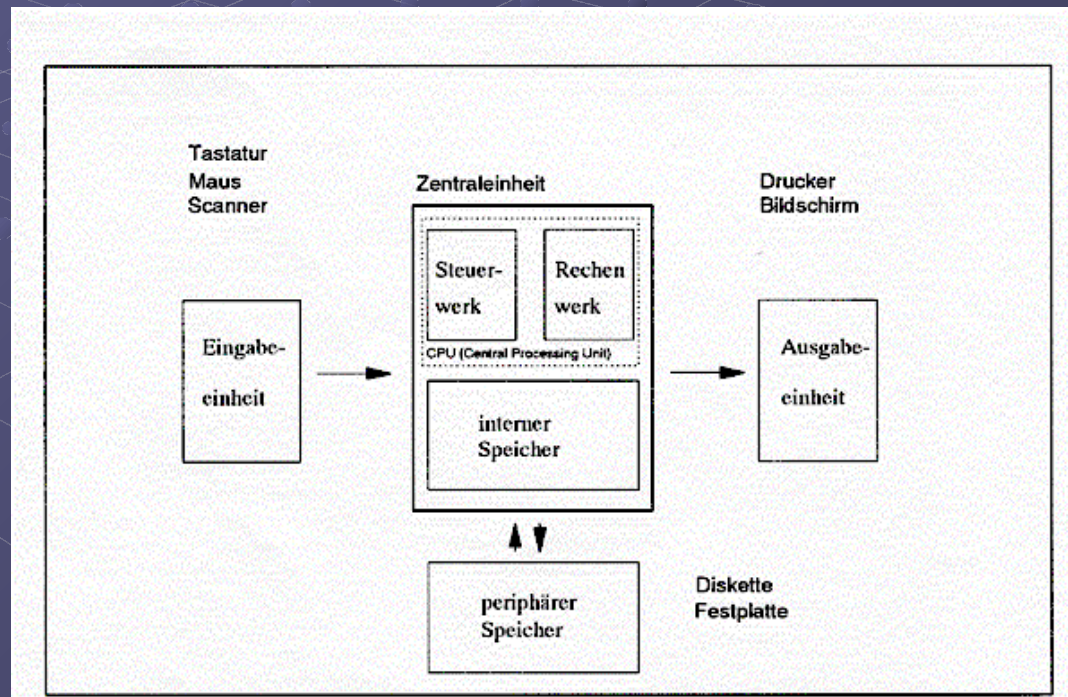


# HARDWARE

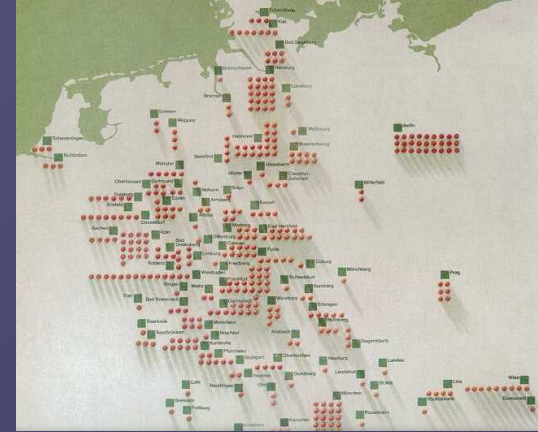
- John von Neumann (1903-1957)

Born Budapest, Lecturer Berlin, since 1930 Princeton Univ.

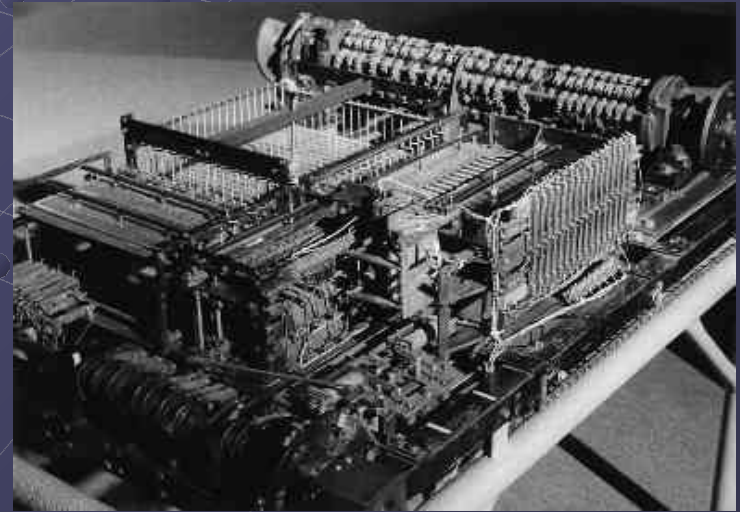
Requirements for the Construction of an electronic computing device(1946)



# History



**Zuse Z4: 1944 Berlin, 1950 Zürich, 1954 Frankreich  
1959 Deutsches Museum München**



**Computing Speed 0.03 MHz**

**Memory 256 byte**





Astronomisches  
Rechen-Institut (ARI)  
at Univ. of  
Heidelberg, Germany



**Siemens 2002  
Computer in 1964  
At ARI**

# History

<http://cdsads.u-strasbg.fr/abs/1960ZA.....50..184V>

Astronomisches Rechen-Institut in Heidelberg  
Mitteilungen Serie A Nr. 14

## Die numerische Integration des $n$ -Körper-Problemes für Sternhaufen I

Von

SEBASTIAN VON HOERNER

Mit 3 Textabbildungen

*(Eingegangen am 10. Mai 1960)*

Astronomisches Rechen-Institut in Heidelberg  
Mitteilungen Serie A Nr. 19

## Die numerische Integration des $n$ -Körper-Problems für Sternhaufen, II.

Von

SEBASTIAN VON HOERNER

Mit 10 Textabbildungen

*(Eingegangen am 19. November 1962)*

<http://cdsads.u-strasbg.fr/abs/1963ZA.....57...47V>

Tabelle 5. *Zahl der gegenseitigen Umläufe, Häufigkeit des Auftretens und kleinster gegenseitiger Abstand  $D_m$  der engsten Paare.* (Alle engsten Paare mit mehr als zwei vollen Umläufen wurden notiert)

Umläufe	Häufigkeit	$D_m$
2—3	11	0.0102
3—5	9	0.0177
5—10	5	0.0070
10—20	2	0,0141
20—50	1	0.0007
50—100	1	0.0035
100—200	1	0.0039

S.v. Hoerner,  
Z.f.Astroph. 1960, 63

Siemens 2002  
N=4,8,12,16 (4 Trx)

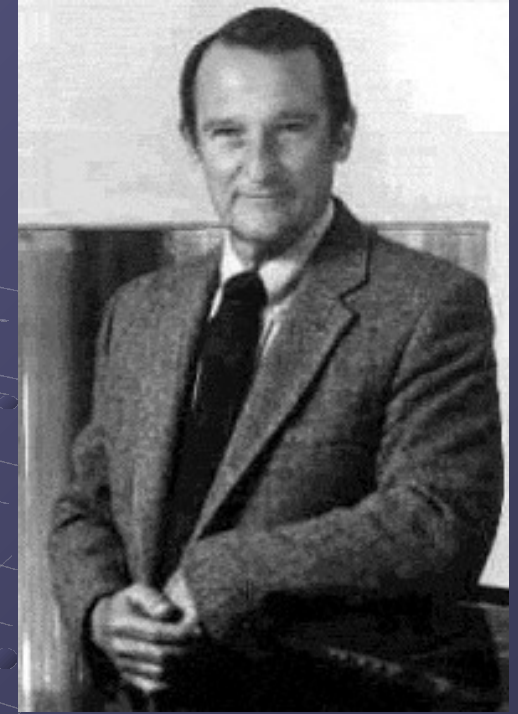
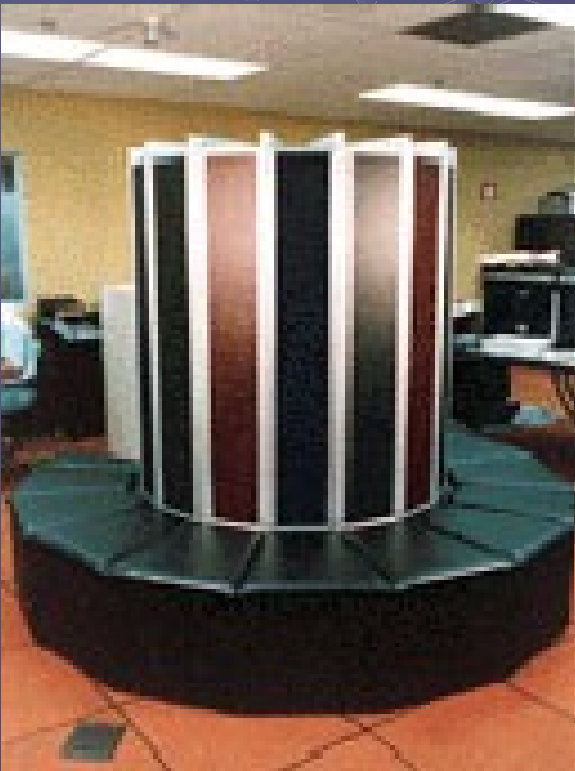
N=16,25 (40 Trx)



# History

## ● Seymour Cray (1925-1996)

“father of supercomputing”



**CRAY1: Vectorregisters (1976)**

**160 Mflop, 80 MHz, 8 MByte RAM**

**CRAY2: (1984)**

**1Gflop, 120MHz, 2GByte RAM**

# History

*Supercomputer  
JUGENE  
IBM Blue Gene  
At FZ Jülich,  
Germany*



*Opening Ceremony June 2008*





# Computational Science...

...after von Neumann...

Exaflop/s?

