

the SILK ROAD PROJECT at NAOC/

丝绸之路计划



ZENTRUM FÜR ASTRONOMIE

Univ. Heidelberg

UNIVERSITÄT HEIDELBERG  
Zukunft. Seit 1386.



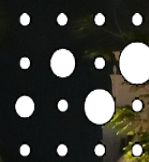
# Introduction to GPU Accelerated Computing

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

<https://astro-silkroad.eu>

<https://wwwstaff.ari.uni-heidelberg.de/spurzem/>



VolkswagenStiftung

[spurzem@ari.uni-heidelberg.de](mailto:spurzem@ari.uni-heidelberg.de)

[spurzem@nao.cas.cn](mailto:spurzem@nao.cas.cn)

中国科学院国家天文台

National Astronomical Observatories, CAS

Picture: Xishuangbanna,  
Yunnan, China by R.Sp.



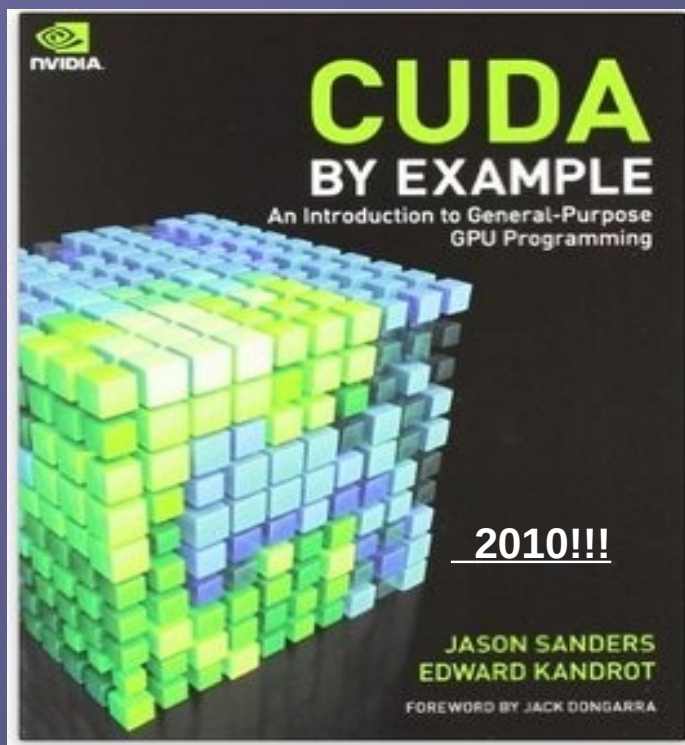
北京大学  
PEKING UNIVERSITY

# Introduction to GPU Accelerated Computing

Feb. 17-21, 2025

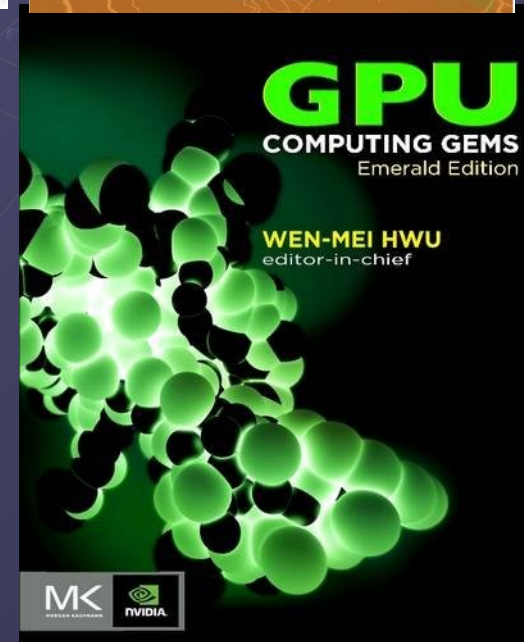
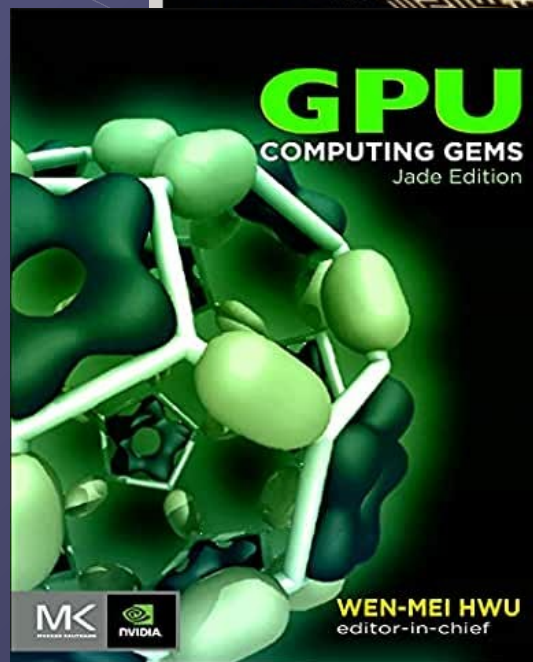
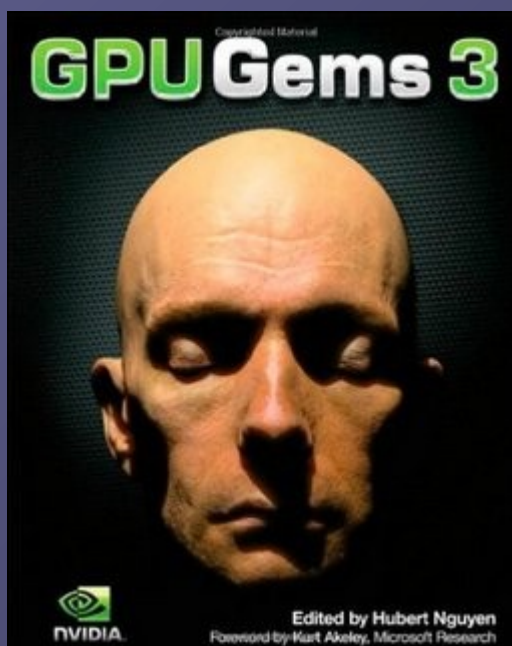
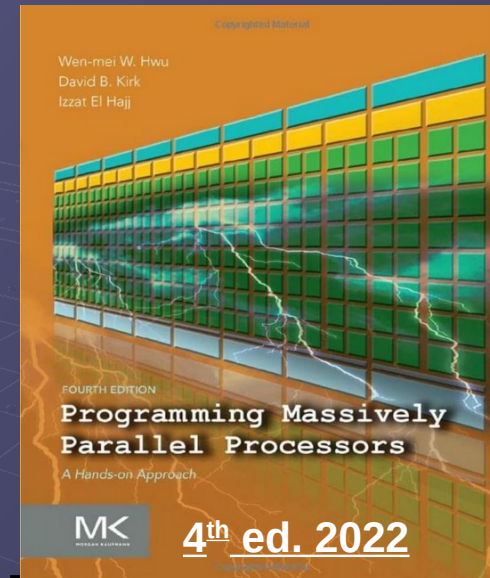
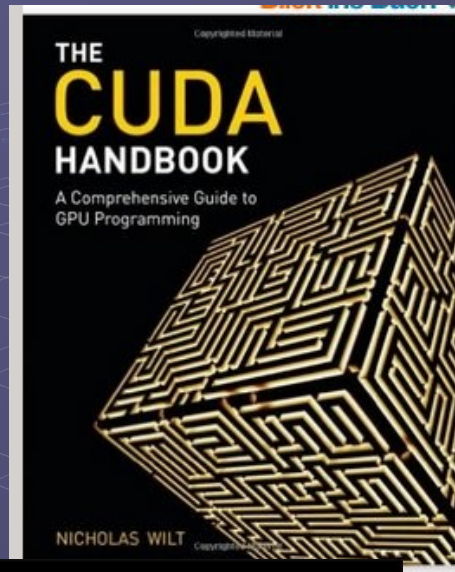
## Table of Contents (subject to adjustment/change):

