

# Introduction

CS121 Parallel Computing  
Fall 2022

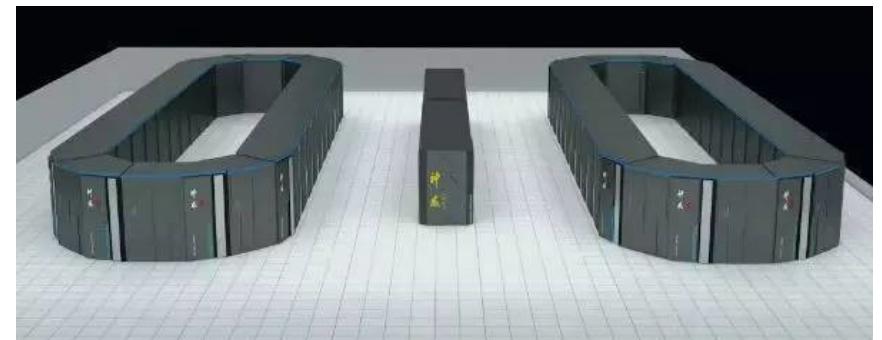
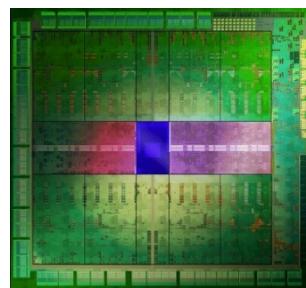
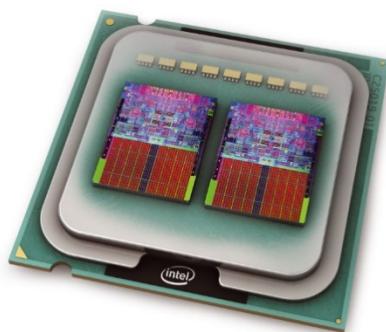
# Course info

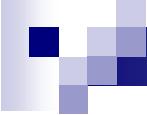
- Instructor Assoc Prof Rui FAN 范睿
- Research Parallel and distributed computing
- Contact fanrui@shanghaitech.edu.cn (English please)
- Office hours Thursdays 4:40-6pm, SIST 1A-504E
- TA 赖余睿, laiyr@shanghaitech.edu.cn
- Recitation TBA
- Website Blackboard and Piazza

Problem sets	20%	<ul style="list-style-type: none"><li>▪ About once every 2 weeks</li></ul>
Labs	20%	<ul style="list-style-type: none"><li>▪ Solve problems using OpenMP and CUDA</li></ul>
Reading project	15% Teams of 2	<ul style="list-style-type: none"><li>▪ Find an interesting research paper from suggested reading list</li><li>▪ Tell me your paper by week 8</li><li>▪ Submit a report and give a 20 minute presentation on reading OR programming project in week 16</li></ul>
Programming project	15% Teams of 2	<ul style="list-style-type: none"><li>▪ Find an interesting problem and write an efficient parallel program for it</li><li>▪ Tell me your problem by week 8</li><li>▪ Submit a report and give a 20 minute presentation on reading OR programming project in week 16</li></ul>
Midterm exam	10%	<ul style="list-style-type: none"><li>▪ At start of week 9</li></ul>
Final exam	20%	

# Parallel computing: what and why

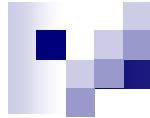
- Parallel computing studies how to use multiple computers together to solve a problem.
- Allows solving complicated problems faster.
  - Ideally, with  $k$  processors we can solve a problem  $k$  times faster.
  - Also more memory to solve larger problems, or same problem with more accuracy.
  - May be more fault tolerant; but also more prone to faults.
- Almost all modern computer systems are parallel.
  - Multicores, GPUs, cloud computing, etc.
- Parallel computing crucial for modern large scale applications, e.g. physical simulations, data mining, machine learning.





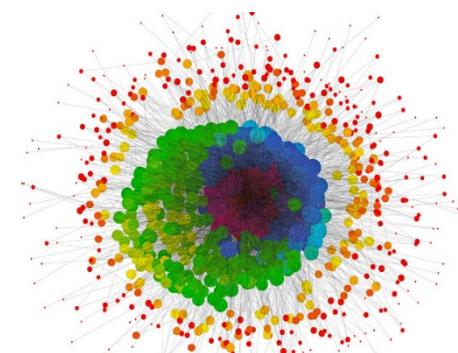
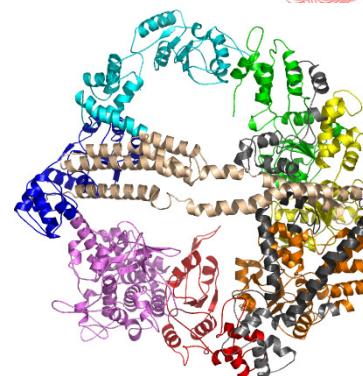
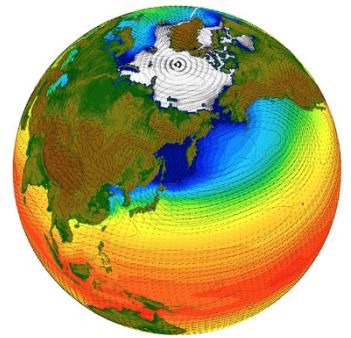
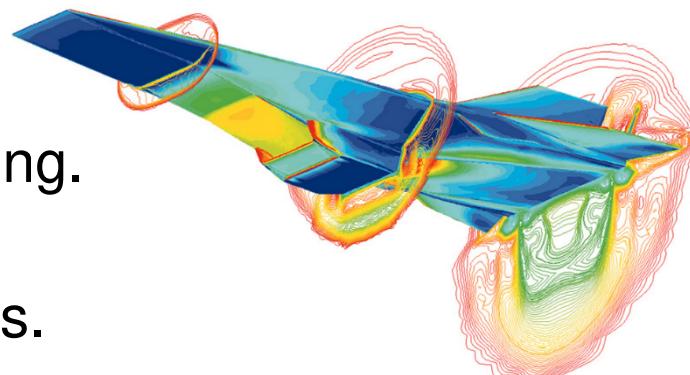
# Course objectives

- To understand the concepts and techniques of parallel computing, and take advantage of the capabilities of modern systems.
  - Parallel hardware models and interaction with parallel software.
  - Power and limitations of parallelism.
  - Efficient parallel algorithms for important problems.