Petaflop/s

Teraflop/s

Gigaflop/s

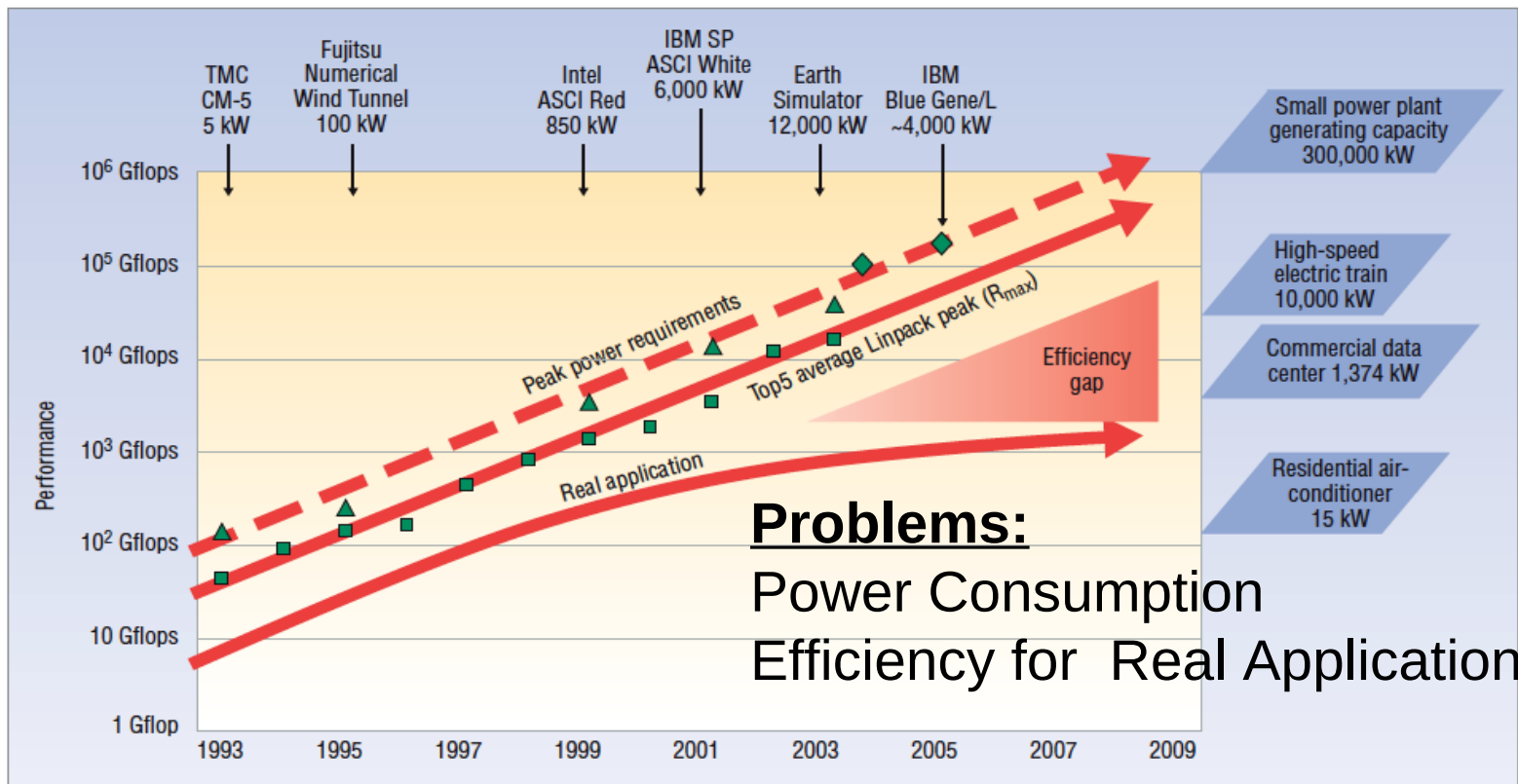


Figure 1. Rising power requirements. Peak power consumption of the top supercomputers has steadily increased over the past 15 years. Thanks to Horst Simon, LBNL/NERSC for this diagram.

GPU Computing

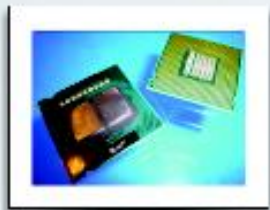
# Special Hardware Accelerators



# SPECIAL HARDWARE

## CPUs

Central Processing Units



General Purpose oriented

1-12 Cores

Up to 4 pipes per core using Vector Units

Fully Programmable, many languages available

Very well studied

Max. 125W per processor

## GPUs

Graphic Processing Units



Graphics oriented

16-512 Cores

Massively Parallel Architecture, specialized instructions for parallel processing

Fully programmable, but limited languages

Algorithms not fully explored

Max. 400W per card

## FPGAs

Field Programmable Gate Arrays



Custom designs, best for processing streaming data

Programmable Logic, Architecture is custom-built for the required application

Requires extensive knowledge to program, development time is longer than CPUs and GPUs

Application interface is custom built on each case

Max. 60W per FPGA

## ASICs

Application Specific Integrated Circuits



Fully custom designs, built for a specific application

Not flexible, cannot be changed once it is built

Development is even more specialized than FPGAs

Power consumption varies with the application, usually best performance per Watt

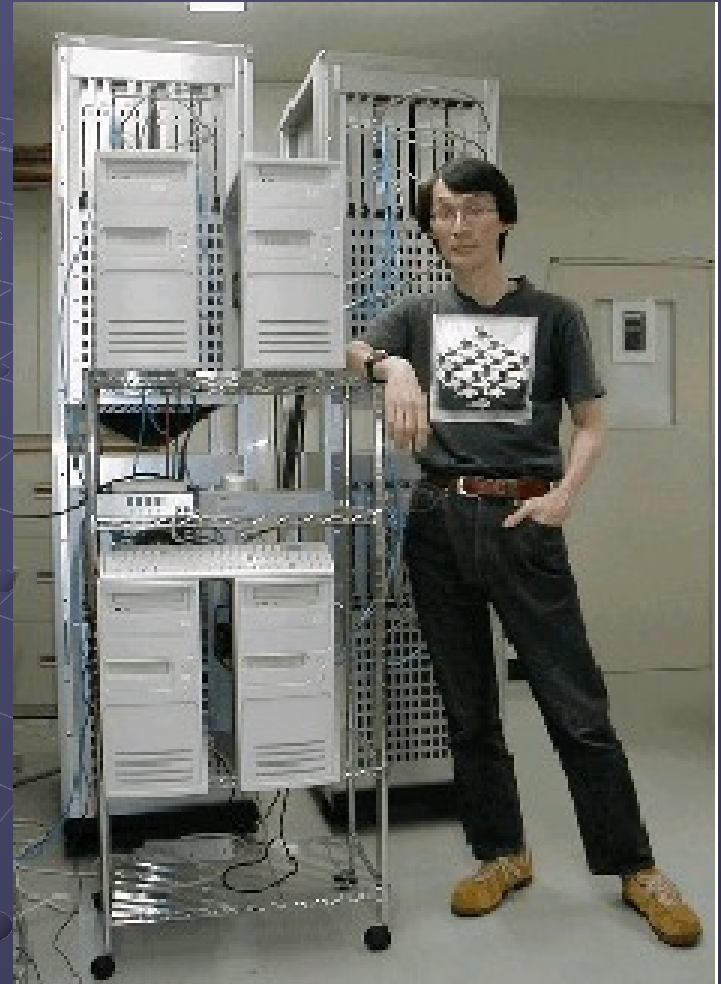
Slide: Guillermo Marcus



# HARDWARE

## GRAPE-6 Gravity/Coulomb Part

- G6 Chip:  $0.25\mu$  2MGate ASIC, 6 Pipelines
- at 90MHz, 31Gflops/chip
- 48Tflops full system (March 2002)
- Plan up to 72Tflops full system (in 2002)
- Installed in Cambridge, Marseille, Drexel, Amsterdam, New York (AMNH), Mitaka (NAO), Tokyo, etc..  
New Jersey, Indiana, Heidelberg





## GRAPE-6



1998, 120  
Gflops

Developers: Junichiro Makino, Toshiyuki Fukushige, Hiroshi Daisaka, Eiichiro Kokubo, Masaki Koga, Makoto Taiji, Ken Namura

[GRAPE-6: Massively-Parallel Special-Purpose Computer for Astrophysical Particle Simulations](#)

[Sales information](#)

## The Green500 List - November 2010

Listed below are the November 2010 The Green500's energy-efficient supercomputers ranked from 1 to 100.

<http://www.green500.org>

Green500 Rank	MFLOPS/W	Site*	Computer*	Total Power (kW)
<u>1</u>	1664.20	IBM Thomas J. Watson Research Center	NNSA/SC Blue Gene/Q Prototype	38.80
<u>2+</u>	1448.03	National Astronomical Observatory of Japan	GRAPE-DR accelerator Cluster, Infiniband	24.59
<u>2</u>	958.35	GSIC Center, Tokyo Institute of Technology	HP ProLiant SL390s G7 Xeon 6C X5670, Nvidia GPU, Linux/Windows	1243.80
<u>3</u>	933.06	NCSA	Hybrid Cluster Core i3 2.93Ghz Dual Core, NVIDIA C2050, Infiniband	36.00

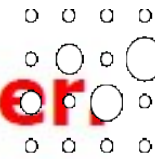
# GPU: NAOC laohu cluster Beijing, China





Heidelberg

# Kepler GPU cluster



VolkswagenStiftung

## Kepler GPU cluster

12 nodes = 12 x 16 = 192 CPU cores (@ 2 GHz)

12 x 64 GB = 768 GB RAM CPU memory

12 GPUs K20m = 12 x 2496 ~ 30k GPU threads

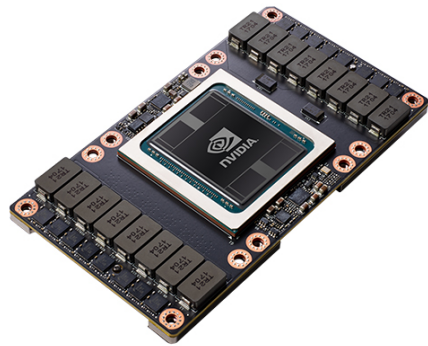
12 x 4.8 GB ~ 57 GB GPU device memory

4 x Xilinx Virtex-6 FPGA (ML 605)

since beg. 2013 operated.



# NVIDIA Volta V100 GPU, 21 billion transistors, 5120 cores



With NVLINK

Without NVLINK



## PERFORMANCE

with NVIDIA GPU Boost\*

DOUBLE-PRECISION

7.8<sub>teraFLOPS</sub>

DOUBLE-PRECISION

7<sub>teraFLOPS</sub>

SINGLE-PRECISION

15.7<sub>teraFLOPS</sub>

SINGLE-PRECISION

14<sub>teraFLOPS</sub>

DEEP LEARNING

125<sub>teraFLOPS</sub>

DEEP LEARNING

112<sub>teraFLOPS</sub>

## INTERCONNECT BANDWIDTH

Bi-Directional

NVLINK

300<sub>GB/s</sub>

PCIe

32<sub>GB/s</sub>

## MEMORY

CoWoS Stacked HBM2

CAPACITY

32/16<sub>GB HBM2</sub>

BANDWIDTH

900<sub>GB/s</sub>

## POWER

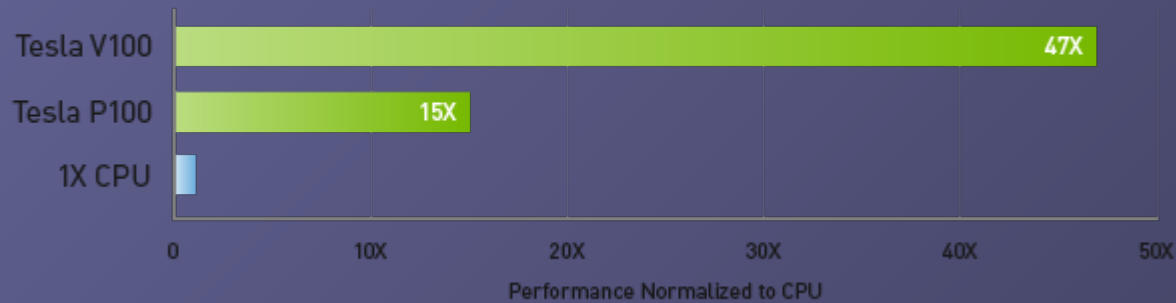
Max Consumption

300<sub>WATTS</sub>

250<sub>WATTS</sub>

# NVIDIA Volta V100 GPU, 21 billion transistors, 5120 cores

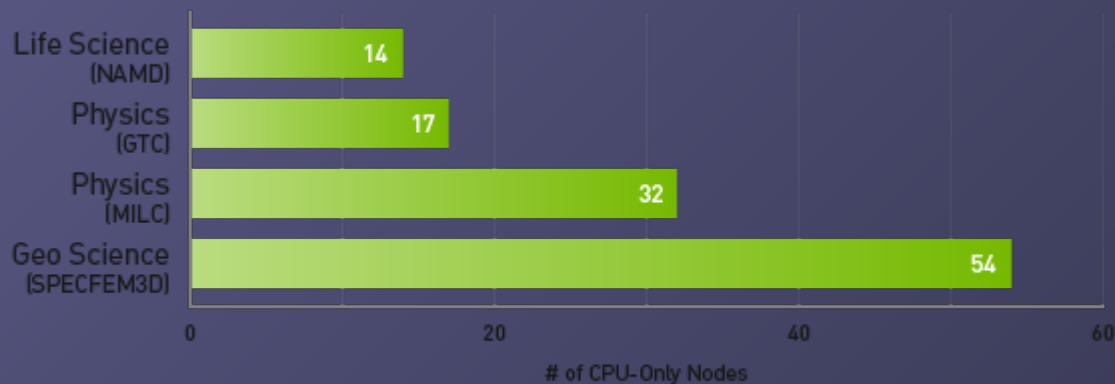
## 47X Higher Throughput Than CPU Server on Deep Learning Inference