1. Monday morning 1: General Introduction Computer Architecture, Many-Core, GPU and others..., Access...
2. Monday morning 2/afternoon: Access to bwUniCluster, CUDA Hello, GPU Properties, First CUDA Scalar, Simple Vector Add
3. Tuesday morning 1: More on GPU Software and Hardware
4. Tuesday morning 2/afternoon: CUDA Vector Add, Scalar Products, Using Blocks and Threads
5. Wednesday morning: Parallelization and Amdahl's Law, GPU Acceleration, Future Architecture
6. Wednesday morning 2/afternoon: CUDA Scalar Products cont'd Events, Histograms, Matrix Multiplication
7. Thursday Morning: Astrophysical N-Body Code
8. Thursday Afternoon: Astrophysical Parallel N-Body Code Using MPI and GPU
9. Friday Morning: CUDA Matrix Mult., Histograms, Wrap-Up, Q+A, Other Lectures (Wen-Mei Hwu)

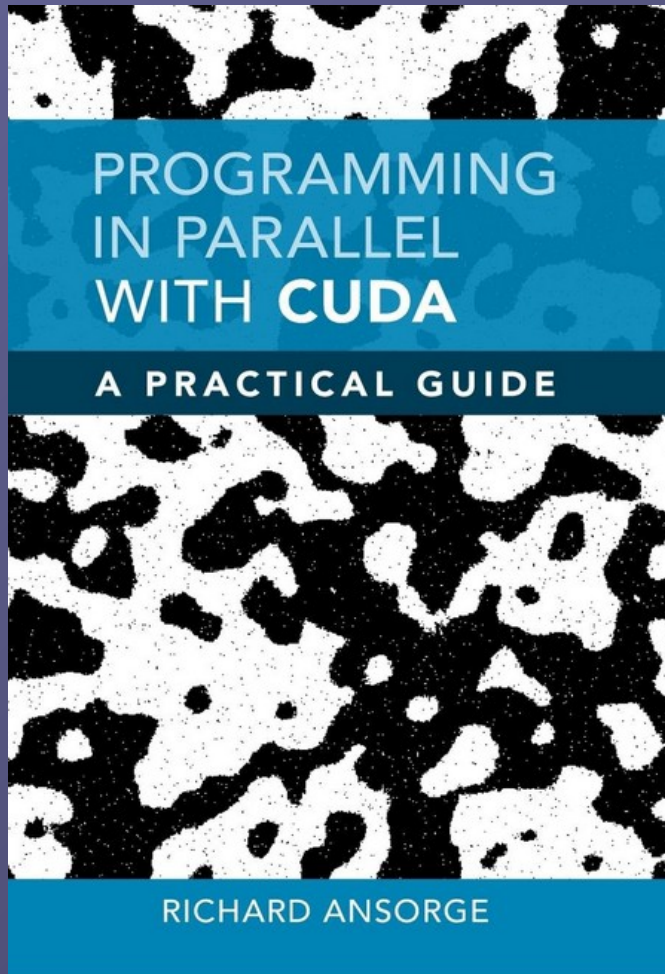


Literature: why NVIDIA? CUDA ... ?

easy to learn!! runs on our training system kepler  
future? SYCL/openCL? HIP / HIPIFY ?



# Literature continued:



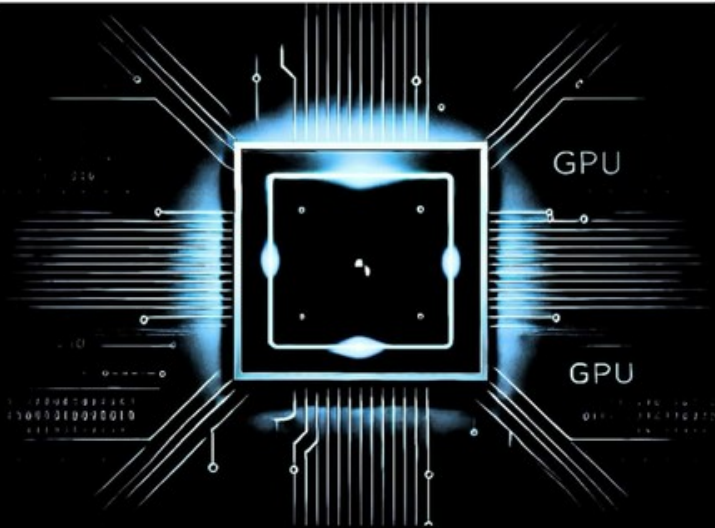
"CUDA is now the dominant language used for programming GPUs, one of the most exciting hardware developments of recent decades. With CUDA, you can use a desktop PC for work that would have previously required a large cluster of PCs or access to a HPC facility. As a result, CUDA is increasingly important in scientific and technical computing across the whole STEM community, from medical physics and financial modelling to big data applications and beyond. This unique book on CUDA draws on the author's passion for and long experience of developing and using computers to acquire and analyse scientific data. The result is an innovative text featuring a much richer set of examples than found in any other comparable book on GPU computing. Much attention has been paid to the C++ coding style, which is compact, elegant and efficient. A code base of examples and supporting material is available online, which readers can build on for their own projects"--

New Book of 2022, Text from Book advertisement in amazon.

# Literature continued:

## GPU

### Architecture



The Ultimate Guide to Building High-  
Performance Computing Systems

Cobbs Walker

“GPU systems have revolutionized the fields of artificial intelligence, data science, and high-performance computing. Their unparalleled ability to handle massive parallel processing tasks has made them indispensable for industries that rely on cutting-edge computational power. From AI model training to scientific simulations and beyond, understanding how to design and optimize GPU architectures is key to maximizing performance and staying ahead in a rapidly evolving tech landscape.”

“Authored by a high-performance computing expert, (Cobbs Walker) provides the most up-to-date, actionable insights on GPU system design. This book is based on years of hands-on experience building and optimizing GPU infrastructures, paired with real-world case studies that demonstrate successful implementations. Whether you're designing AI systems, working on complex simulations, or building GPU-driven applications, the expertise shared here is reliable and practical.”

For deeper interest in GPU hardware.

Text from Book advertisement in amazon.

# Introduction to GPU Accelerated Computing

Feb. 17-21, 2025

## “Table of Contents” what we will NOT cover:

- Artificial Intelligence / Machine Learning
- Graphics Rendering / Ray Tracing with GPU
- Using Tensor Cores for 3D simulations
- Other Languages such as HIP (AMD), OpenCL...

We will solely use CUDA for High Performance Computing on GPU (many simple examples, one real Application).

What you learn here will give you a good start for all the applications not covered!