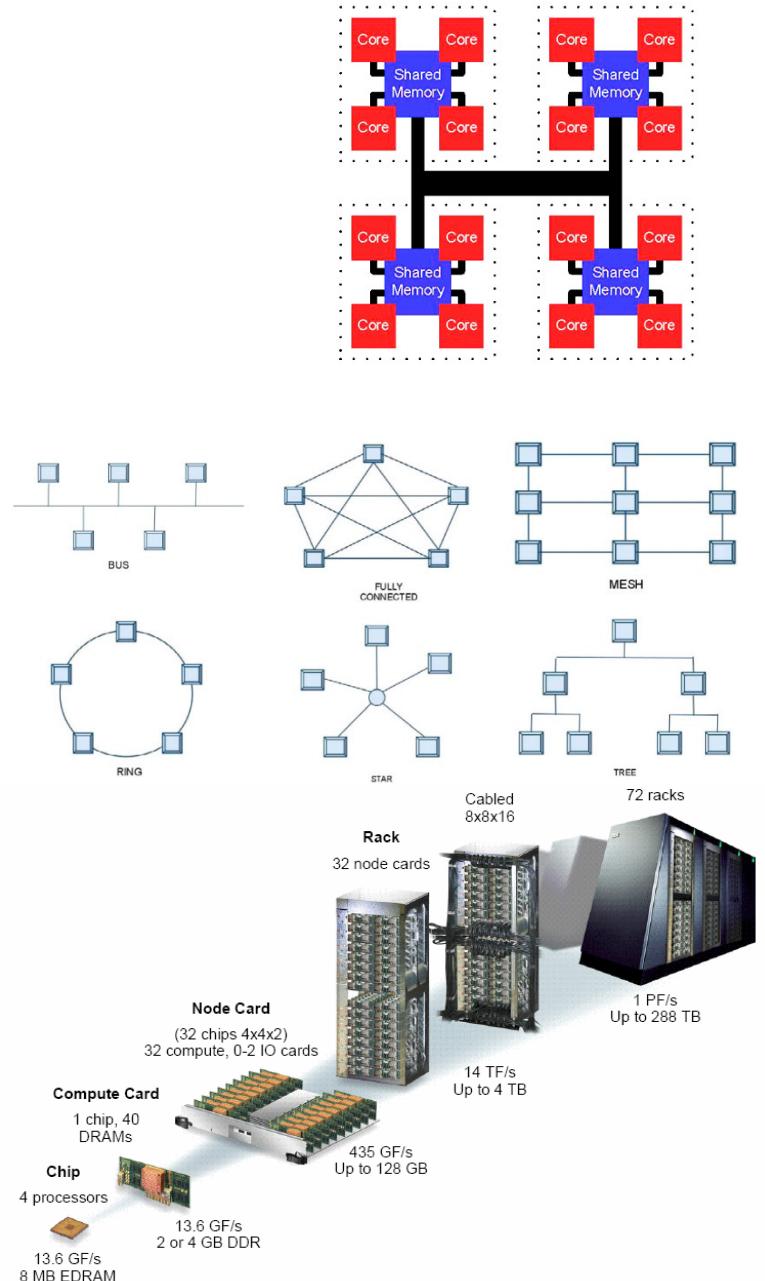
# Applications

- Fluid dynamics, weather prediction, climate modeling.
  - DNA, protein, drug structures and interactions.
  - Quantum / atomic simulations, cosmological simulations.
  - Cryptoanalysis.
  - Big data analytics.
  - Simulating financial and social behaviors.
  - Machine learning and AI.
  - Simulating the human brain.



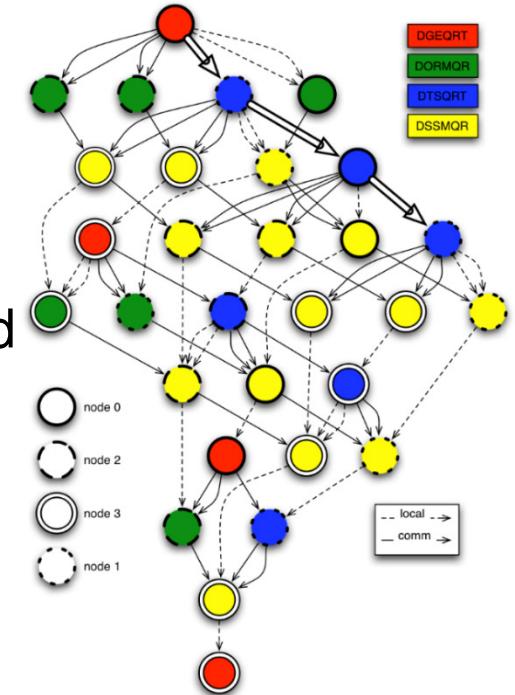
# Parallel hardware

- Efficient parallel computing requires synergy between parallel hardware and software.
- Parallel system consists of multiple independent processors communicating over an interconnect.
- Unlike sequential (von Neumann) architecture, many parallel hardware designs.
  - Different types of processors (multicores, manycores, FPGA, etc.).
  - Heterogeneous designs combine multiple architectures, e.g. multicores and GPUs.
  - Different interconnect designs.
  - Communicate through shared memory, or message passing over network.
- Parallelism exists at many layers.
  - Instruction, core, chip, node, rack, etc.



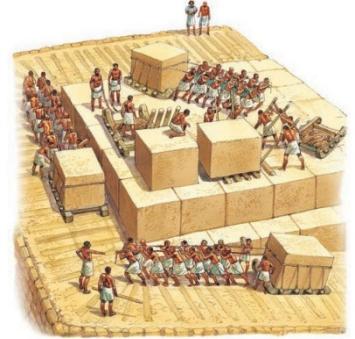
# Parallel software

- Break a large problem into subproblems (tasks) that can be solved (somewhat) independently.
- OS and scheduler allocate tasks to different processors.
  - Respect dependencies between tasks.
- Parallel software must be matched to the hardware.
  - Similar amounts of concurrency in software and hardware.
  - Hardware must adequately handle software communication pattern.
  - No single hardware model suffices.
  - Parallel software is often not portable.
- PRAM model tries to abstract parallel hardware.
  - Useful for understanding inherent parallelism.
  - Unrealistically discounts cost of communication.



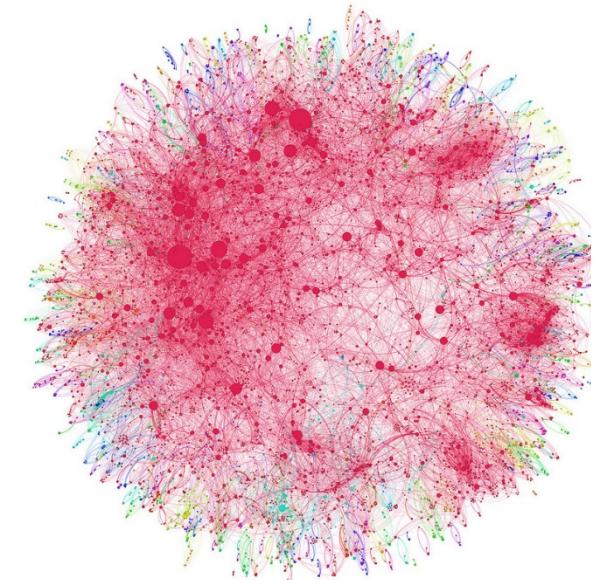
# Challenges

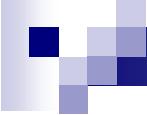
- Harnessing power of the masses.
  - Easier said than done...
- Communication
  - Processors compute faster than they can communicate.
  - Problem gets worse as number of processors increases.
  - Main bottleneck to parallel computing.
- Synchronization
  - Tasks may interfere with each other, so can't be done at same time.
- Scheduling
  - Track and enforce dependencies.
  - Find good allocation of tasks to processors.
    - Data locality, heterogeneous processors
  - Maximize utilization and performance.



# Challenges

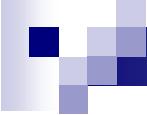
- Structured vs unstructured
  - Structured problems can be solved with custom hardware.
  - Unstructured problems more general, but less efficient.
- Inherent limitations
  - Some problems are not (or don't seem to be) parallelizable.
    - Ex Binary search, Dijkstra's shortest paths algorithm.
  - Other problems require clever algorithms to become parallel.
    - Ex Fibonacci series ( $a_n = a_{n-1} + a_{n-2}$ ).
- The human factor
  - Hard to keep track of concurrent events and dependencies.
  - Parallel algorithms are hard(er) to design and debug.





# Course outline

- Parallel architectures
  - Shared memory
  - Distributed memory
  - Manycore
- Parallel languages
  - OpenMP, MPI, CUDA, MapReduce
- Algorithm design techniques
  - Decomposition, load balancing, scheduling
- Parallel algorithms
  - Dense and sparse matrix algorithms, sorting, search, graph algorithms, PRAM algorithms, etc.