Workload: ResNet-50 | CPU: 1X Xeon E5-2690v4 @ 2.6 GHz | GPU: Add 1X Tesla P100 or V100

## 1 GPU Node Replaces Up To 54 CPU Nodes

Node Replacement: HPC Mixed Workload



CPU Server: Dual Xeon Gold 6140@2.30GHz, GPU Servers: same CPU server w/ 4x V100 PCIe | CUDA Version: CUDA 9.x | Dataset: NAMD (STMV), GTC (mpi#proc.in), MILC (APEX Medium), SPECFEM3D (four\_material\_simple\_model) | To arrive at CPU node equivalence, we use measured benchmark with up to 8 CPU nodes. Then we use linear scaling to scale beyond 8 nodes.



# Top 10 List November 2010

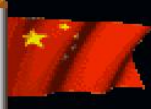
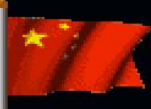
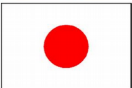

From [www.top500.org](http://www.top500.org) - list of fastest

supercomputers in the world...  
... last year Nov. 2010:

## ► China Grabs Supercomputing Leadership Spot in Latest Ranking of World's Top 500 Supercomputers

Thu, 2010-11-11 22:42

MANNHEIM, Germany; BERKELEY, Calif.; and KNOXVILLE, Tenn.—The 36<sup>th</sup> edition of the closely watched TOP500 list of the world's most powerful supercomputers confirms the rumored takeover of the top spot by the Chinese Tianhe-1A system at the National Supercomputer Center in Tianjin, achieving a performance level of 2.57 petaflop/s (quadrillions of calculations per second).

1	National Supercomputing Center in Tianjin China		<b>Tianhe-1A</b> - NUDT TH MPP, X5670 2.93Ghz 6C, NVIDIA GPU, FT-1000 8C NUDT	<b><u>GPU</u></b>
2	DOE/SC/Oak Ridge National Laboratory United States		<b>Jaguar</b> - Cray XT5-HE Opteron 6-core 2.6 GHz Cray Inc.	
3	National Supercomputing Centre in Shenzhen (NSCS) China		<b>Nebulae</b> - Dawning TC3600 Blade, Intel X5650, NVidia Tesla C2050 GPU Dawning	<b><u>GPU</u></b>
4	GSIC Center, Tokyo Institute of Technology Japan		<b>TSUBAME 2.0</b> - HP ProLiant SL390s G7 Xeon 6C X5670, Nvidia GPU, Linux/Windows NEC/HP	<b><u>GPU</u></b>
5	DOE/SC/LBNL/NERSC United States		<b>Hopper</b> - Cray XE6 12-core 2.1 GHz Cray Inc.	
6	Commissariat a l'Energie Atomique (CEA) France	<b>FR</b>	<b>Tera-100</b> - Bull bullx super-node S6010/S6030 Bull SA	
7	DOE/NNSA/LANL United States		<b>Roadrunner</b> - BladeCenter QS22/LS21 Cluster, PowerXCell 8i 3.2 Ghz / Opteron DC 1.8 GHz, Voltaire Infiniband IBM	
8	National Institute for Computational Sciences/University of Tennessee United States		<b>Kraken XT5</b> - Cray XT5-HE Opteron 6-core 2.6 GHz Cray Inc.	
9	Forschungszentrum Juelich (FZJ) Germany		<b>JUGENE</b> - Blue Gene/P Solution IBM	
10	DOE/NNSA/LANL/SNL United States		<b>Cielo</b> - Cray XE6 8-core 2.4 GHz Cray Inc.	

# NCSA director: GPU is future of supercomputing

by Brooke Crothers



Font size



Print



E-mail



Share



6 comments

Tweet

99



Share

25

2

Digg ↑

The director of the National Center for Supercomputing Applications has seen the future of supercomputing and it can be summed up in three letters: GPU.

Thom Dunning, who directs the NCSA and the Institute for Advanced Computing Applications and Technologies at the famed supercomputing facilities on the campus of University of Illinois at Urbana-Champaign, says high-performance computing will begin to move toward graphics processing units or GPUs. Not coincidentally, **this is exactly what China has done to achieve the world's fastest speeds with its "Tianhe-1A"** supercomputer. That computer combines about 7,000 Nvidia GPUs with 14,000 Intel CPUs: the only hybrid CPU-GPU system in the world of that scale.

"What we're really seeing in the efforts in China as well as the ones we have in the U.S. is that GPUs are what the future will look like," said Dunning in a phone interview Thursday. "What we're seeing is the beginning of something that's going to be happening all over the world."

NCSA already has a small CPU-GPU hybrid system. "It's something we have been working on for a number of years. We have a CPU-GPU cluster for the NCSA academic community. Made up of Intel CPUs and Nvidia GPUs. A 50 teraflop machine," he said. (Note that **Oak Ridge National Laboratories is also installing a hybrid system now.**)



Thom Dunning directs the Institute for Advanced Computing Applications and Technologies and the NCSA.

# Intel MIC Hardware

INSPUR, NAOC - 2013.XI.26



icpc ... "-mmic" ...  $61 \times 4 = 244$  x 1.1 GHz omp cores !!!  
Full fp64 !!!



# Intel MIC Hardware

## Intel® Xeon Phi™ Coprocessor Family Reference Table

SKU #	Form Factor, Thermal	Peak Double Precision	Max # of Cores	Clock Speed (GHz)	GDDR5 Memory Speeds (GT/s)	Peak Memory BW	Memory Capacity (GB)	Total Cache (MB)	Board TDP (Watts)	Process
SE10P <small>(special edition)</small>	PCIe Card, Passively Cooled	1073 GF	61	1.1	5.5	352	8	30.5	300	22nm
SE10X <small>(special edition)</small>	PCIe Card, No Thermal Solution	1073 GF	61	1.1	5.5	352	8	30.5	300	
5110P	PCIe Card, Passively Cooled	1011 GF	60	1.053	5.0	320	8	30	225	
3100 Series	PCIe Card, Actively Cooled	>1 TF	Disclosed at 3100 series launch (H1'13)		5.0	240	6	28.5	300	
	PCIe Card, Passively Cooled	> 1 TF			5.0	240	6	28.5	300	



PCIe Card, Actively Cooled



PCIe Card, Passively Cooled

Current Generation:  
Knights Landing  
14nm

## Intel MIC hardware / Recent Processors



### Intel® Xeon Phi™ Processor 7290

- 36 MB L2 Cache
- 72 Cores
- 72 Threads
- 1.70 GHz Max Turbo Frequency



### Intel® Xeon Phi™ Processor 7290F

- 36 MB L2 Cache
- 72 Cores
- 72 Threads
- 1.70 GHz Max Turbo Frequency



# Supercomputer from China: 96/33 Pflop/s Linpack Wuxi/Guangzhou/Tianjin National Supercomputing Center Taihu 10 mill. cores

**Tianhe-2 (MilkyWay-2) - TH-IVB-FEP Cluster, Intel Xeon E5-2692 12C 2.200GHz, TH Express-2, Intel Xeon Phi 31S1P**



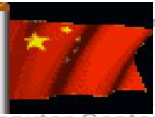



32000 Intel Xeon 12 core,  
48000 Intel Phi Accelerators 57 Core,  
now Chinese processor



Test of Taihu planned;  
But:  
Local cluster with new  
GPUs at NAOC gives  
much more resources.



# Top 10 List November 2018 (from [www.top500.org](http://www.top500.org) )

Rank	Site	System	Cores	Rmax (TFlop/s)	Rpeak (TFlop/s)	Power (kW)
1	DOE/SC/Oak Ridge National Laboratory United States	<b>Summit</b> - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband IBM	2,397,824	143,500.0	200,794.9	9,783
	<b>USA</b>			<b><u>GPU Volta</u></b>		
2	DOE/NNSA/LLNL United States	<b>Sierra</b> - IBM Power System S922LC, IBM POWER9 22C 3.1GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband IBM / NVIDIA / Mellanox	1,572,480	94,640.0	125,712.0	7,438
	<b>USA</b>			<b><u>GPU Volta</u></b>		
3	National Supercomputer Center Wuxi China	<b>Sunway TaihuLight</b> - Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway NRPC	10,649,600	93,014.6	125,435.9	15,371
				<b><u>Chinese Processor</u></b>		
4	National Super Computer Center Guangzhou China	<b>Tianhe-2A</b> - TH-IVB-FEP Cluster, Intel Xeon E5-2692v2 12C 2.2GHz, TH Express-2, Matrix-2000 NUDT	4,981,760	61,444.5	100,678.7	18,482
				<b><u>Chinese Processor</u></b>		
5	Swiss National Supercomputing Centre (CSCS) Switzerland	<b>Piz Daint</b> - Cray XC50, Xeon E5-2690v3 12C 2.6GHz, Aries interconnect , NVIDIA Tesla P100 Cray Inc.	387,872	21,230.0	27,154.3	2,384
	<b>Swiss</b>			<b><u>GPU Pascal</u></b>		
6	DOE/NNSA/LANL/SNL United States	<b>Trinity</b> - Cray XC40, Xeon E5-2698v3 16C 2.3GHz, Intel Xeon Phi 7250 68C 1.4GHz, Aries interconnect Cray Inc.	979,072	20,158.7	41,461.2	7,578
	<b>USA</b>			<b><u>Xeonφ</u></b>		
7	National Institute of Advanced Industrial Science and Technology (AIST) Japan	<b>AI Bridging Cloud Infrastructure (ABCI)</b> - PRIMERGY CX2570 M4, Xeon Gold 6148 20C 2.4GHz, NVIDIA Tesla V100 SXM2, Infiniband EDR Fujitsu	391,680	19,880.0	32,576.6	1,649
				<b><u>GPU Volta</u></b>		
8	Leibniz Rechenzentrum Germany	<b>SuperMUC-NG</b> - ThinkSystem SD530, Xeon Platinum 8174 24C 3.1GHz, Intel Omni-Path Lenovo	305,856	19,476.6	26,873.9	
				<b><u>GPU Volta (part)</u></b>		
9	DOE/SC/Oak Ridge National Laboratory United States	<b>Titan</b> - Cray XK7, Opteron 6274 16C 2.200GHz, Cray Gemini interconnect, NVIDIA K20x Cray Inc.	560,640	17,590.0	27,112.5	8,209
	<b>USA</b>			<b><u>GPU Kepler</u></b>		
10	DOE/NNSA/LLNL United States	<b>Sequoia</b> - BlueGene/Q, Power BQC 16C 1.60 GHz, Custom IBM	1,572,864	17,173.2	20,132.7	7,890
	<b>USA</b>			<b><u>Blue Gene</u></b>		

Saudi-A.