# 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**

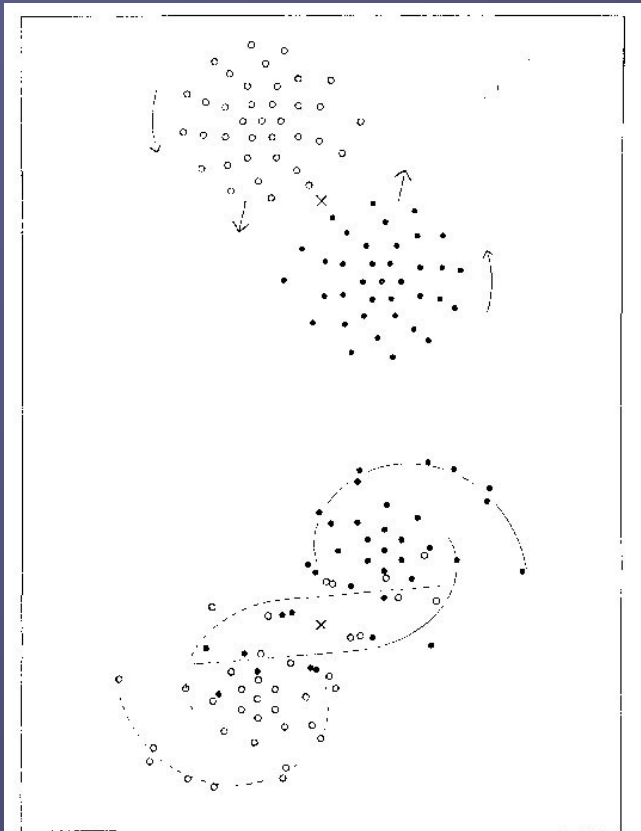


FIG. 4b

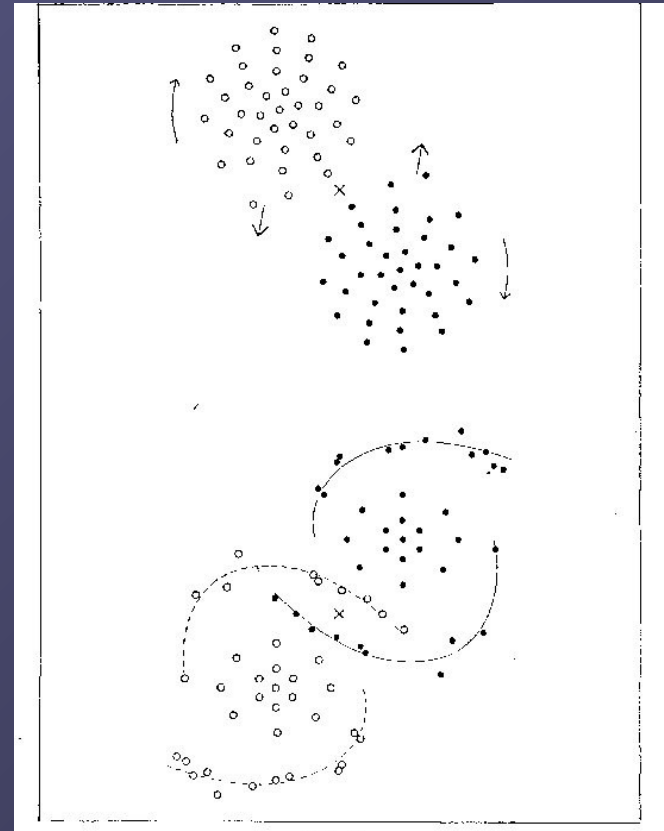


FIG. 4a

# HARDWARE

...before von Neumann...

● Konrad Zuse (1910-1995) Berlin

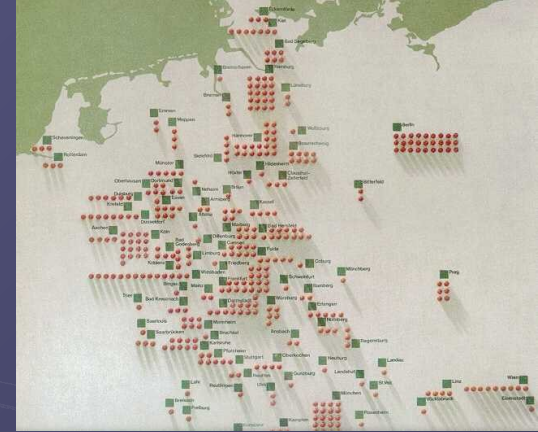


Invented freely programmable Computer



Z1 in parental flat 1936

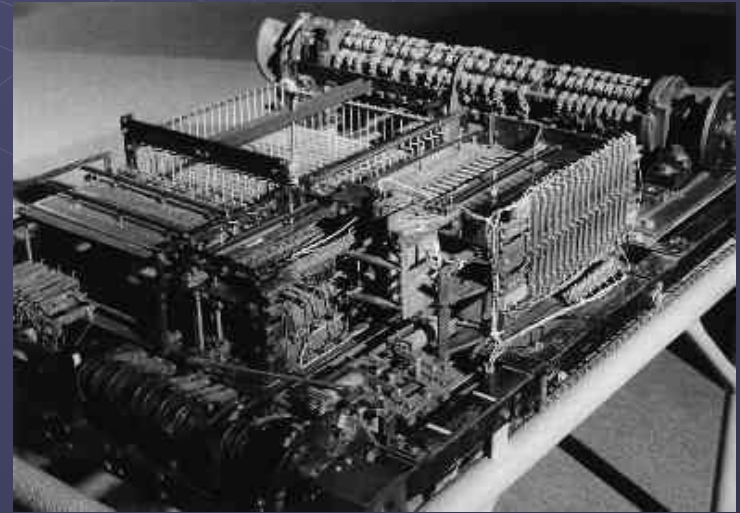
# History



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



**Computing Speed 0.03 MHz**



**Memory 256 byte**

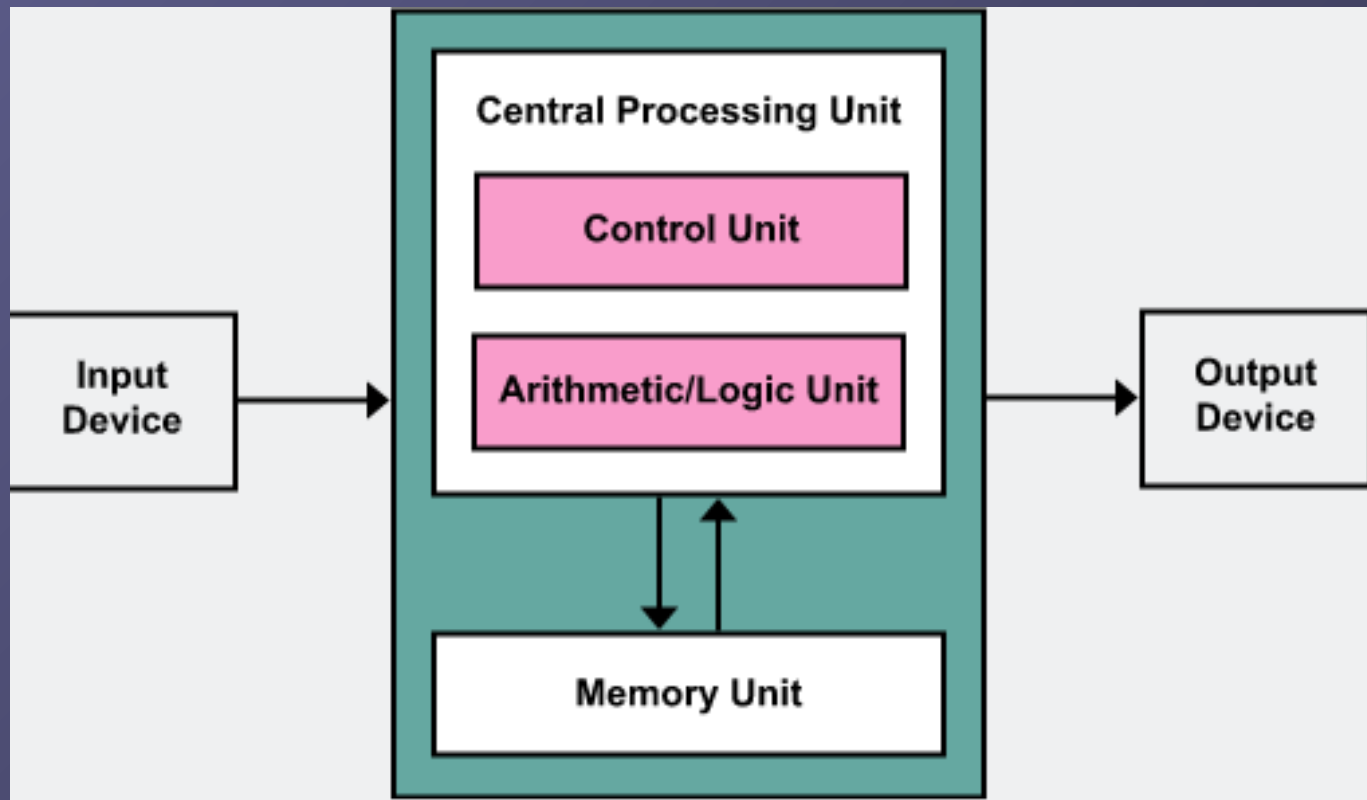
# HARDWARE

- John von Neumann (1903-1957)

Born Budapest, Lecturer Berlin, since 1930 Princeton Univ.