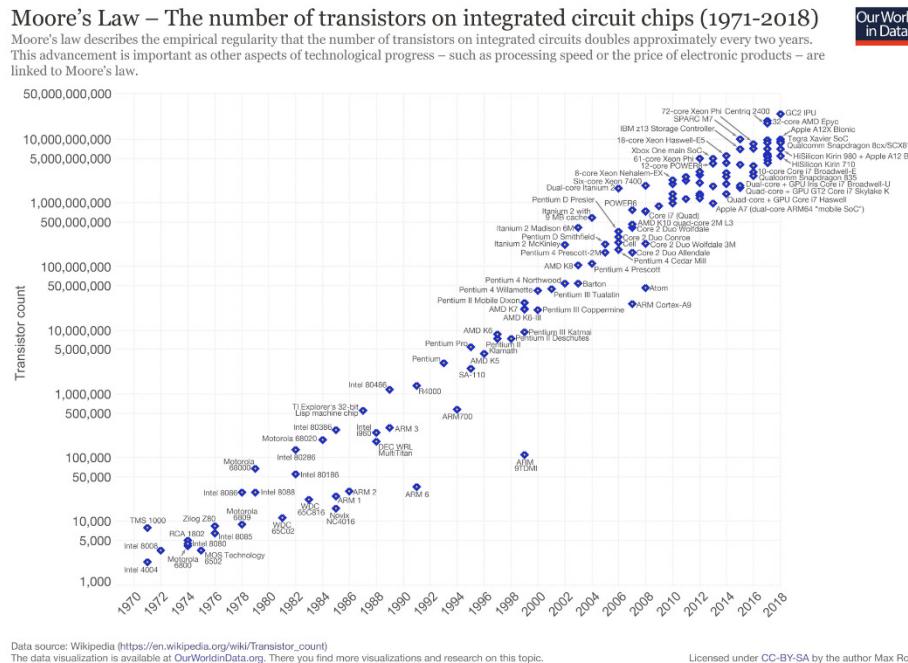


# A brief history

- Research and theory started in the early 60's.
  - Cray-1 reached 160 MFLOPS in 1976.
- Commercially successful supercomputers (Cray, Thinking Machines, etc.) started in 1980's.
  - Used expensive custom processors.
- In 1990's massively parallel processors (MPPs) and clusters became dominant.
  - MPPs use commercial (OTS) processors with custom interconnects.
  - Clusters use OTS processors and interconnects running Linux.
    - Cheap, easy to build and relatively powerful.
    - Most data centers today are clusters.
- Fastest supercomputer today is Fujitsu Fugaku MPP.
  - Runs at 442 PFLOPS, about 3M times faster than a workstation.
- Apart from supercomputers, progress in parallel computing stalled in 1990's until mid 2000's.

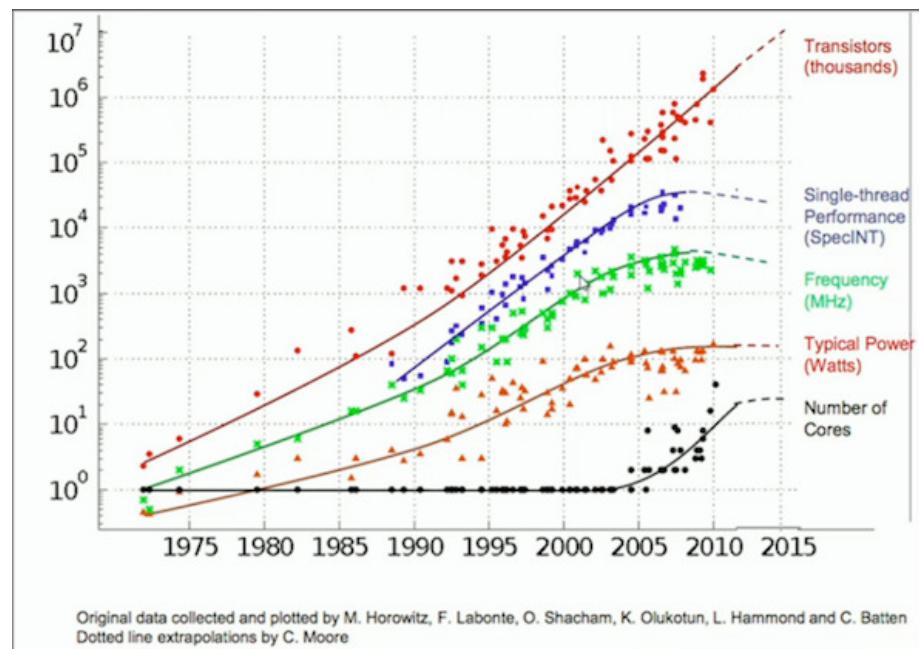
# Moore's Law and parallel computing

- In 1965, Gordon Moore, co-founder of Intel, predicted transistor count would double every 18 months.
  - Held true for the last 50 years!
- Until mid 2000's, this implied single processor performance doubled at same rate.
- This held back development of parallel computers, since in the time to develop one, single processor performance would improve dramatically.
- But since ca. 2005, parallel processing has become essential for taking advantage of Moore's Law.



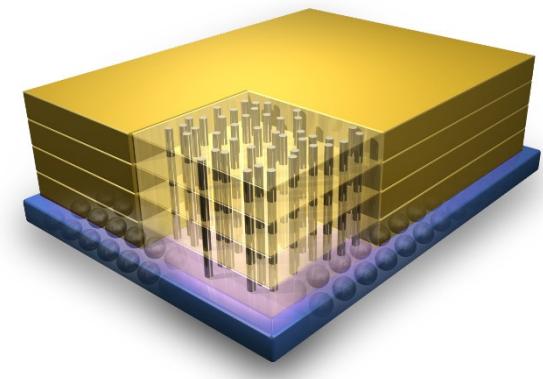
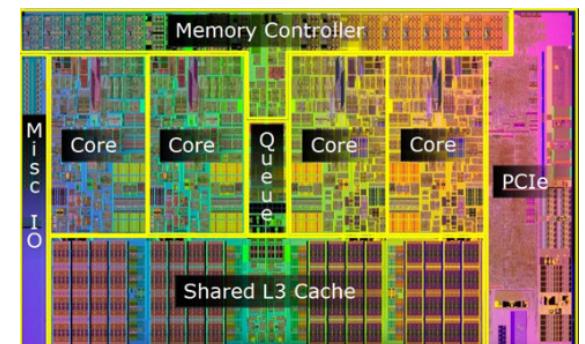
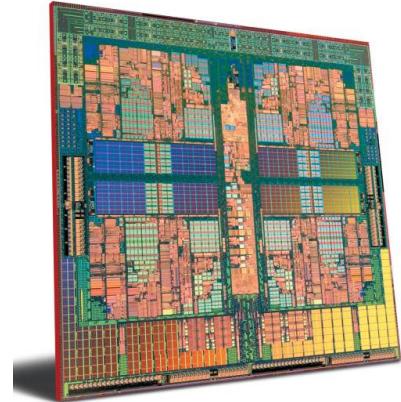
# Moore's Law and performance

- Transistor properties, e.g. size and clock speed, do not scale equally.
- Higher single processor clock speeds is increasingly difficult to achieve.
  - Heat
  - Power consumption
  - Current leakage



# Moore's Law revisited

- Multicore technology addresses (lack of) clock speed scaling.
  - Link multiple processing cores together on same chip.
  - More efficient to replace a single high speed processor with multiple slower processors.
  - Another approach is to stack chips in a 3D structure.
- Developing software for multicores has been harder than scaling hardware.
  - Software developers with parallel computing skills are in high demand.



# The state of the art

- Parallel computers today mainly based on four processor architectures.
  - Multicores
    - Small / moderate number ( $\leq 128$ ) of fast, general purpose cores.
    - Ex AMD EPYC, Intel Xeon, IBM Power.
  - Manycores
    - Large number (10K's) of simple cores.
    - Ex Nvidia Ampere GPU, Intel Xeon Phi, Sunway SW26010Pro.
  - FPGA (field programmable gate arrays)
    - Reconfigurable hardware customized for specific problems.
  - ASIC (application specific integrated circuits)
    - Specially built hardware for specific problems.
    - Ex Google TPU, Graphcore IPU, IBM TrueNorth.
- In addition to processing speed, energy efficiency also increasing important.
  - Biggest datacenters consume over 100 MW of power,  $\sim 50K$  homes.
  - Biggest supercomputers consume  $\sim 20$ MW of power.
  - Best supercomputers achieve 50 GFLOPS / W.