# TOP500 List Refreshed, US Edged Out of Third Place ■ ■ ■ ■

TOP500 Team | June 19, 2017 00:22 CEST

FRANKFURT, Germany; BERKELEY, Calif.; and KNOXVILLE, Tenn.— The 49th edition of the TOP500 list was released today in conjunction with the opening session of the ISC High Performance conference, which is taking place this week in Frankfurt, Germany. The list ranks the world's most powerful supercomputers based on the Linpack benchmark and is released twice per year.

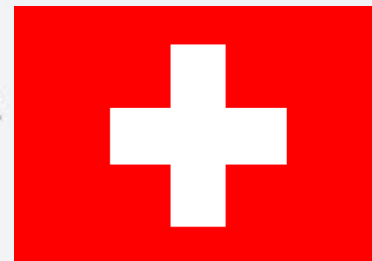
[Read more](#)

## System

Sunway TaihuLight - Sunway MPP, Sunway SW26010 260C 1.450  
Sunway , NRCPC  
National Supercomputing Center in Wuxi  
China

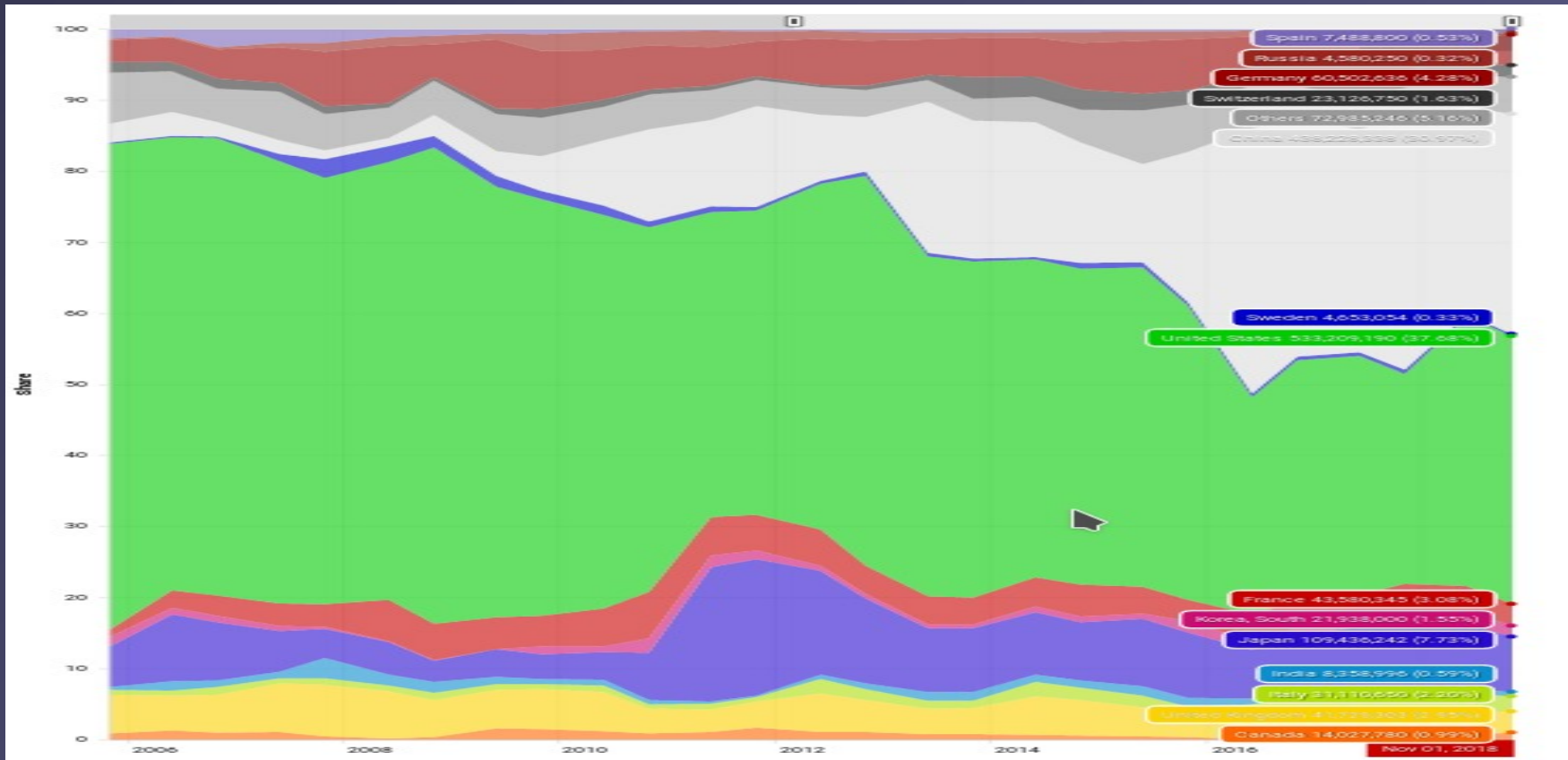
Tianhe-2 (MilkyWay-2) - TH-IVB-FEP Cluster, Intel Xeon E5-265  
2.200GHz, TH Express-2, Intel Xeon Phi 31S1P , NUDT  
National Super Computer Center in Guangzhou  
China

Piz Daint - Cray XC50, Xeon E5-2690v3 12C 2.6GHz, Aries interco  
NVIDIA Tesla P100 , Cray Inc.  
Swiss National Supercomputing Centre (CSCS)  
Switzerland



■ ■ ■ ■  
**By  
Switzerland**

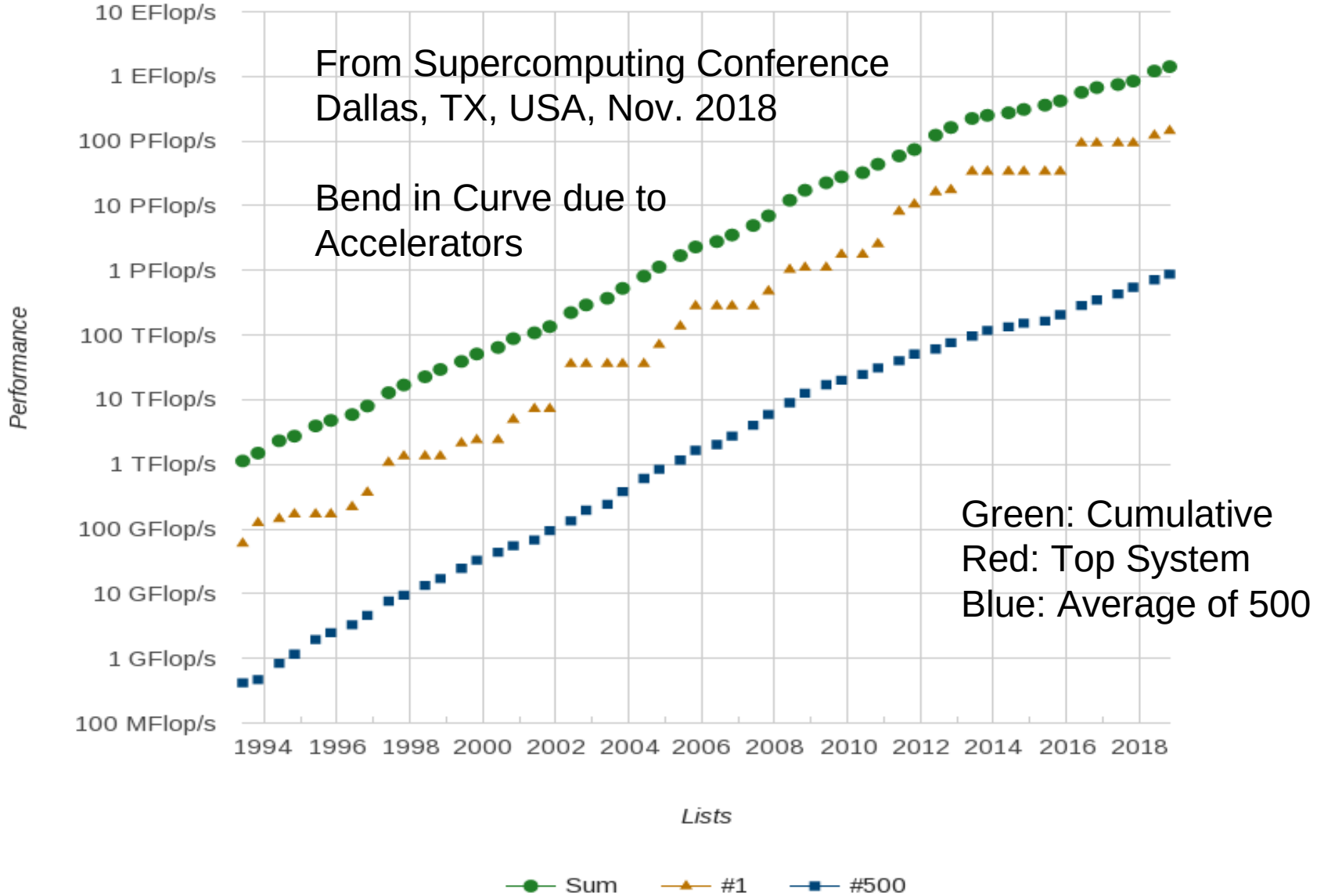
# Top 500 List November 2018 – Performance Share of Countries