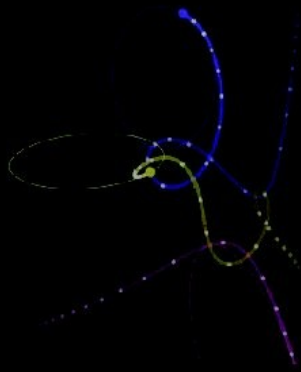
Fundamental Architecture of an electronic computing device(1946)

Source: [https://en.wikipedia.org/wiki/Von\\_Neumann\\_architecture#/media/File:Von\\_Neumann\\_Architecture.svg](https://en.wikipedia.org/wiki/Von_Neumann_architecture#/media/File:Von_Neumann_Architecture.svg)





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)

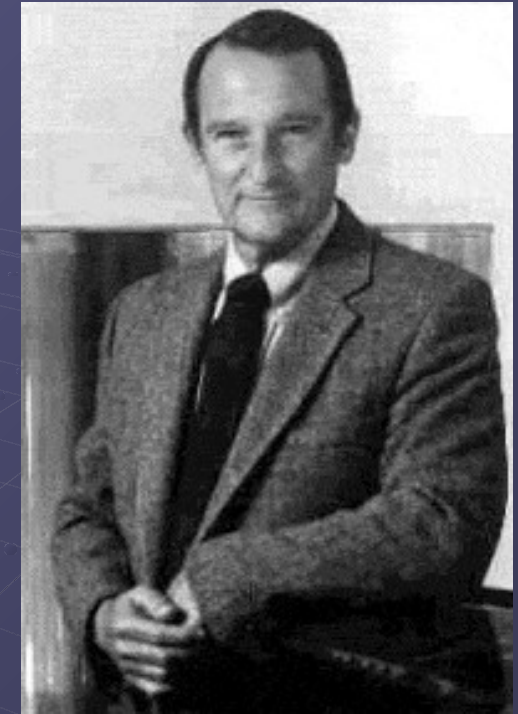
N=16,25 (40 Trx)

# History

## ● Seymour Cray (1925-1996)

“father of supercomputing”

[https://en.wikipedia.org/wiki/Women\\_in\\_computing](https://en.wikipedia.org/wiki/Women_in_computing)



**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

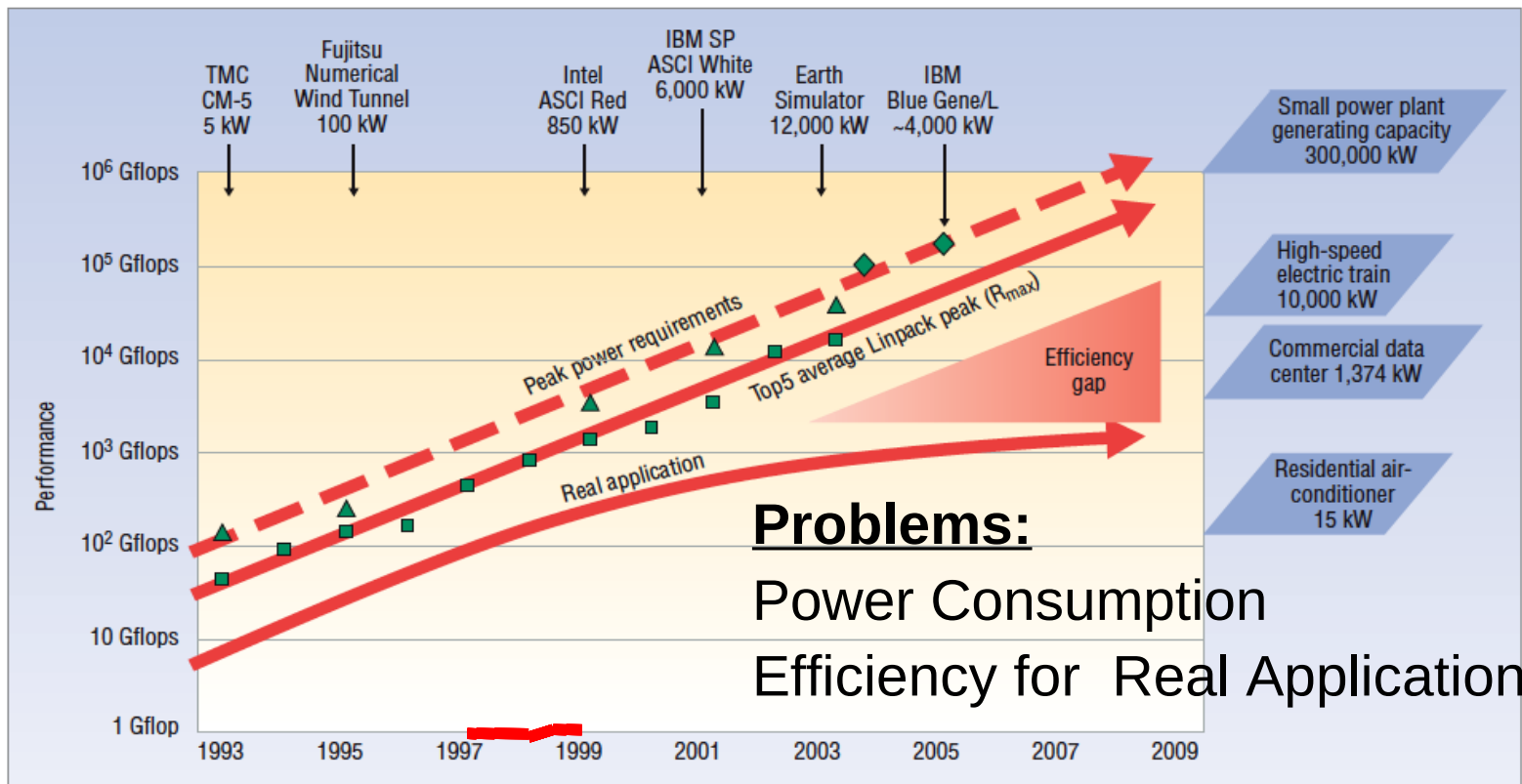


Figure 1. Rising power requirements. Peak power consumption of the top supercomputers has steadily increased over the past 15 years. Thanks to Horst Simon, LBNL/NERSC for this diagram.