# Top 500 list

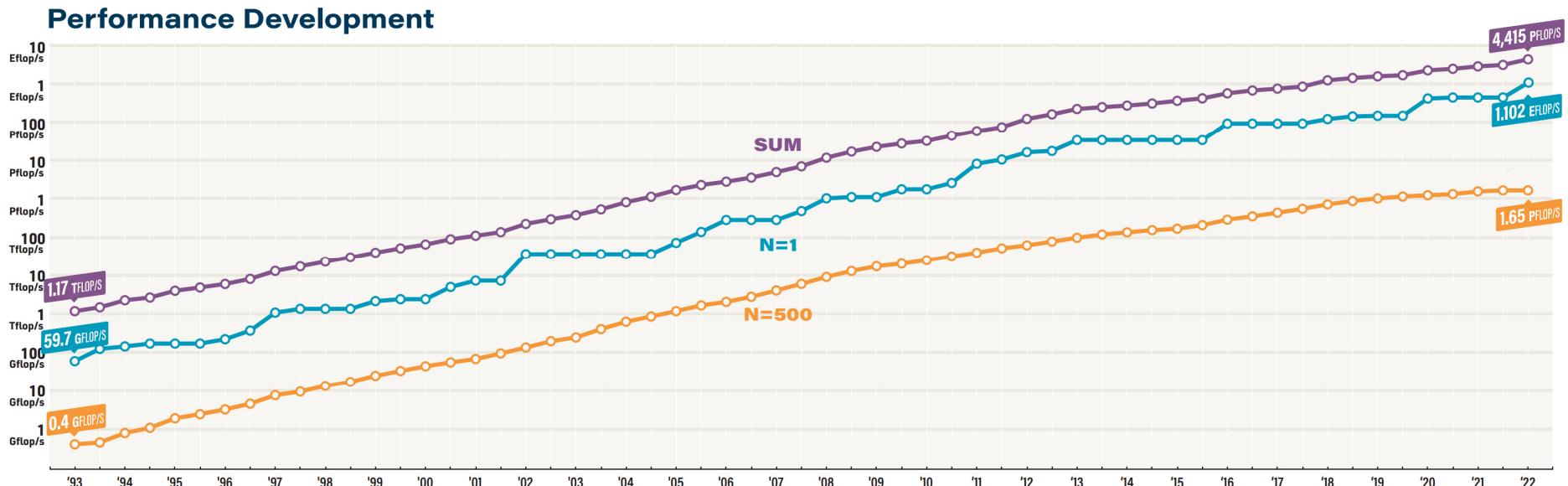
- Biannual ranking of fastest 500 supercomputers in the world.
  - Speed measured in floating point operations per second.
  - Uses high-performance LINPACK to solve a dense linear system  $Ax = b$ .
    - Compute intensive, but doesn't stress memory system.
    - May not represent performance on real-world problems.

JUNE 2022	SYSTEM	SPECS	SITE	COUNTRY	CORES	R <sub>MAX</sub> PFLOP/S	POWER MW
1	Frontier	HPE Cray EX235a, AMD Opt 3rd Gen EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-10	DOE/SC/ORNL	USA	8,730,112	1,102.0	21.3
2	Fugaku	Fujitsu A64FX (48C, 2.2GHz), Tofu Interconnect D	RIKEN R-CCS	Japan	7,630,848	442.0	29.9
3	LUMI	HPE Cray EX235a, AMD Opt 3rd Gen EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-10	EuroHPC/CSC	Finland	1,268,736	151.9	2.94
4	Summit	IBM POWER9 (22C, 3.07GHz), NVIDIA Volta GV100 (80C), Dual-Rail Mellanox EDR Infiniband	DOE/SC/ORNL	USA	2,414,592	148.6	10.1
5	Sierra	IBM POWER9 (22C, 3.1GHz), NVIDIA Tesla V100 (80C), Dual-Rail Mellanox EDR Infiniband	DOE/NNSA/LLNL	USA	1,572,480	94.6	7.44

Mega	Giga	Tera	Peta	Exa
$10^6$	$10^9$	$10^{12}$	$10^{15}$	$10^{18}$

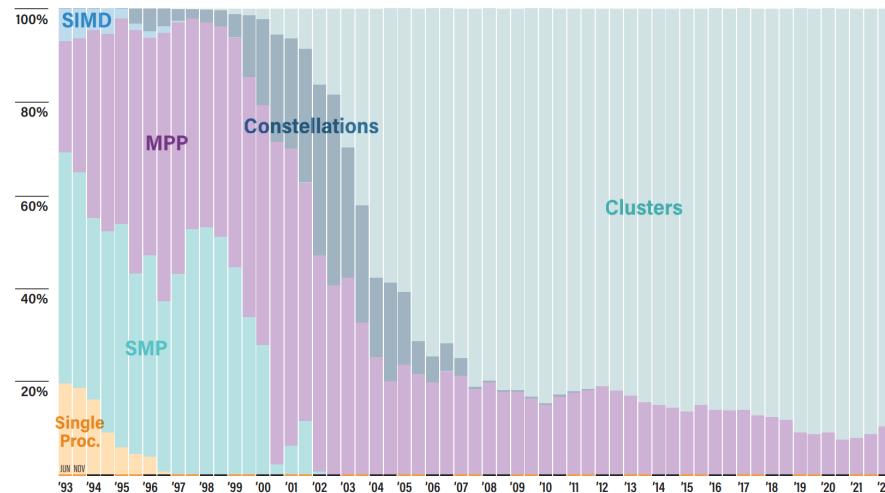
- For comparison, Intel multicore achieves ~50 GFLOPS / core, and GPU achieves ~ 10 TFLOPS / board.