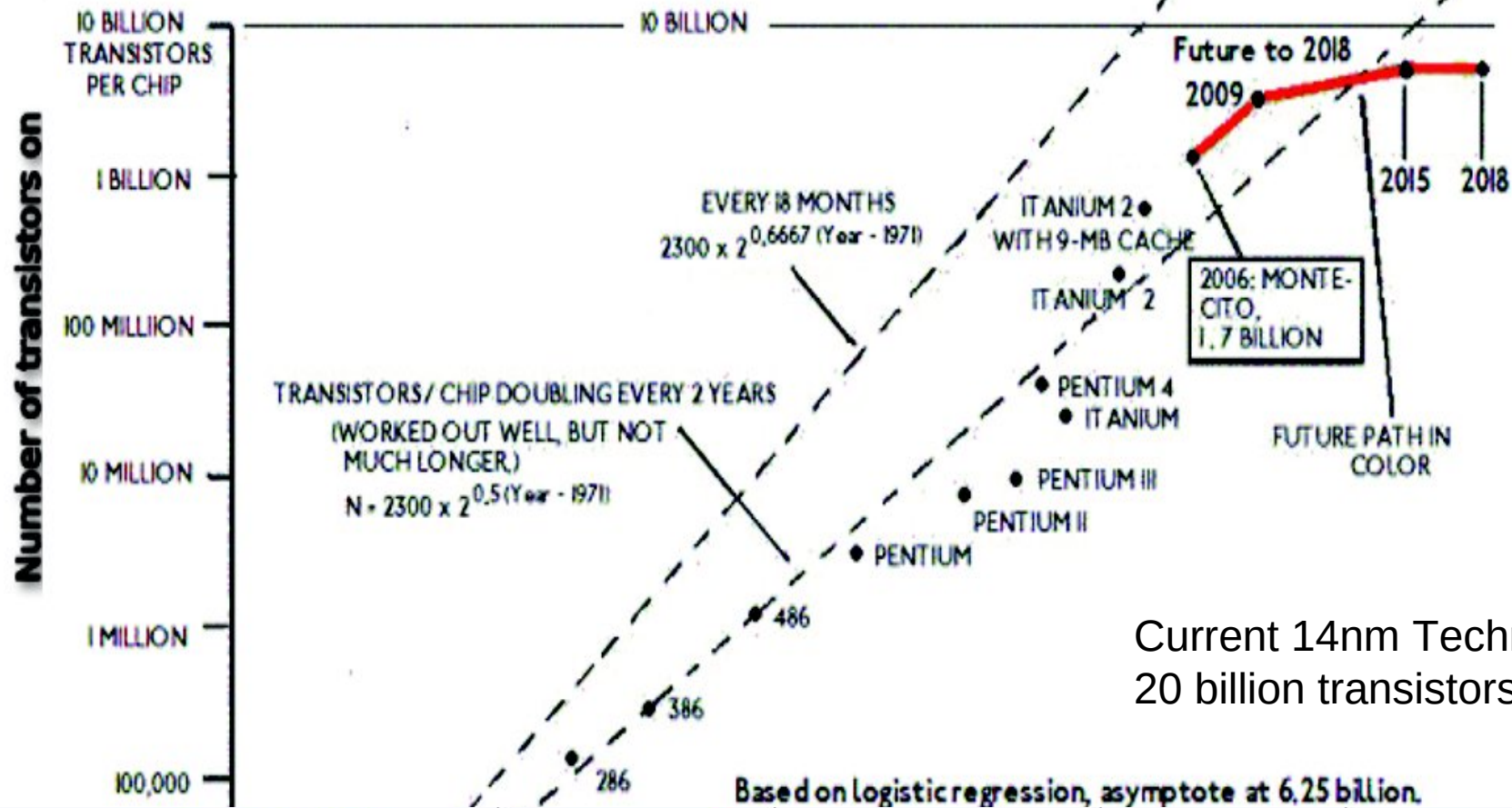


# Performance Development

# Moore's Law?



# Moore's Law Ending (Red Line): Delayed products, Delayed 45nm / 32 nm, Reduced Capex



Current 14nm Technology  
20 billion transistors

Based on logistic regression, asymptote at 6.25 billion.

22-core Xeon Broadwell-E5	7,200,000,000 <sup>[36]</sup>	2016	Intel	14 nm	456 mm <sup>2</sup>
SPARC M7	10,000,000,000 <sup>[37]</sup>	2015	Oracle	20 nm	
24-core AMD EPYC 7401P	19,200,000,000	2017	AMD	14 nm	195 mm <sup>2</sup>



by Clayton Kallmark  
Dedicated to  
Professor Frederick E. Terman



# GREEN 500 list – Power Efficiency (Gflops/Watts), see also <http://www.top500.org>

Rank	TOP500 Rank	System	Cores	Rmax (TFlop/s)	Power (kW)	Power Efficiency (GFlops/watts)
1	375	<b>Shoubu system B</b> - ZettaScaler-2.2, Xeon D-1571 16C 1.3GHz, Infiniband EDR, PEZY-SC2, PEZY Computing / Exascaler Inc. Advanced Center for Computing and Communication, RIKEN Japan	953,280	1,063.3	60	17.604
<b><u>Japan</u></b>						
2	374	<b>DGX SaturnV Volta</b> - NVIDIA DGX-1 Volta36, Xeon E5-2698v4 20C 2.2GHz, Infiniband EDR, NVIDIA Tesla V100, Nvidia NVIDIA Corporation United States	22,440	1,070.0	97	15.113
<b><u>GPU Volta</u></b>						
3	1	<b>Summit</b> - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband, IBM DOE/SC/Oak Ridge National Laboratory United States	2,397,824	143,500.0	9,783	14.668
<b><u>GPU Volta</u></b>						
4	7	<b>AI Bridging Cloud Infrastructure (ABCI)</b> - PRIMERGY CX2570 M4, Xeon Gold 6148 20C 2.4GHz, NVIDIA Tesla V100 SXM2, Infiniband EDR, Fujitsu National Institute of Advanced Industrial Science and Technology (AIST) Japan	391,680	19,880.0	1,649	14.423
<b><u>GPU Volta</u></b>						
5	22	<b>TSUBAME3.0</b> - SGI ICE XA, IP139-SXM2, Xeon E5-2680v4 14C 2.4GHz, Intel Omni-Path, NVIDIA Tesla P100 SXM2, HPE GSIC Center, Tokyo Institute of Technology Japan	135,828	8,125.0	792	13.704
<b><u>GPU Pascal</u></b>						
6	2	<b>Sierra</b> - IBM Power System S922LC, IBM POWER9 22C 3.1GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband, IBM / NVIDIA / Mellanox DOE/NNSA/LLNL United States	1,572,480	94,640.0	7,438	12.723
<b><u>GPU Volta</u></b>						
7	446	<b>AIST AI Cloud</b> - NEC 4U-8GPU Server, Xeon E5-2630Lv4 10C 1.8GHz, Infiniband EDR, NVIDIA Tesla P100 SXM2, NEC National Institute of Advanced Industrial Science and Technology Japan	23,400	961.0	76	12.681
<b><u>GPU Pascal</u></b>						
8	411	<b>MareNostrum P9 CTE</b> - IBM Power System AC922, IBM POWER9 22C 3.1GHz, Dual-rail Mellanox EDR Infiniband, NVIDIA Tesla V100, IBM Barcelona Supercomputing Center Spain	19,440	1,018.0	86	11.865
<b><u>GPU Volta</u></b>						
9	38	<b>Advanced Computing System(PreE)</b> - Sugon TC8600, Hygon Dhyana 32C 2GHz, Deep Computing Processor, 200Gb 6D-Torus, Sugon Sugon China	163,840	4,325.0	380	11.382
<b><u>China</u></b>						
10	20	<b>Taiwania 2</b> - QCT QuantaGrid D52G-4U/LC, Xeon Gold 6154 18C 3GHz, Mellanox InfiniBand EDR, NVIDIA Tesla V100 SXM2, Quanta Computer / Taiwan Fixed Network / ASUS Cloud National Center for High Performance Computing Taiwan	170,352	9,000.0	798	11.285
<b><u>GPU Volta</u></b>						



GPU Computing

More on GPU

# Graphics Processors (GPU) as General Purpose Supercomputers (GPGPU)



2008...

GeForce 9800 GTX, 128 Stream Proc., 512 MB

GeForce 9800 GX2, 256 Stream Proc., 1 GB

GeForce 9800 GT, 64 Stream Proc., 512 MB

[...]

2009: Tesla ~200 Proc., 4GB

2010: Fermi ~400 Proc., 4GB

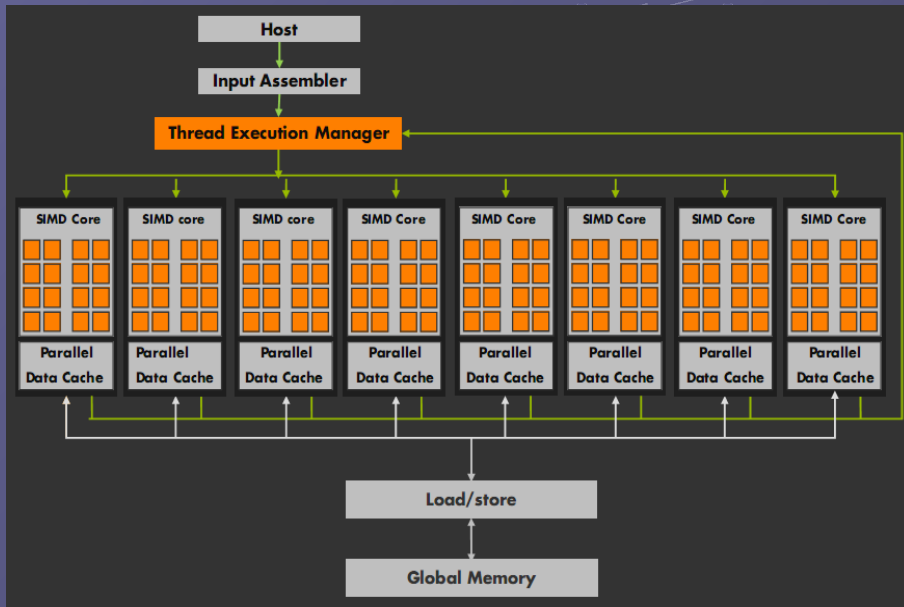
2013: Kepler K20, ~2500 Procs., 6GB

2016: Kepler K80, ~5000 Procs.

2017/18: Pascal, Volta > 5000 Procs.

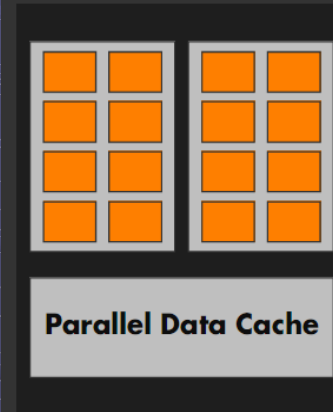


# Hardware around 2006



## Each core

- 8 functional units
- SIMD 16/32 "warp"
- 8-10 stage pipeline
- Thread scheduler
- 128-512 threads/core
- 16 KB shared memory