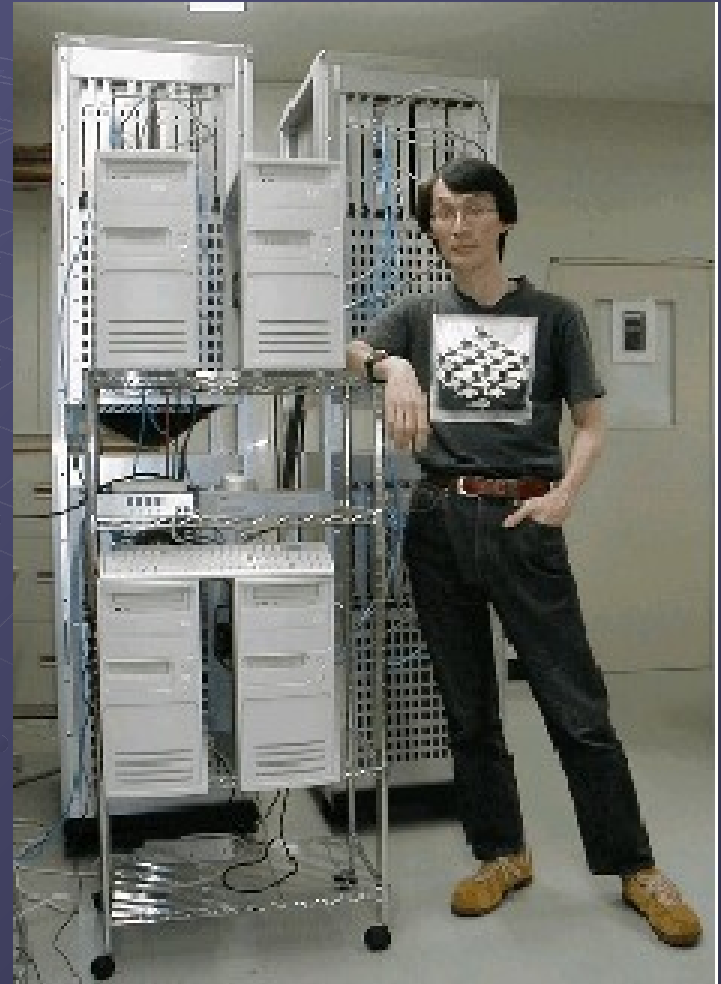
GPU Computing

# Special Hardware Accelerators

# 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



# GPU: NAOC laohu cluster Beijing, China



# BwUniCluster 2.0

The bwUniCluster 2.0 is the joint high-performance computer system of Baden-Württemberg's Universities and Universities of Applied Sciences for general purpose and teaching and located at the Scientific Computing Center (SCC) at Karlsruhe Institute of Technology (KIT). The bwUniCluster 2.0 complements the four bwForClusters and their dedicated scientific areas.

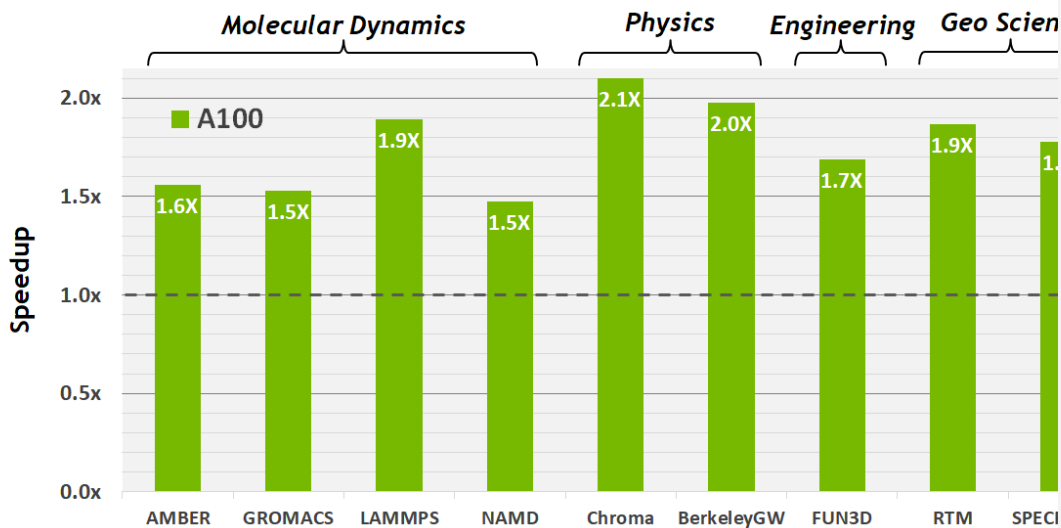


Total Number of Nodes: 848

GPU Nodes: 39 (NVIDIA Ampère A100, Volta V100)

# NVIDIA Ampere A100 GPU, 54 billion transistors, 6920 cores

## ACCELERATING HPC



All results are measured  
Except BerkeleyGW, V100 used is single V100 SXM2. A100 used is single A100 SXM4  
More apps detail: AMBER based on PME-Cellulose, GROMACS with STMV (h-bond), LAMMPS with Atomic Fluid LJ-2.5, NAMD with v3.0a1 STMV\_NVE  
Chroma with szsc121\_24\_128, FUN3D with dpw, RTM with Isotropic Radius 4 1024<sup>3</sup>, SPECFEM3D with Cartesian four material model  
BerkeleyGW based on Chi Sum and uses 8xV100 in DGX-1, vs 8xA100 in DGX A100

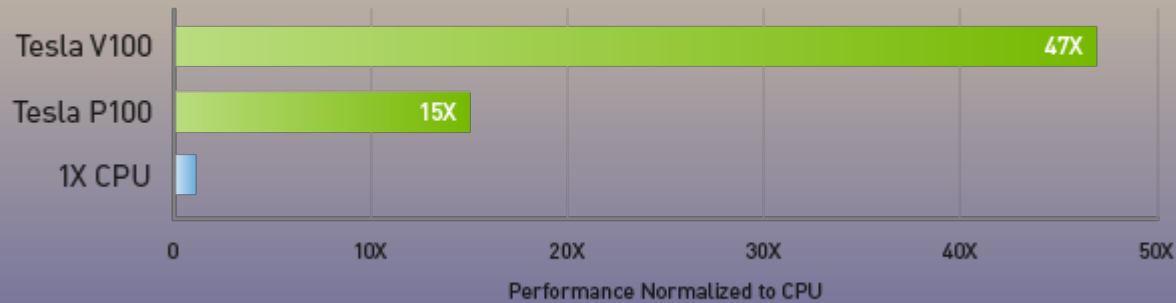
## 11X More HPC Performance in Four Years

### Top HPC Apps



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

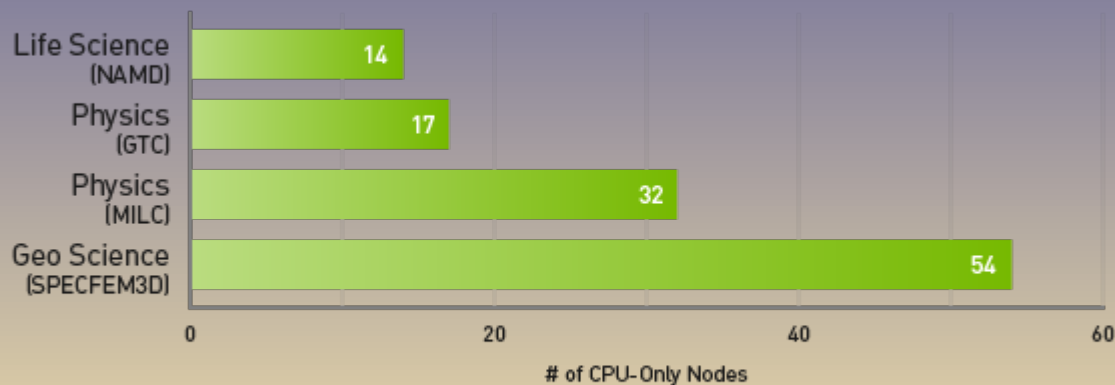
## 47X Higher Throughput Than CPU Server on Deep Learning Inference



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.