# Top 500 – Trends

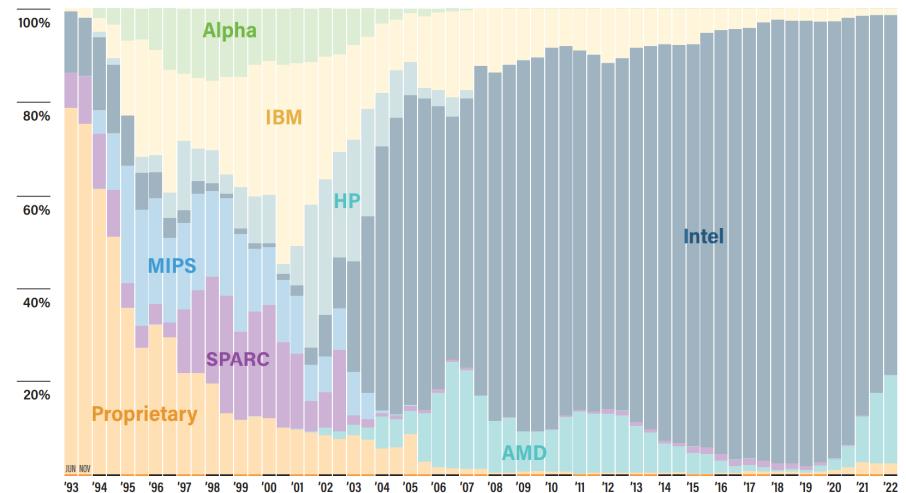


# Top 500 – Architecture

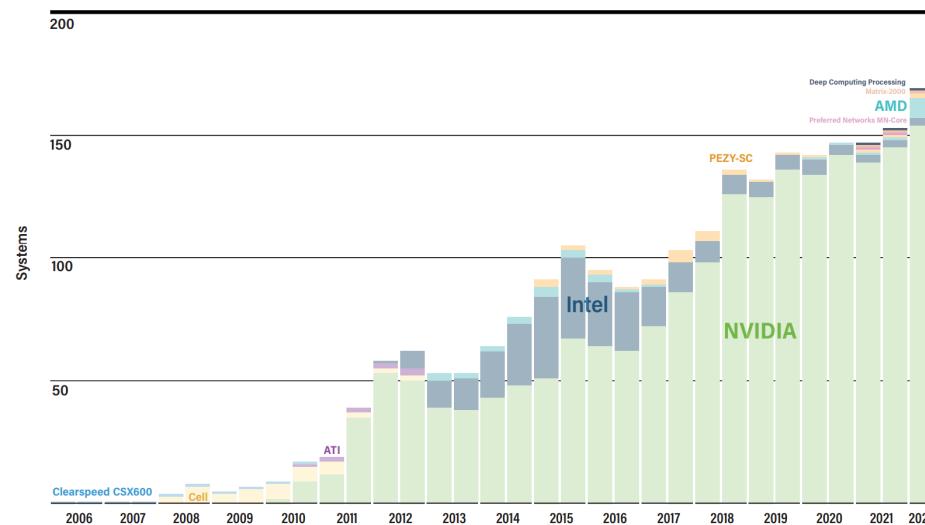
**Architectures**



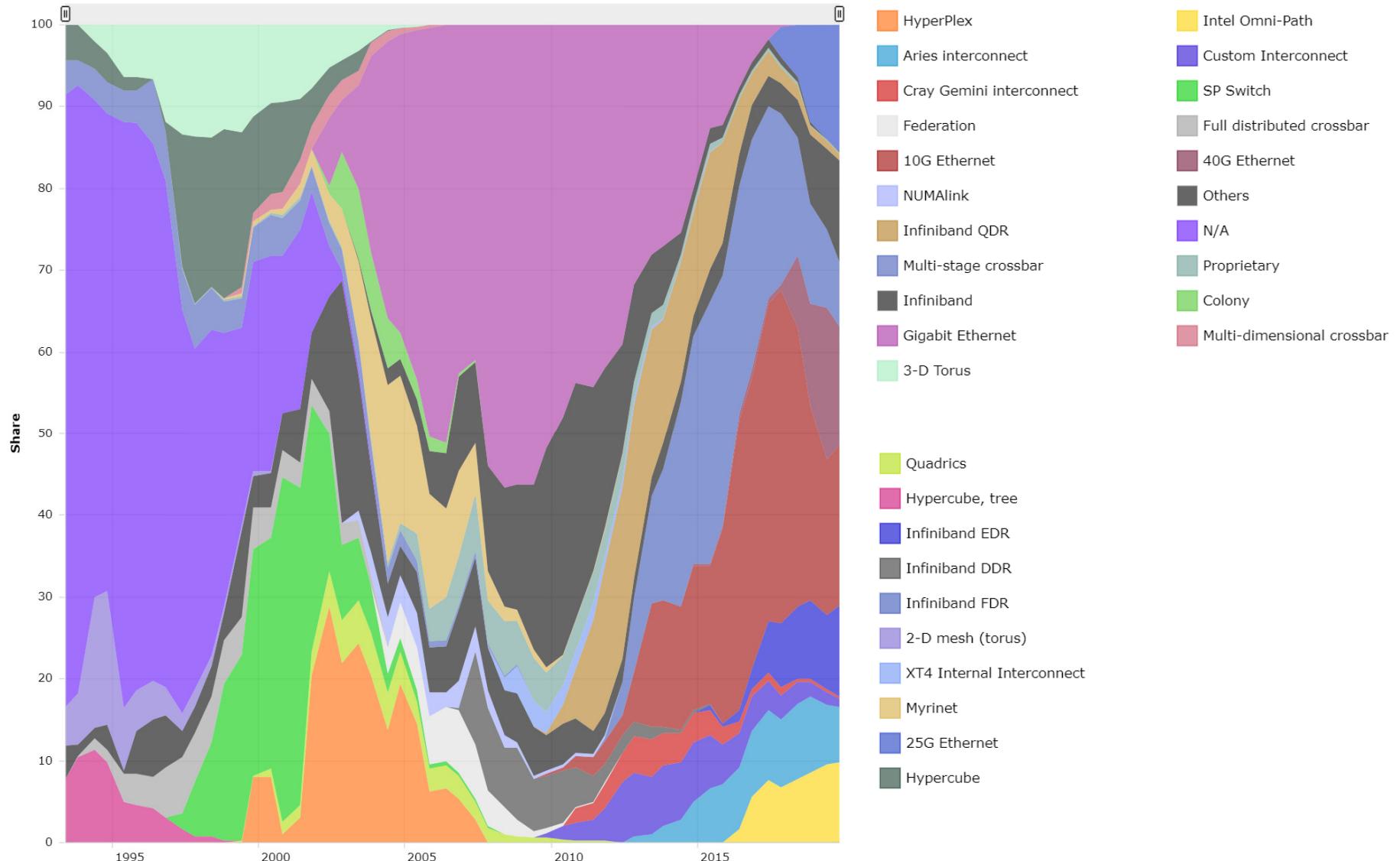
**Chip Technology**

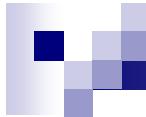


**Accelerators/Co-processors**

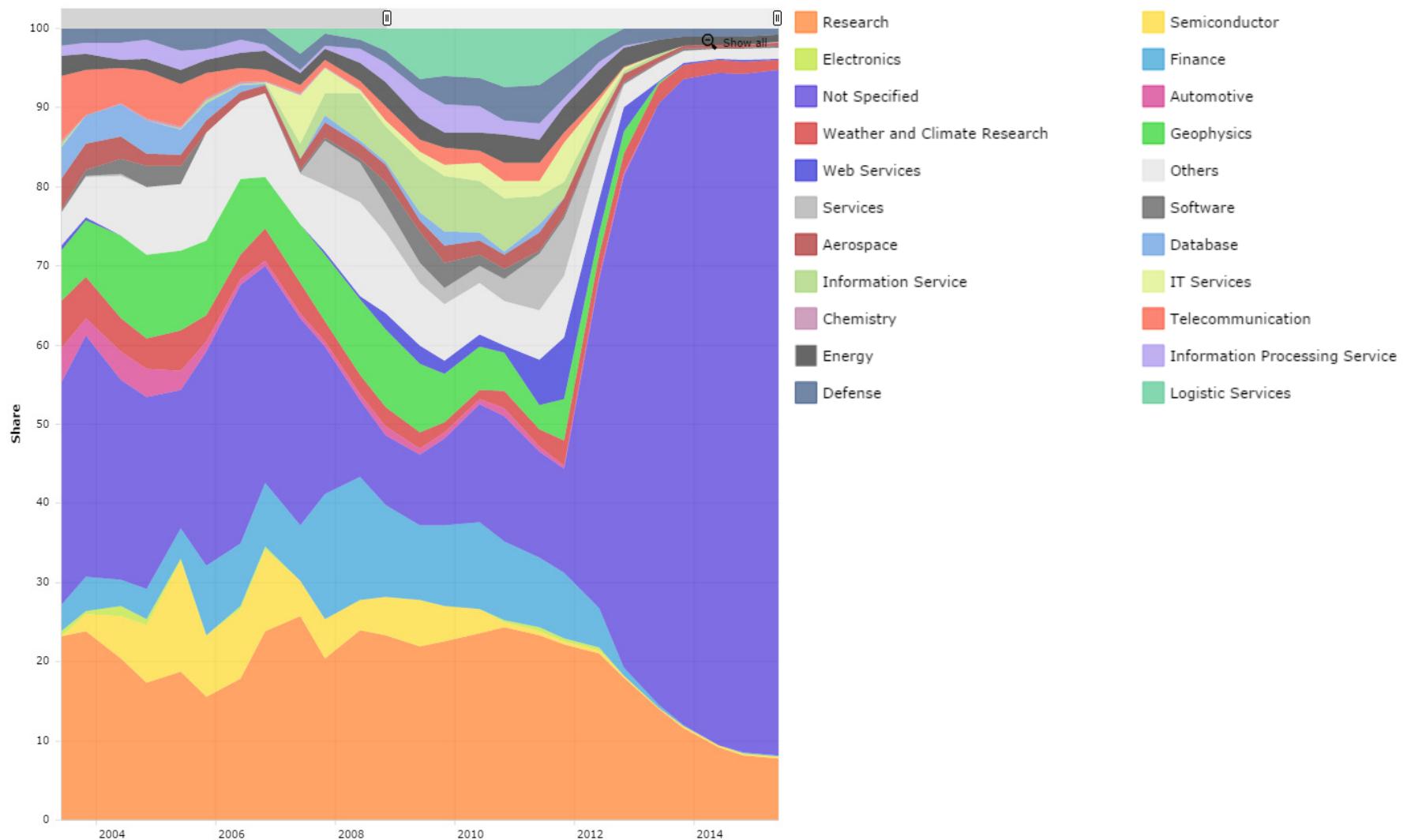


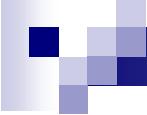
# Top 500 – Interconnect





# Top 500 – Applications





# Other performance measures

- LINPACK does compute-intensive operations on structured dense matrices.
  - Uniform control flow, predictable and coalesced memory accesses.
  - Ideal for physical simulations.
- Data-intensive applications today have instruction divergence, branching and random memory accesses.
- New benchmarks give more complete performance picture
  - HPCG performs sparse matrix operations.
  - Graph500 performs breadth-first search.
- A computer's performance can differ dramatically depending on benchmark.