## Total #threads/chip

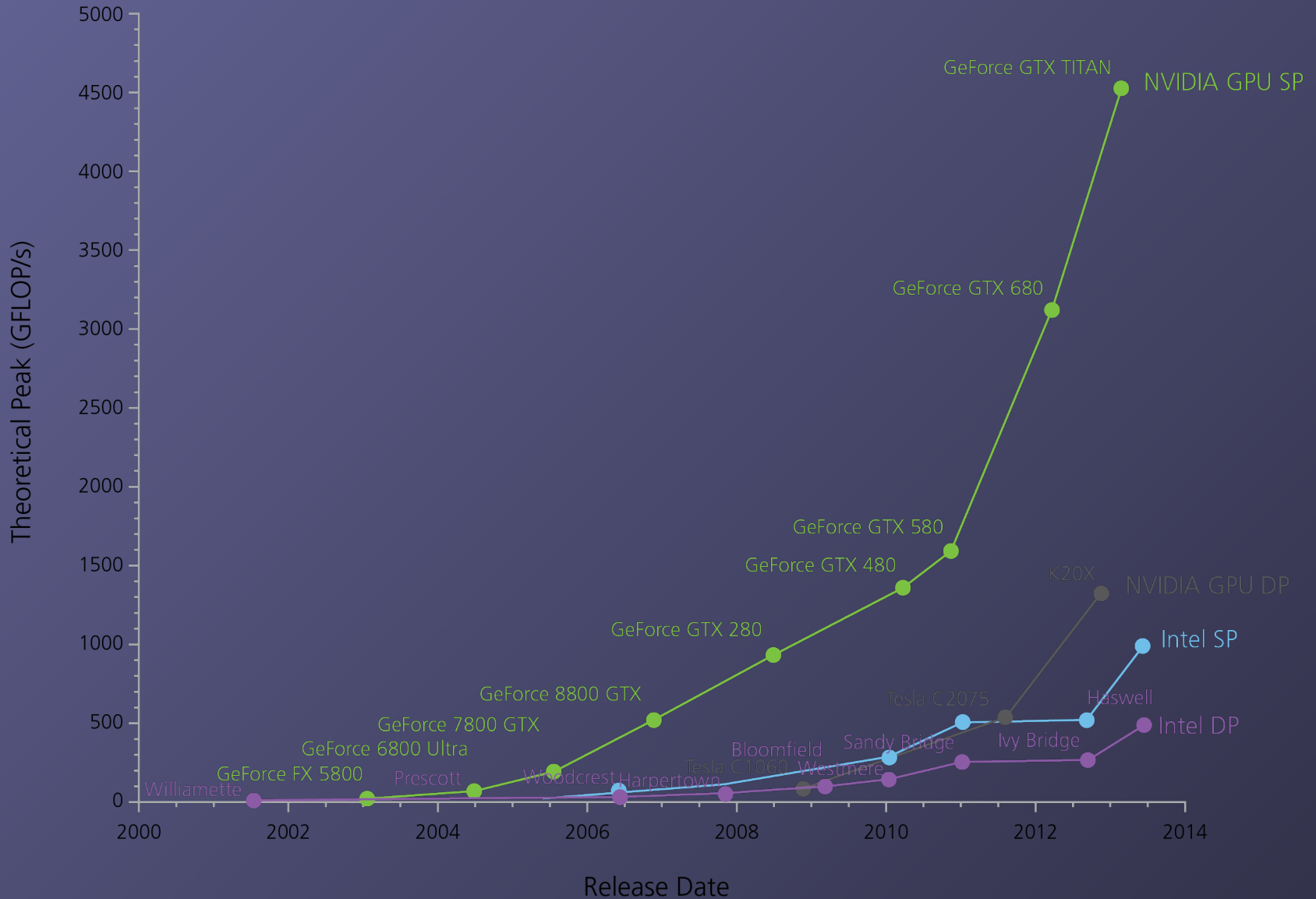
$$16 * 512 = 8K$$

## GeForce 8800 GTX:

$$575 \text{ MHz} * 128 \text{ processors} * 2 \text{ flop/inst} * 2 \text{ inst/clock} = 333 \text{ Gflops}$$



# CPU vs. GPU speedup timeline

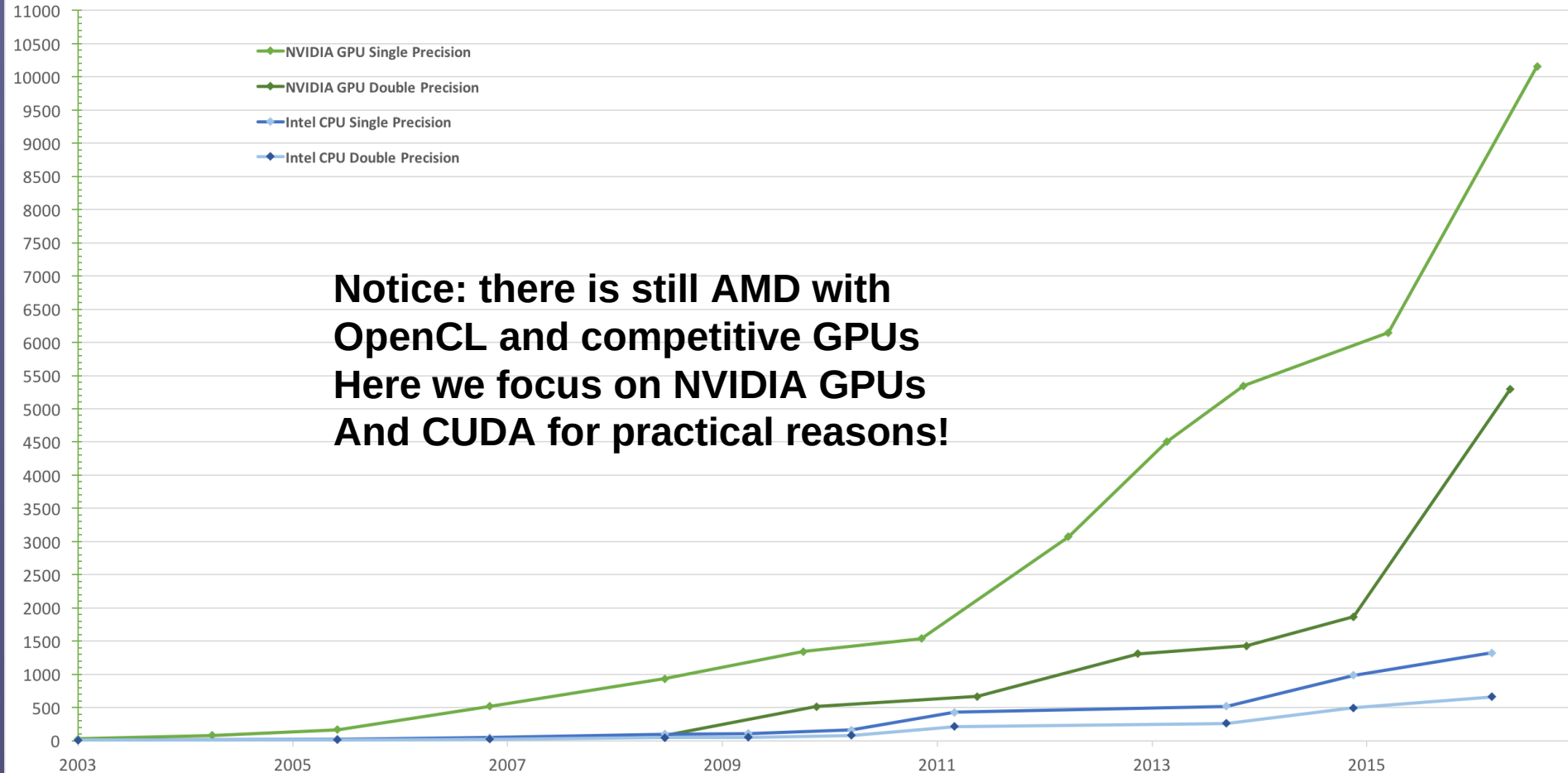


# Floating Point Operations per Second for CPU and GPU:

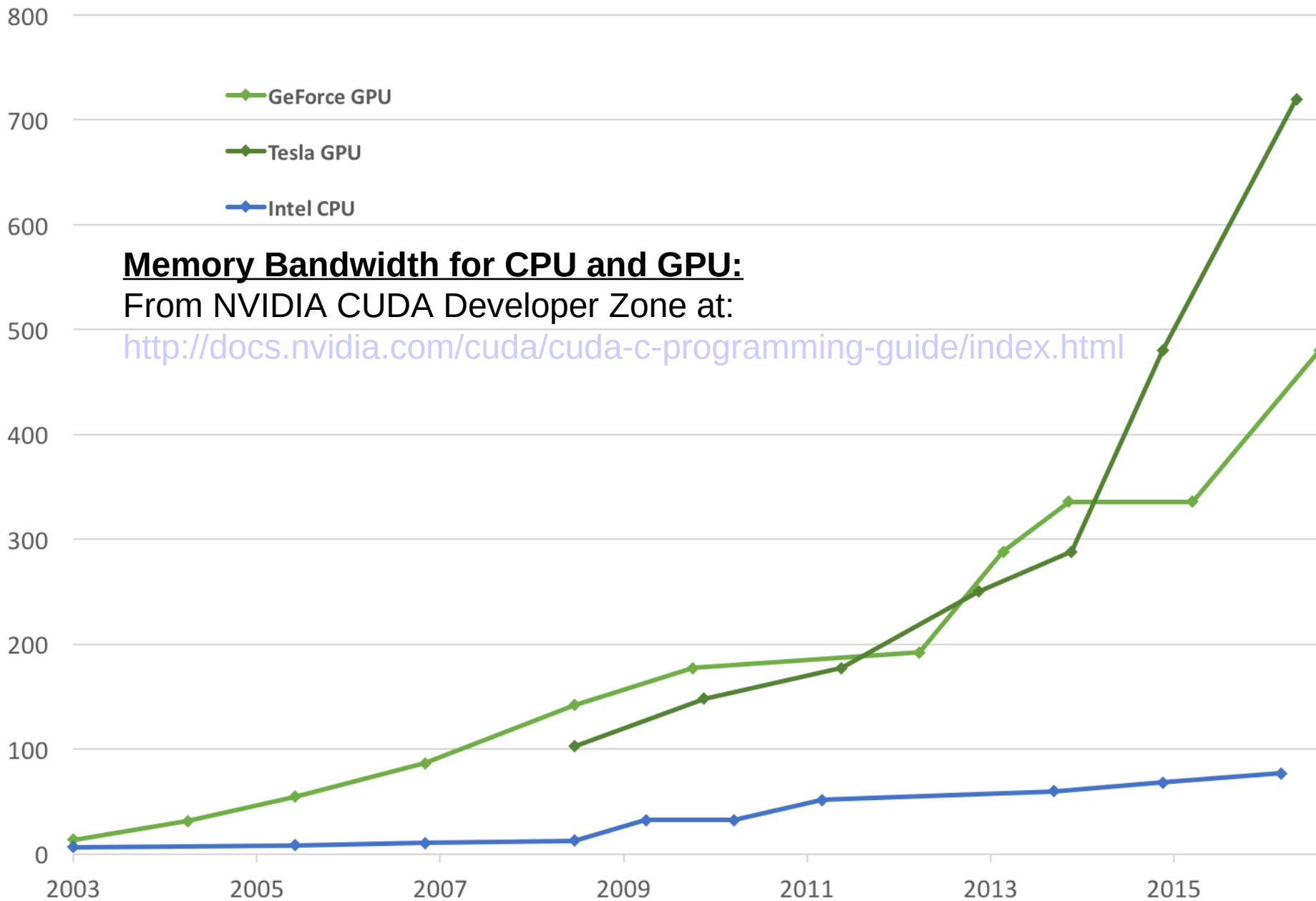
From NVIDIA CUDA Developer Zone at:

<http://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html>

Theoretical GFLOP/s at base clock

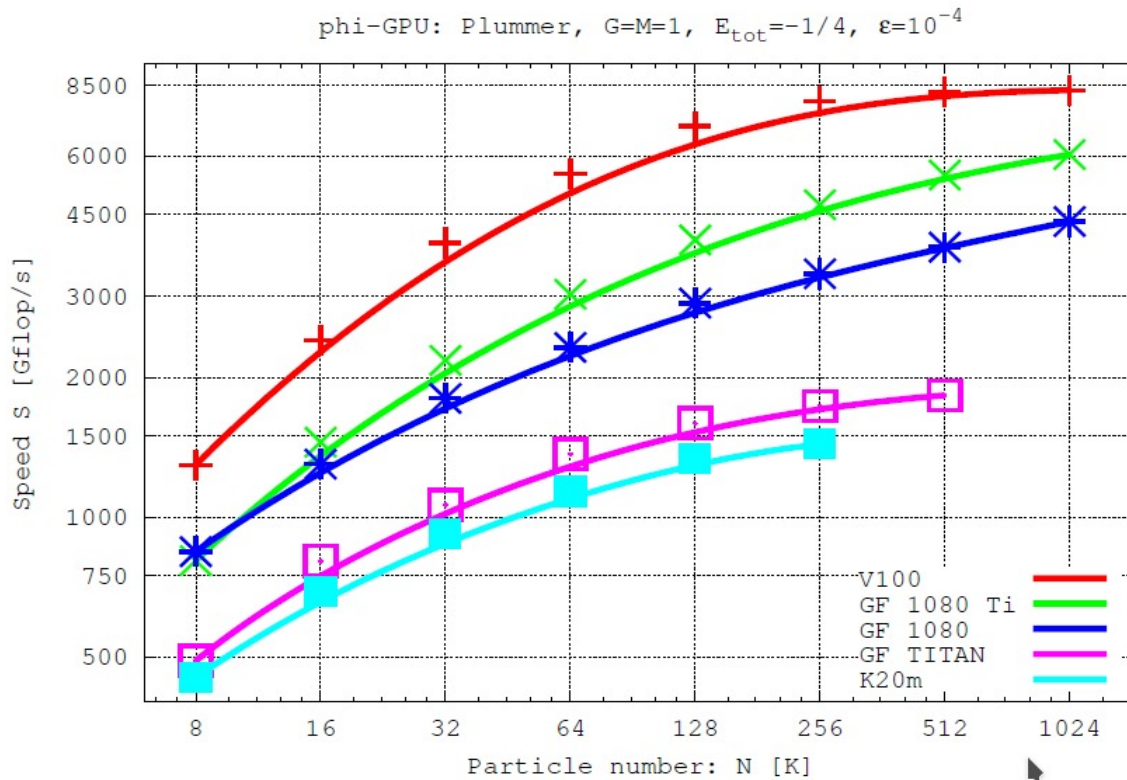


## Theoretical Peak GB/s





# Kepler, Pascal, Volta, Scaling, it works...



Volta V100

Pascal GF1080

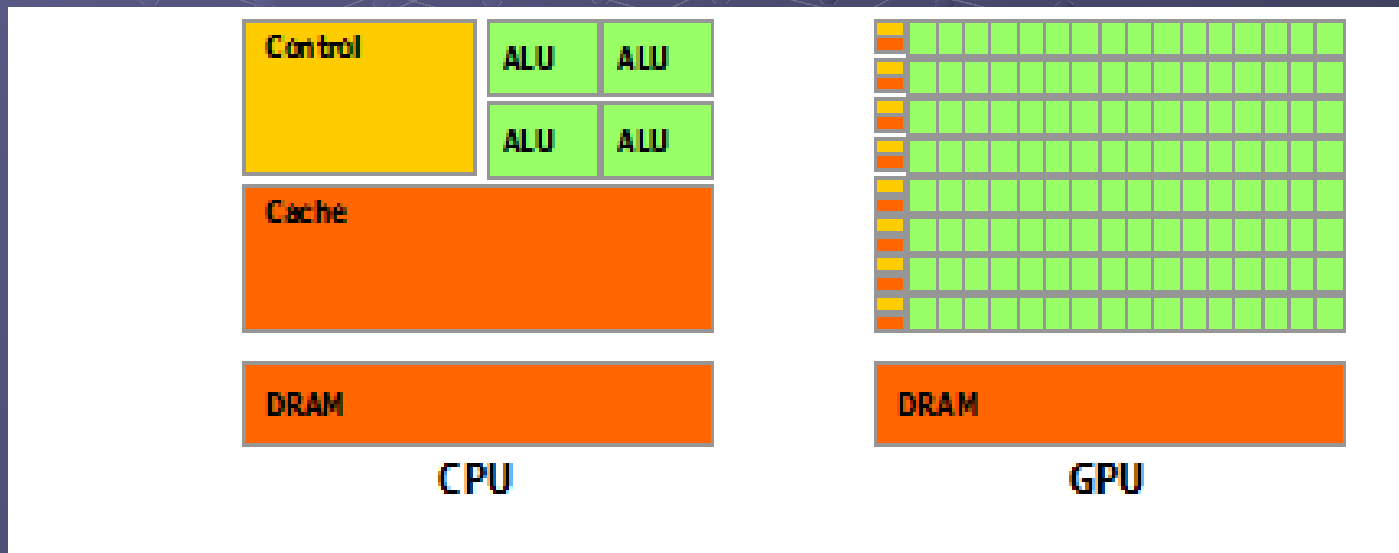
Kepler K20m

Spurzem, Berczik,  
et al., 2013,  
LNCS Supercomputing,  
2013, pp. 13-25,  
Springer.  
(updated unpublished)

**Fig. 4.** Here we report a preliminary result from a benchmark test of our code on one Kepler K20 card; we compare with the performance on Fermi C2050 (used in the Mole-8.5 cluster), and the oldest Tesla C1060 GPU (used in the laohu cluster of 2009) - the latter is used as a normalization reference. We plot the speed ratio of our usual benchmarking simulation used in the previous figures, as a function of particle number. From this we see the sustained performance of a Kepler K20 would be about 1.4 - 1.5 Tflop/s.

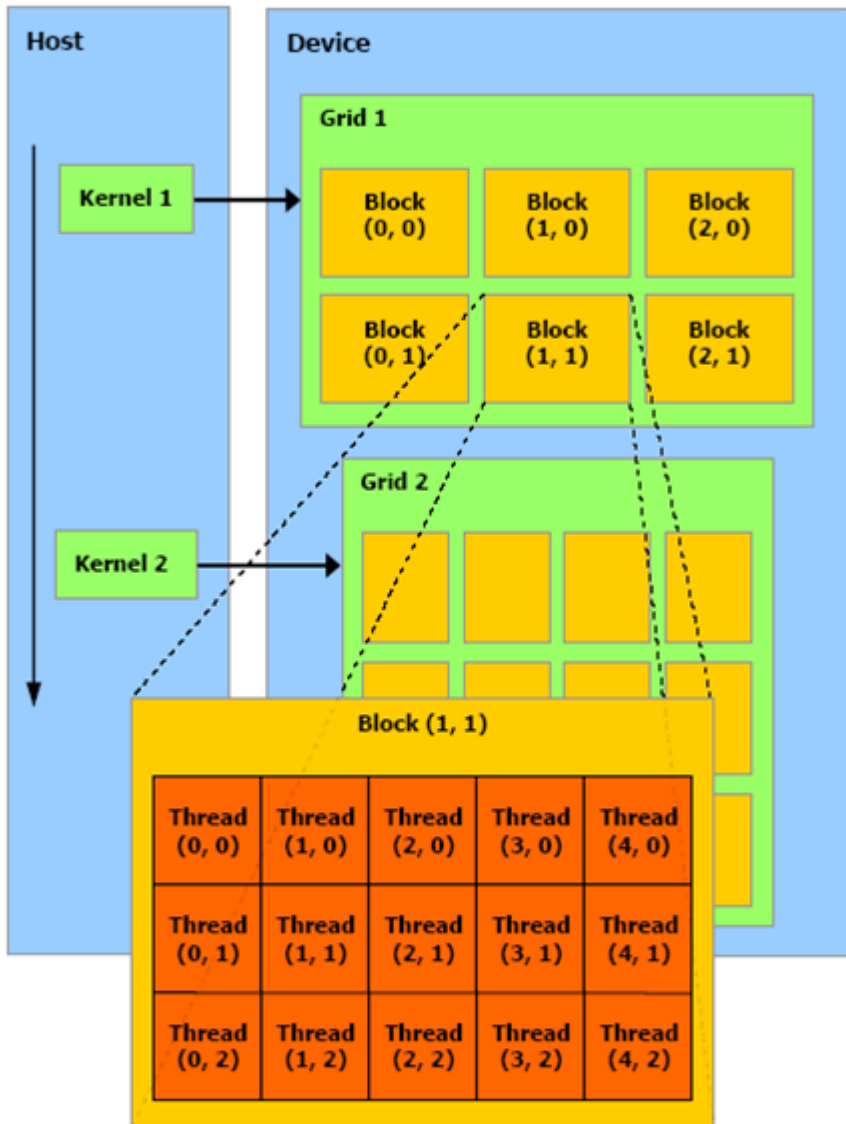
**X = first GPU of laohu 2010**

# CPU and GPU; from CUDA NVIDIA Developer Zone at <http://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html>

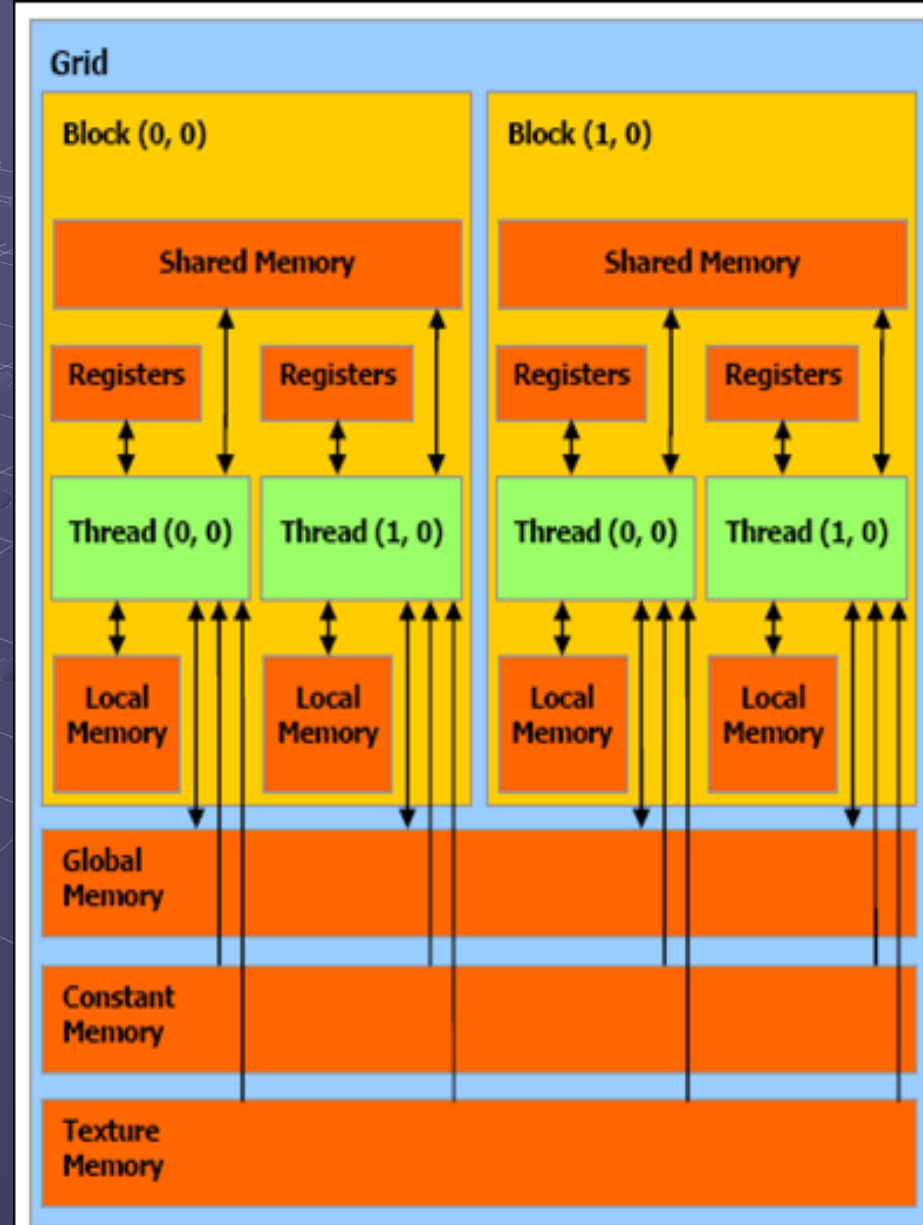


**“The GPU devotes more transistors to computing”  
“favours data parallel operations”**

# GPU Structure From: [http://geco.mines.edu/tesla/cuda\\_tutorial\\_mio/](http://geco.mines.edu/tesla/cuda_tutorial_mio/)