# NVIDIA Ampere A100 GPU, 54 billion transistors, 6920 cores (Hopper H100, ...)



With NVLINK

Without NVLINK



	A100 80GB PCIe	A100 80GB SXM
FP64	9.7 TFLOPS	
FP64 Tensor Core	19.5 TFLOPS	
FP32	19.5 TFLOPS	
Tensor Float 32 (TF32)	156 TFLOPS   312 TFLOPS*	
BFLOAT16 Tensor Core	312 TFLOPS   624 TFLOPS*	
FP16 Tensor Core	312 TFLOPS   624 TFLOPS*	
INT8 Tensor Core	624 TOPS   1248 TOPS*	
GPU Memory	80GB HBM2e	80GB HBM2e
GPU Memory Bandwidth	1,935 GB/s	2,039 GB/s
Max Thermal Design Power (TDP)	300W	400W ***



# AMD Instinct™ MI250X

## GPU Specifications

**GPU Architecture:** CDNA2

**Stream Processors:** 14,080

**Peak Half Precision (FP16) Performance:**  
383 TFLOPs

**Peak Single Precision Matrix (FP32) Performance:**  
95.7 TFLOPs

**Peak Single Precision (FP32) Performance:**  
47.9 TFLOPs

**Peak INT4 Performance:** 383 TOPs

**Peak bfloat16:** 383 TFLOPs

**Lithography:** TSMC 6nm FinFET

**Compute Units:** 220

**Peak Engine Clock:** 1700 MHz

**Peak Double Precision Matrix (FP64) Performance:**  
95.7 TFLOPs

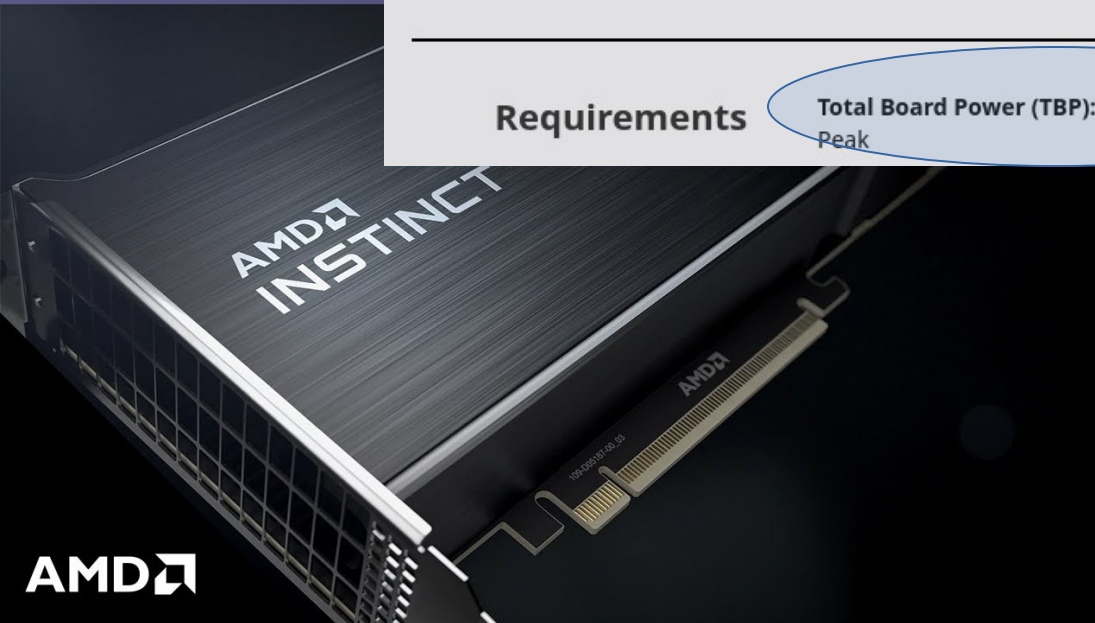
**Peak Double Precision (FP64) Performance:**  
47.9 TFLOPs

**Peak INT8 Performance:** 383 TOPs

**OS Support:** Linux x86\_64

## Requirements

**Total Board Power (TBP):** 500W | 560W Peak

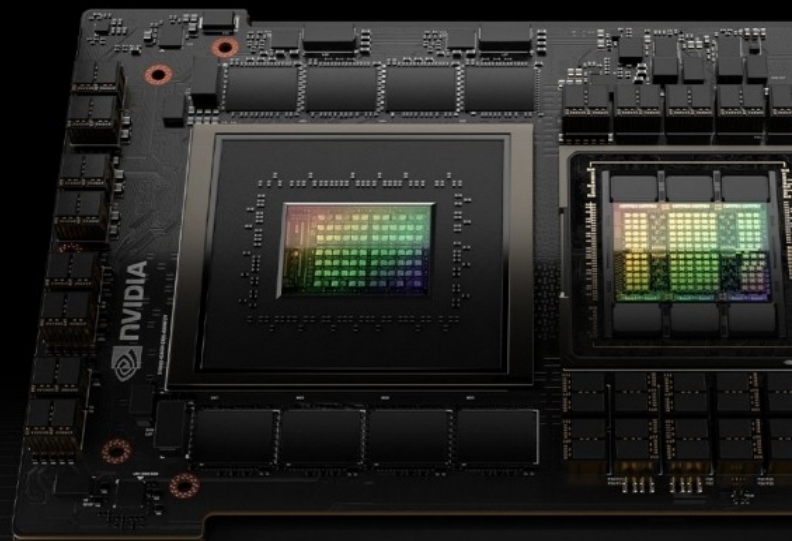


AMD Instinct  
MI250X GPU

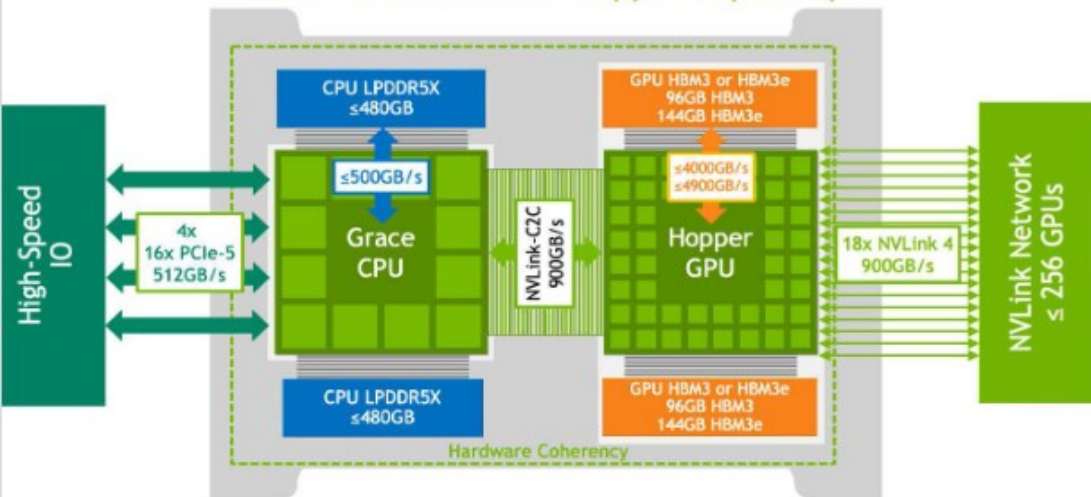
Nov.2023 Lists:  
Used in:

Frontier (#1 US)  
And LUMI (#5 FIN)

New “Grace Hopper GH200 superchip” ; GPU + CPU on one platform; used in new Jupiter supercomputer at JSC Jülich.



### NVIDIA GH200 Grace Hopper Superchip



Hopper GPU  
16896 CUDA cores  
528 tensor cores  
34 Tflop/s double prec.  
67 Tflop/s single prec.  
67 Tflop/s tensor core  
double prec.

72 Armv9 CPU cores  
480 GB memory

<https://www.nvidia.com/en-us/data-center/grace-hopper-superchip/>



From <https://www.top500.org/>

Nov. 2023 List



USA



USA



USA



USA

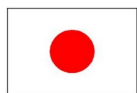


Italy

Rank	System	Cores	Rmax (PFlop/s)	Rpeak (PFlop/s)	Power (kW)
1	<b>El Capitan</b> - HPE Cray EX255a, AMD 4th Gen EPYC 24C 1.8GHz, AMD Instinct MI300A, Slingshot-11, TOSS, HPE DOE/NNSA/LLNL United States	11,039,616	1,742.00	2,746.38	29,551
		<u><i>GPU AMD Instinct</i></u>			
2	<b>Frontier</b> - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE Cray OS, HPE DOE/SC/Oak Ridge National Laboratory United States	9,066,176	1,353.00	2,055.72	24,607
		<u><i>GPU AMD Instinct</i></u>			
3	<b>Aurora</b> - HPE Cray EX - Intel Exascale Compute Blade, Xeon CPU Max 9470 52C 2.4GHz, Intel Data Center GPU Max, Slingshot-11, Intel DOE/SC/Argonne National Laboratory United States	9,264,128	1,012.00	1,980.01	38,698
		<u><i>Intel Data Center GPU</i></u>			
4	<b>Eagle</b> - Microsoft NDv5, Xeon Platinum 8480C 48C 2GHz, NVIDIA H100, NVIDIA Infiniband NDR, Microsoft Azure Microsoft Azure United States	2,073,600	561.20	846.84	
		<u><i>GPU NVIDIA Hopper</i></u>			
5	<b>HPC6</b> - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, RHEL 8.9, HPE Eni S.p.A. Italy	3,143,520	477.90	606.97	8,461
		<u><i>GPU AMD Instinct</i></u>			