# HPCG

New HPCG results announced at ISC 2022

Rank	Site	Computer	Cores	HPL Rmax (Pflop/s)	TOP500 Rank	HPCG (Pflop/s)	Fraction of Peak
1	RIKEN Center for Computational Science <b>Japan</b>	<b>Supercomputer Fugaku</b> — A64FX 48C 2.2GHz, Tofu interconnect D	7,630,848	442.01	2	16.00	3.0%
2	DOE/SC/Oak Ridge National Laboratory <b>United States</b>	<b>Summit</b> — IBM POWER9 22C 3.07GHz, Dual-rail Mellanox EDR Infiniband, NVIDIA Volta GV100	2,414,592	148.60	4	2.926	1.5%
3	EuroHPC/CSC <b>Finland</b>	<b>LUMI</b> — AMD Optimized 3rd Generation EPYC 64C 2GHz, Slingshot-11, AMD Instinct MI250X	1,110,144	151.90	3	1.936	0.9%
4	DOE/SC/LBNL/NERSC <b>United States</b>	<b>Perlmutter</b> — AMD EPYC 7763 64C 2.45GHz, Slingshot-10, NVIDIA A100 SXM4 40 GB	761,856	70.87	7	1.905	2.0%
5	DOE/NNSA/LLNL <b>United States</b>	<b>Sierra</b> — IBM POWER9 22C 3.1GHz, Dual-rail Mellanox EDR Infiniband, NVIDIA Volta GV100	1,572,480	94.64	5	1.796	1.4%
6	NVIDIA Corporation <b>United States</b>	<b>Selene</b> — AMD EPYC 7742 64C 2.25GHz, Mellanox HDR Infiniband, NVIDIA A100	555,520	63.46	8	1.623	2.0%
7	Forschungszentrum Juelich (FZJ) <b>Germany</b>	<b>JUWELS Booster Module</b> — AMD EPYC 7402 24C 2.8GHz, Mellanox HDR InfiniBand/ParTec ParaStation ClusterSuite, NVIDIA A100	449,280	44.12	11	1.275	1.8%
8	Saudi Aramco <b>Saudi Arabia</b>	<b>Dammam-7</b> — Xeon Gold 6248 20C 2.5GHz, InfiniBand HDR 100, NVIDIA Tesla V100 SXM2	672,520	22.40	18	0.881	1.6%
9	Eni S.p.A. <b>Italy</b>	<b>HPC5</b> — Xeon Gold 6252 24C 2.1GHz, Mellanox HDR Infiniband, NVIDIA Tesla V100	669,760	35.45	12	0.860	1.7%
10	Information Technology Center, The University of Tokyo <b>Japan</b>	<b>Wisteria/BDEC-01 (Odyssey)</b> — A64FX 48C 2.2GHz, Tofu interconnect D	368,640	22.12	20	0.818	3.2%

# Graph500

RANK	PREVIOUS RANK	MACHINE	VENDOR	TYPE	NETWORK	INSTALLATION SITE	LOCATION	COUNTRY	YEAR	APPLICATION	USAGE	NUMBER OF NODES	NUMBER OF CORES	MEMORY	IMPLEMENTATION	SCALE	GTEPS
1	1	Supercomputer Fugaku	Fujitsu	Fujitsu A64FX	Tofu Interconnect D	RIKEN Center for Computational Science (R-CCS)	Kobe Hyogo	Japan	2020	Various scientific and industrial fields	Academic and industry	158976	7630848	5087232	Custom	41	102955
2	2	Sunway TaihuLight	NRCPC	Sunway MPP	Sunway	National Supercomputing Center in Wuxi	Wuxi	China	2015	research	research	40768	10599680	1304580 gigabytes	Custom	40	23755.7
3	3	Wisteria/BDEC-01 (Odyssey)	Fujitsu	PRIMEHPC FX1000	Tofu interconnect D	Information Technology Center The University of Tokyo	Kashiwa Chiba	Japan	2021	University	Research	7680	368640	245760	Custom	37	16118
4	4	TOKI-SORA	Fujitsu	PRIMEHPC FX1000	Tofu interconnect D	Japan Aerospace eXploration Agency (JAXA)	Tokyo	Japan	2020	Research	CFD	5760	276480	184320	Custom	36	10813
5	5	LUMI-C	HPE	HPE Cray EX	HPE Slingshot-10	EuroHPC/CSC	Kajaani	Finland	2021	Research	Various	1492	190976	381952	custom	38	8467.71
6	6	OLCF Summit (CPU-Only)	IBM	IBM POWER9		Oak Ridge National Laboratory	Oak Ridge TN	United States	2018	Government	Scientific Research	2048	86016	1048576	Custom	40	7665.7
7	7	SuperMUC-NG	Lenovo	ThinkSystem SD530 Xeon Platinum 8174 24C 3.1GHz Intel Omni-Path		Leibniz Rechenzentrum	Garching	Germany	2018	Academic	Research	4096	196608	393216	custom-amd-heavy-diropt	39	6279.47
8	8	Lise	Atos	Bull Intel Cluster Intel Xeon Platinum 9242 48C 2.3GHz Intel Omni-Path	Intel Omni-Path	Zuse Institute Berlin (ZIB)	Berlin	Germany	2019	Research	Academic	1270	121920	502272	custom the same code used on Fugaku	38	5423.94
9	new	DepGraph Supernode	HUST &amp; Nvidia	DepGraph (+GPU Tesla A100)	Custom	National Engineering Research Center for Big Data Technology and System	Wuhan	China	2022	Research	Academic	1	128	512	Custom (Implementation of Yu Zhang)	33	4623.379
10	9	NERSC Cori - 1024 haswell partition	Cray	XC40	Aries	NERSC/LBNL	DOE/SC/LBNL/NERSC	United States	2017	Government	basic science and simulation	1024	32768	133376	custom-amd-heavy-diropt	37	2562.16



**HPCL**

State Key Laboratory of  
High Performance Computing

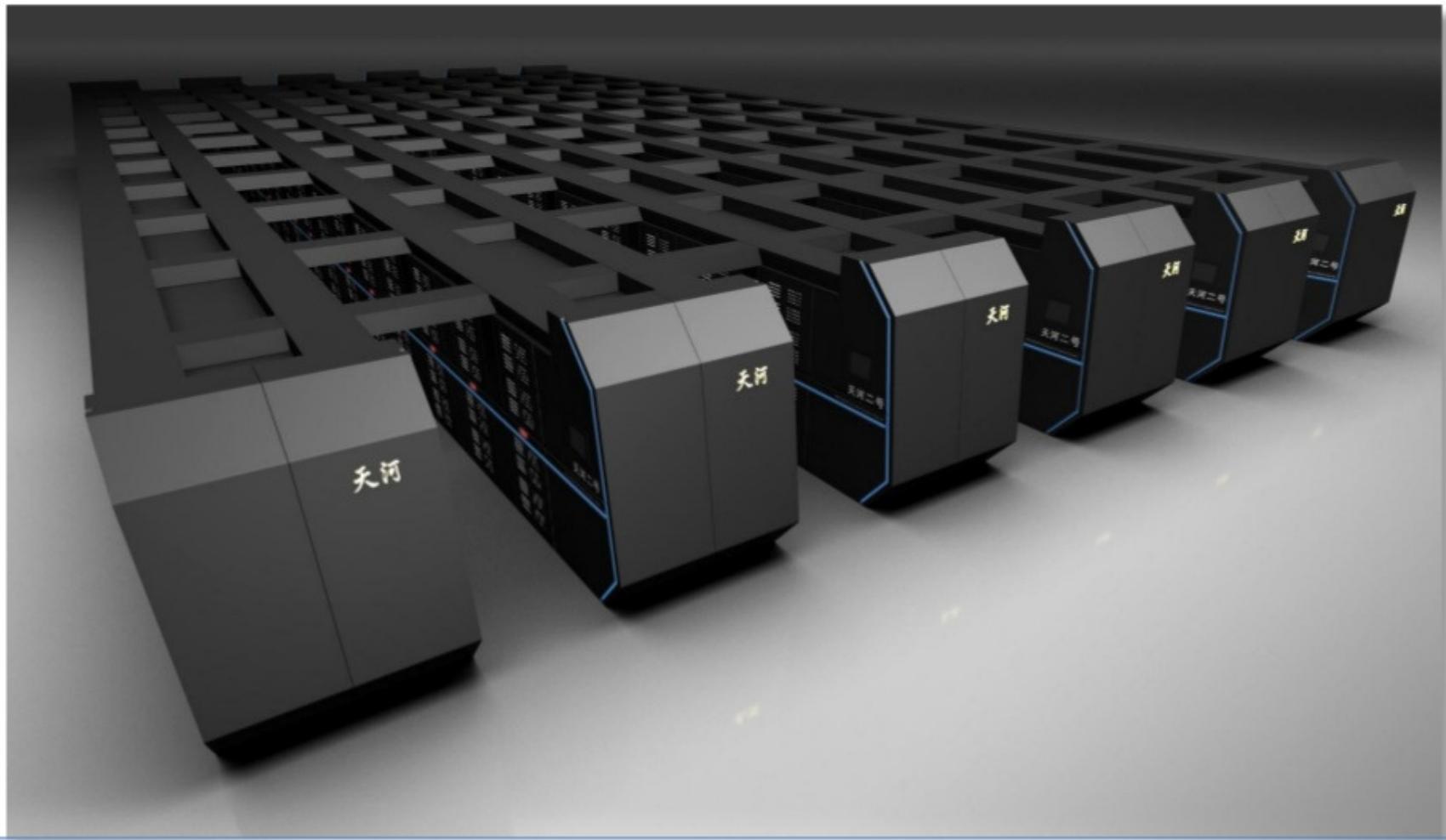
# Overview of Tianhe-2 (MilkyWay-2) Supercomputer