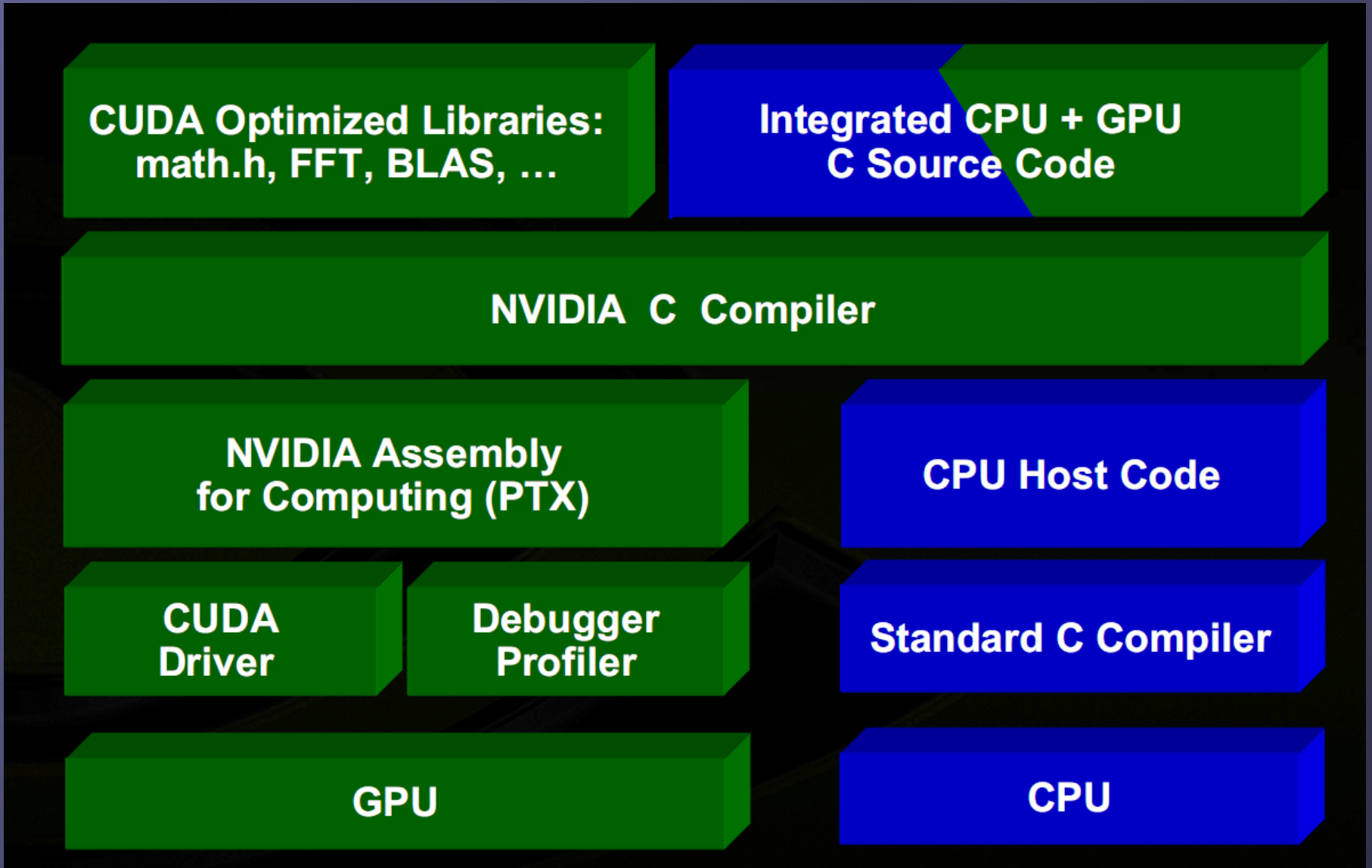


The host issues a succession of kernel invocations to the device. Each kernel is executed as a batch of threads organized as a grid of thread blocks





# CUDA



# Simple CUDA example

## CPU C program

```
void addMatrix(float *a, float *b,
              float *c, int N)
{
    int i, j, index;
    for (i = 0; i < N; i++) {
        for (j = 0; j < N; j++) {
            index = i + j * N;
            c[index]=a[index] + b[index];
        }
    }
}

void main()
{
    .....
    addMatrix(a, b, c, N);
}
```

## CUDA C program

```
__global__ void addMatrix(float *a, float *b,
                          float *c, int N)
{
    int i=blockIdx.x*blockDim.x+threadIdx.x;
    int j=blockIdx.y*blockDim.y+threadIdx.y;
    int index = i + j * N;
    if ( i < N && j < N)
        c[index]= a[index] + b[index];
}

void main()
{
    ..... // allocate & transfer data to GPU
    dim3 dimBlk (blocksize, blocksize);
    dim3 dimGrd (N/dimBlk.x, N/dimBlk.y);
    addMatrix<<<dimGrd,dimBlk>>>(a, b, c,N);
}
```

# GPU Computing Applications

Source: <http://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html>

## Libraries and Middleware

cuDNN TensorRT	cuFFT, cuBLAS, cuRAND, cuSPARSE	CULA MAGMA	Thrust NPP	VSIPL, SVM, OpenCurrent	PhysX, OptiX, iRay	MATLAB Mathematica
-------------------	------------------------------------	------------	---------------	----------------------------	-----------------------	-----------------------

## Programming Languages

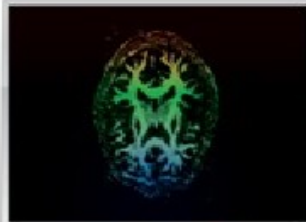
C	C++	Fortran	Java, Python, Wrappers	DirectCompute	Directives (e.g., OpenACC)
---	-----	---------	---------------------------	---------------	-------------------------------

## CUDA-enabled NVIDIA GPUs

Turing Architecture (Compute capabilities 7.x)	DRIVE/JETSON AGX Xavier	GeForce 2000 Series	Quadro RTX Series	Tesla T Series
Volta Architecture (Compute capabilities 7.x)	DRIVE/JETSON AGX Xavier			Tesla V Series
Pascal Architecture (Compute capabilities 6.x)	Tegra X2	GeForce 1000 Series	Quadro P Series	Tesla P Series
Maxwell Architecture (Compute capabilities 5.x)	Tegra X1	GeForce 900 Series	Quadro M Series	Tesla M Series
Kepler Architecture (Compute capabilities 3.x)	Tegra K1	GeForce 700 Series GeForce 600 Series	Quadro K Series	Tesla K Series
	EMBEDDED	CONSUMER DESKTOP, LAPTOP	PROFESSIONAL WORKSTATION	DATA CENTER

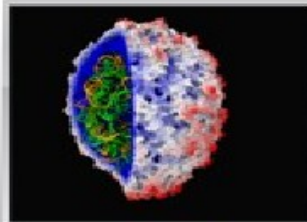


# Speedups using GPU vs. CPU



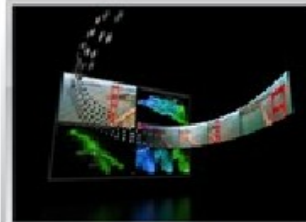
**146X**

Interactive visualization of volumetric white matter connectivity<sup>1</sup>



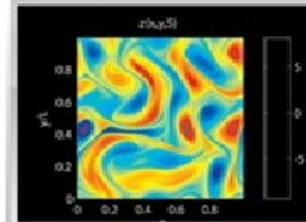
**36X**

Ionic placement for molecular dynamics simulation on GPU<sup>2</sup>



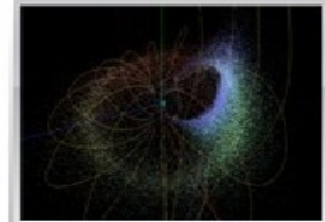
**18X**

Transcoding HD video stream to H.264 for portable video<sup>3</sup>



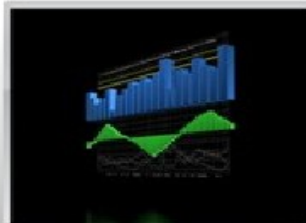
**17X**

Simulation in Matlab using mex file CUDA function<sup>4</sup>



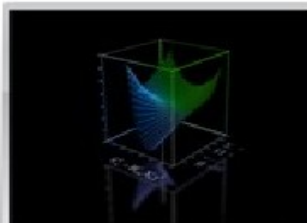
**100X**

Astrophysics N-body simulation<sup>5</sup>



**149X**

Financial simulation of LIBOR model with swaptions<sup>6</sup>



**47X**

GLAME@lab: M-script API for linear Algebra operations on GPU<sup>7</sup>



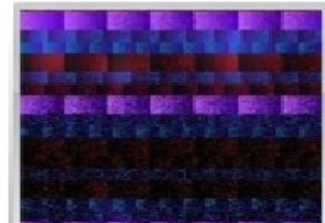
**20X**

Ultrasound medical imaging for cancer diagnostics<sup>8</sup>



**24X**

Highly optimized object oriented molecular dynamics<sup>9</sup>



**30X**

Cmatch exact string matching - find similar proteins & gene sequences<sup>10</sup>



# Towards Peta-Scale Green Computation

— applications of the GPU supercomputers in CAS

<http://www.nvidia.com/gtc2010-content>



## GPU TECHNOLOGY CONFERENCE

GTC 2010 | Sept 20-23, 2010

San Jose Convention Center, San Jose, California

Watch the Keynote Recordings

Algorithms & Numerical Techniques

Astronomy & Astrophysics

Audio Processing

Cloud Computing

Computational Fluid Dynamics

Computer Graphics

Computer Vision

Databases & Data Mining

Digital Content Creation

Embedded & Automotive

Energy Exploration

Film

Finance

General Interest

GPU Accelerated Internet

High Performance Computing

Imaging

Life Sciences

Machine Learning & Artificial

Intelligence

Medical Imaging & Visualization

Mobile & Tablet & Phone

Molecular Dynamics

Neuroscience

Physics Simulation

Programming Languages & Techniques

Quantum Chemistry

Ray Tracing

Signal Processing

Stereoscopic 3D

Tools & Libraries

Video Processing

Wei Ge  
Xiaowei Wang

Inst. of Proc. Eng.



Yunquan Zhang  
Inst. of Software



Rainer Spurzem  
Nat. Astro. Obs.  
Chn.



Long Wang  
SC Center