From <https://www.top500.org/>  
Nov. 2024 List



**Japan**



**USA**



**Finland  
(EuroHPC)**



**USA**

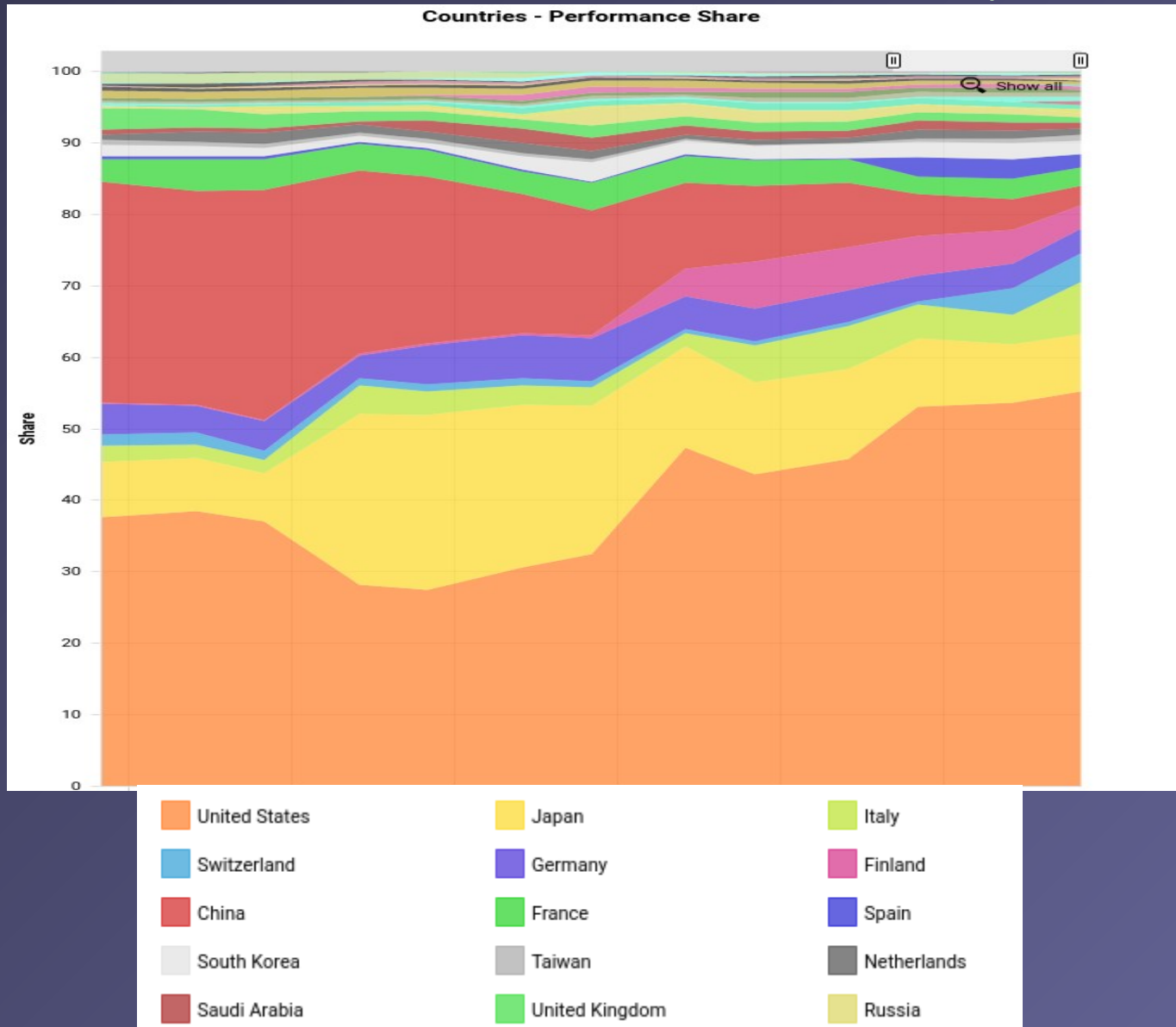


**USA**

Rank	System	Cores	Rmax (PFlop/s)	Rpeak (PFlop/s)	Power (kW)
6	<b>Supercomputer Fugaku</b> - Supercomputer Fugaku, A64FX 48C 2.2GHz, Tofu interconnect D, Fujitsu RIKEN Center for Computational Science Japan	7,630,848	442.01	537.21	29,899
<b><u>Fujitsu Arm</u></b>					
7	<b>Alps</b> - HPE Cray EX254n, NVIDIA Grace 72C 3.1GHz, NVIDIA GH200 Superchip, Slingshot-11, HPE Cray OS, HPE Swiss National Supercomputing Centre (CSCS) Switzerland	2,121,600	434.90	574.84	7,124
<b><u>GPU NVIDIA GH200 Grace</u></b>					
8	<b>LUMI</b> - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE EuroHPC/CSC Finland	2,752,704	379.70	531.51	7,107
<b><u>GPU AMD Instinct</u></b>					
9	<b>Leonardo</b> - BullSequana XH2000, Xeon Platinum 8358 32C 2.6GHz, NVIDIA A100 SXM4 64 GB, Quad-rail NVIDIA HDR100 Infiniband, EVIDEN EuroHPC/CINECA Italy	1,824,768	241.20	306.31	7,494
<b><u>GPU NVIDIA Ampere</u></b>					
10	<b>Tuolumne</b> - HPE Cray EX255a, AMD 4th Gen EPYC 24C 1.8GHz, AMD Instinct MI300A, Slingshot-11, TOSS, HPE DOE/NNSA/LLNL United States	1,161,216	208.10	288.88	3,387
<b><u>GPU AMD Instinct</u></b>					

# Top 500 List November 2023 –

## Performance Share of Countries [From https://www.top500.org](https://www.top500.org)



# LUMI

## Supercomputer, Kajaani, Finland

Using only  
Hydroelectric  
Power and its  
Heat used for  
heating buildings.

No. 5 in top500  
No. 7 in green500

2.2 million cores  
~12.000 AMD GPUs



EuroHPC and LUMI consortium:  
Finland, Belgium, Czech Republic, Denmark, Estonia,  
Iceland, Norway, Poland, Sweden, and Switzerland.

RIKEN, Kobe, JAPAN

# FUGAKU



*Nature's Secrets*

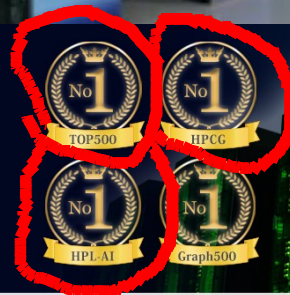
# 富岳

Mt. Fuji

## The world's fastest Super Computer 2020 /2021

7.6 million cores, 442 Pflop/s

source :nytimes



Fugaku extends its reign as champion of supercomputers

JUWELS Booster 936 nodes (AMD CPU, 4x Ampere GPU)  
~450.000 AMD cores, 25 million NVIDIA Ampere GPU cores  
~ 70 Pflop/s SP ~ 44 Pflop/s DP  
No. 18 in top500 list, No. 3 in green500 list

Jülich Wizard for European Leadership Science



Watch out for new Exascale System at Jülich (JSC): JEDI / JUPITER !