Yutong Lu

School of Computer Science, National University of Defense Technology;  
State Key Laboratory of High Performance Computing, China

[ytlu@nudt.edu.cn](mailto:ytlu@nudt.edu.cn)

# Motivation



**Tianhe-2 (Milkyway-2) Supercomputer**

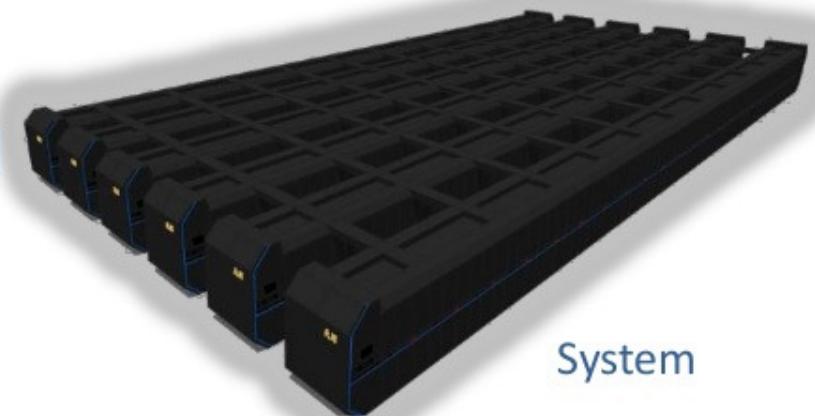
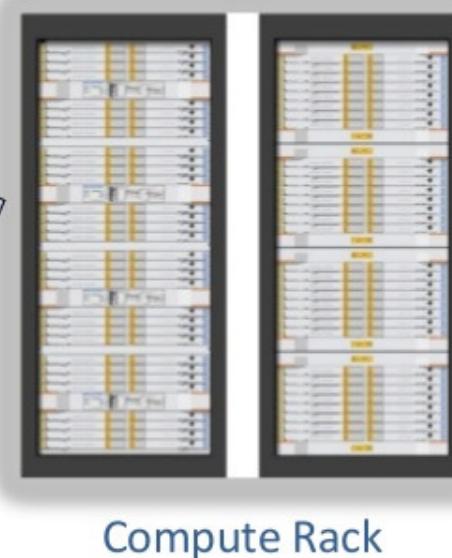
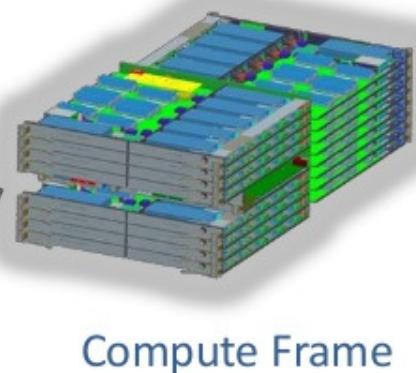
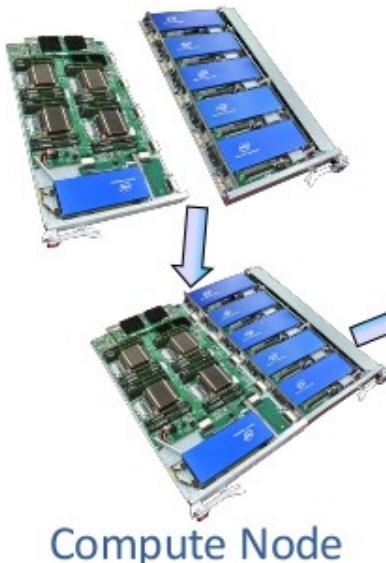
## ■ Hybrid Architecture

### ◆ Xeon CPU & Xeon Phi

Items	Configuration
Processors	32000 Intel Xeon CPUs + 48000 Xeon Phis + 4096 FT CPUs Peak performance is 54.9PFlops, HPL
Interconnect	Proprietary high-speed interconnection network TH Express-2
Memory	1.4PB in total
Storage	Global shared parallel storage system, 12.4PB
Cabinets	125+13+24=162 compute/communication/storage Cabinets
Power	17.8 MW (1902MFlops/W)
Cooling	Closed Air cooling system

# From Chips to Entire System

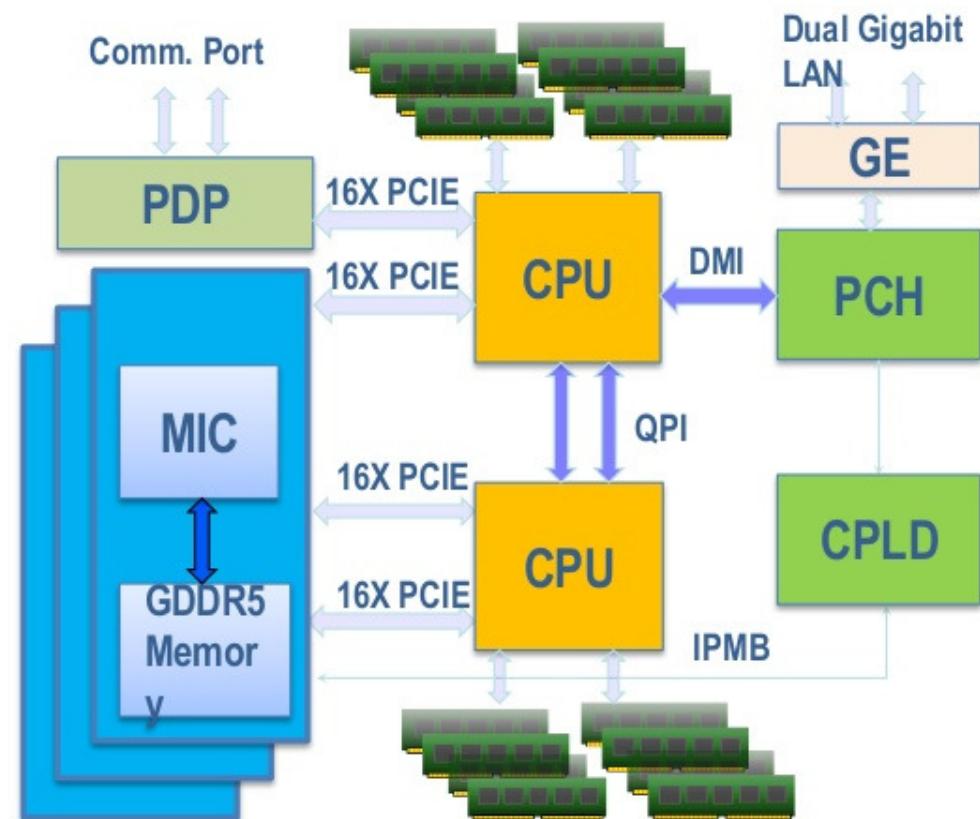
- ◆ 16000 compute nodes in total
- ◆ Frame: 32 compute Nodes
- ◆ Rack: 4 Compute Frames
- ◆ Whole System: 125 Racks