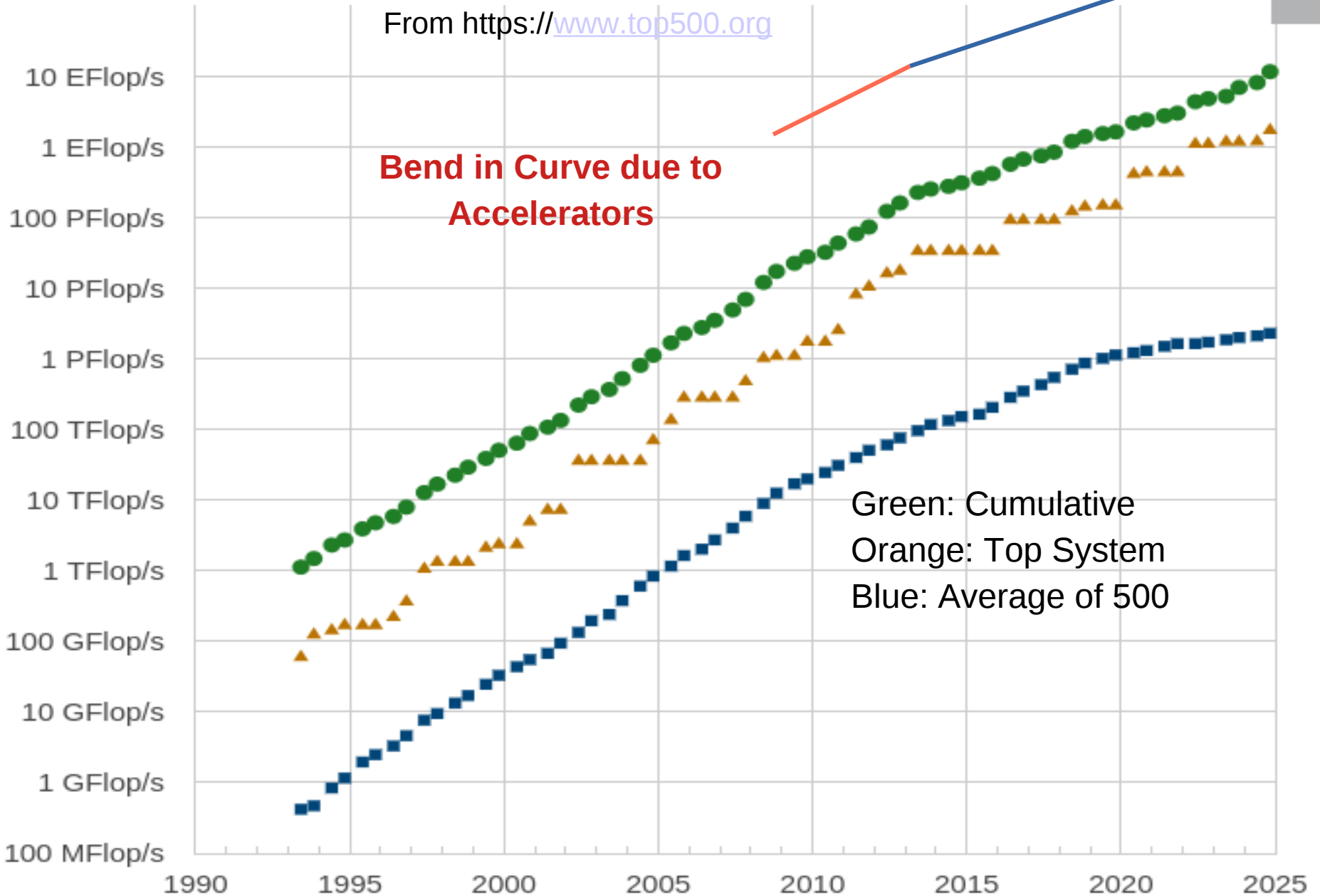
# Performance Development

From <https://www.top500.org>

Moore's Law?

**Bend in Curve due to Accelerators**

Green: Cumulative  
Orange: Top System  
Blue: Average of 500



# GREEN 500 list Nov. 2024

Power Efficiency

(Gflops/Watts),

see also top500 webpage

right: 1-5

below: 6-10

									Energy
									Efficiency
Rank	TOP500 Rank	System	Cores	Rmax (PFlop/s)	Power (kW)	GFlops/watts			
	1	<b>JEDI</b> - BullSequana XH3000, Grace Hopper Superchip 72C 3GHz, NVIDIA GH200 Superchip, Quad-Rail NVIDIA InfiniBand NDR200, ParTec/EVIDEN EuroHPC/FZJ Germany	19,584	4.50	67	72.733			<u><b>GPU NVIDIA Grace Hopper</b></u>
	2	<b>ROMEO-2025</b> - BullSequana XH3000, Grace Hopper Superchip 72C 3GHz, NVIDIA GH200 Superchip, Quad-Rail NVIDIA InfiniBand NDR200, Red Hat Enterprise Linux, EVIDEN ROMEO HPC Center - Champagne-Ardenne France	47,328	9.86	160	70.912			<u><b>GPU NVIDIA Grace Hopper</b></u>
6	18	<b>JETI - JUPITER Exascale Transition Instrument</b> - BullSequana XH3000, Grace Hopper Superchip 72C 3GHz, NVIDIA GH200 Superchip, Quad-Rail NVIDIA InfiniBand NDR200, RedHat Linux and Modular Operating System, ParTec/EVIDEN EuroHPC/FZJ Germany	391,680	83.14	1,311	67.9			<u><b>GPU NVIDIA Grace Hopper</b></u>
7	69	<b>Helios GPU</b> - HPE Cray EX254n, NVIDIA Grace 72C 3.1GHz, NVIDIA GH200 Superchip, Slingshot-11, HPE Cyfronet Poland	39,950	11.11	657	66.9			<u><b>GPU NVIDIA Grace Hopper</b></u>
8	369	<b>Henri</b> - ThinkSystem SR670 V2, Intel Xeon Platinum 8362 32C 2.8GHz, NVIDIA H100 80GB PCIe, Infiniband HDR, Lenovo Flatiron Institute United States	8,288	2.88	44	65.3			<u><b>GPU NVIDIA Hopper</b></u>
9	338	<b>HoreKa-Teal</b> - ThinkSystem SD665-N V3, AMD EPYC 9354 32C 3.25GHz, Nvidia H100 94Gb SXM5, Infiniband NDR200, Lenovo Karlsruhe Institut für Technologie (KIT) Germany	13,616	3.12	50	62.9			<u><b>GPU NVIDIA Hopper</b></u>
10	49	<b>rzAdams</b> - HPE Cray EX255a, AMD 4th Gen EPYC 24C 1.8GHz, AMD Instinct MI300A, Slingshot-11, TOSS, HPE DOE/NNSA/LLNL United States	129,024	24.38	388	62.8			<u><b>GPU AMD Instinct</b></u>
	3	<b>Adastra 2</b> - HPE Cray EX255a, AMD 4th Gen EPYC 24C 1.8GHz, AMD Instinct MI300A, Slingshot-11, RHEL, HPE Grand Equipement National de Calcul Intensif - Centre Informatique National de l'Enseignement Suprieur (GENCI-CINES) France	16,128	2.53	37	69.098			<u><b>GPU AMD Instinct</b></u>
	4	<b>Isambard-AI phase 1</b> - HPE Cray EX254n, NVIDIA Grace 72C 3.1GHz, NVIDIA GH200 Superchip, Slingshot-11, HPE University of Bristol United Kingdom	34,272	7.42	117	68.835			<u><b>GPU NVIDIA Grace Hopper</b></u>
	5	<b>Capella</b> - Lenovo ThinkSystem SD665-N V3, AMD EPYC 9334 32C 2.7GHz, Nvidia H100 SXM5 94Gb, Infiniband NDR200, AlmaLinux 9.4, MEGWARE TU Dresden, ZIH Germany	85,248	24.06	445	68.053			<u><b>GPU NVIDIA Hopper</b></u>