# Compute Node

## ■ Neo-Heterogeneous Compute Node

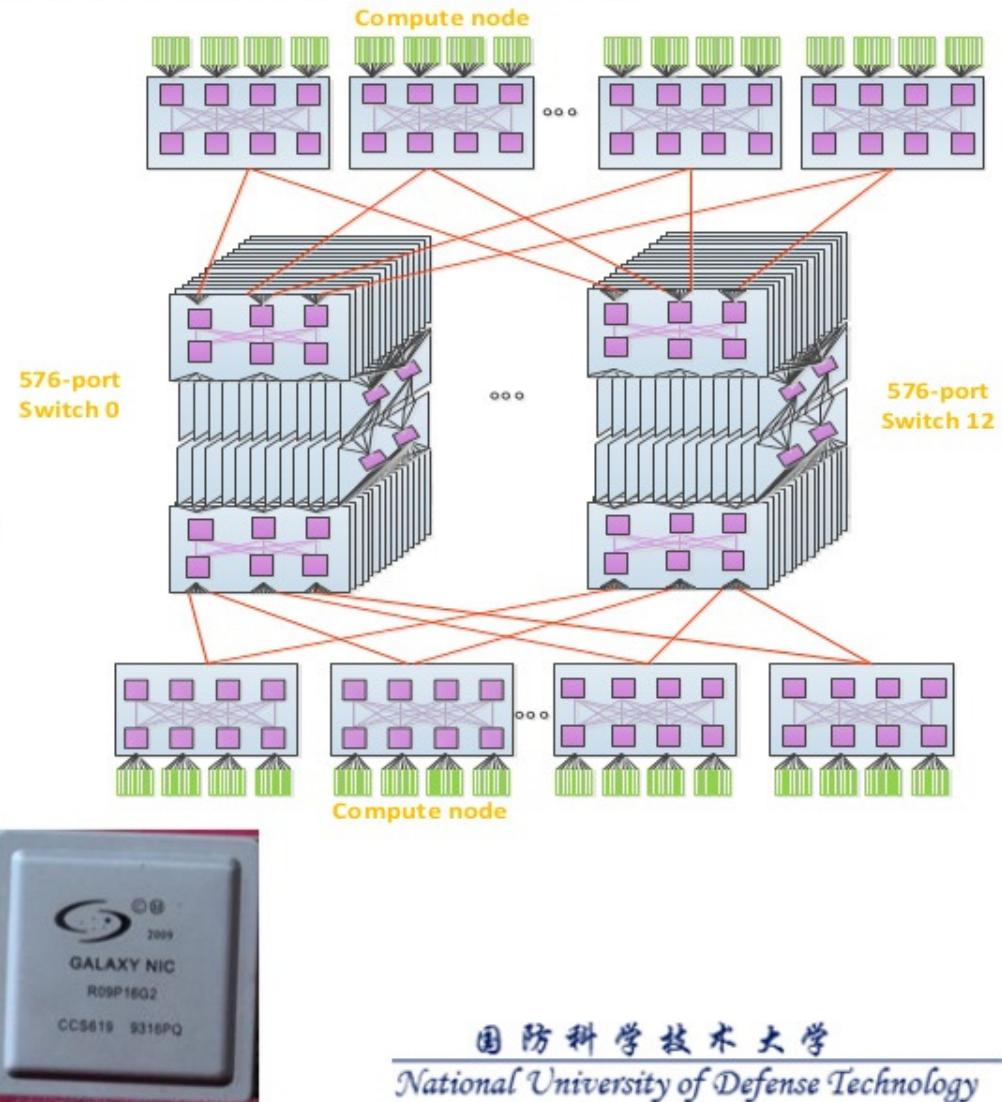
- ◆ Similar ISA, different ALU
- ◆ 2 Intel Ivy Bridge CPU + 3 Intel Xeon Phi
- ◆ 16 Registered ECC DDR3 DIMMs, 64GB
- ◆ 3 PCI-E 3.0 with 16 lanes
- ◆ PDP Comm. Port
- ◆ Dual Gigabit LAN
- ◆ Peak Perf. : 3.432Tflops



# Interconnection network

## ■ TH Express-2 interconnection network

- ◆ Fat-tree topology using 13 576-port top level switches
- ◆ Opto-electronic hybrid transport tech.
- ◆ Proprietary network protocol
- ◆ NRC +NIC





# HPC Software stack

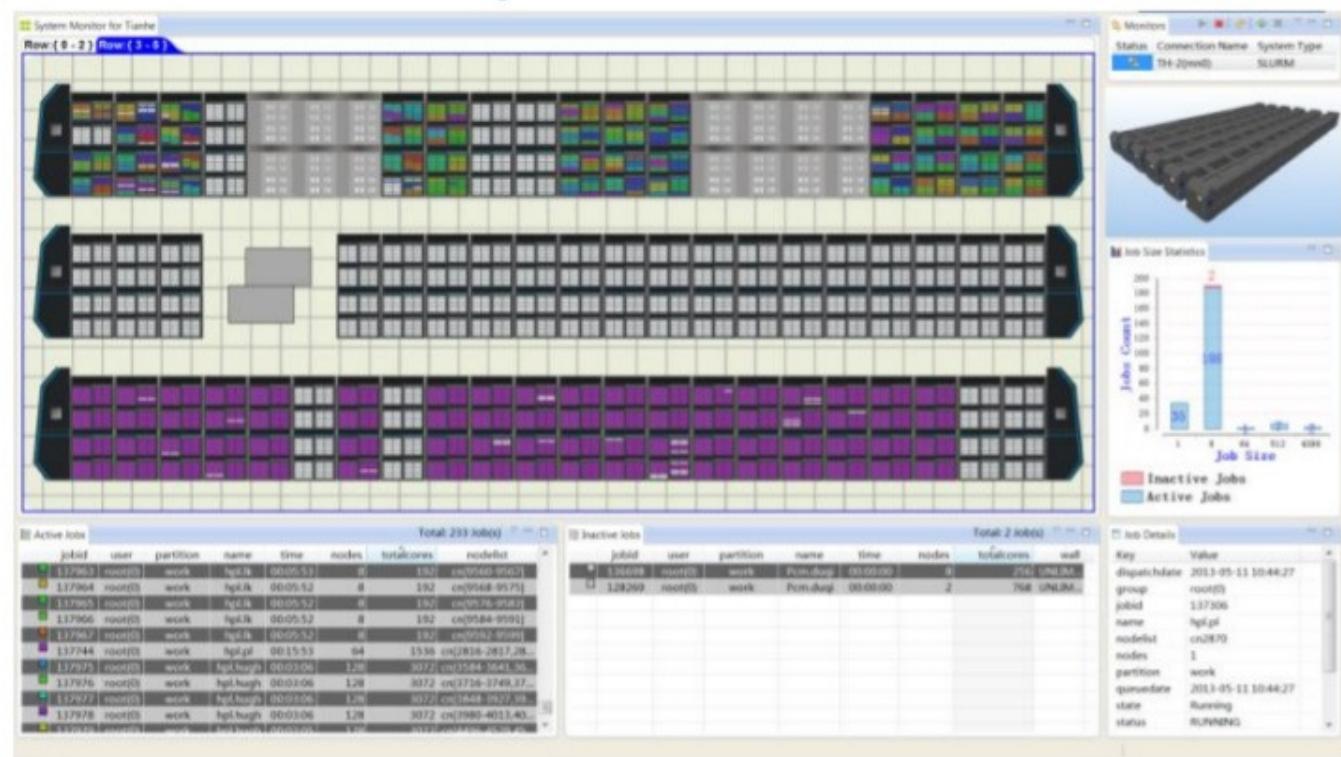


## ■ Operating System

- ◆ Kylin Linux

## ■ Resource manage system

- ◆ Power-aware resource allocation
- ◆ Multiple custom schedule policies

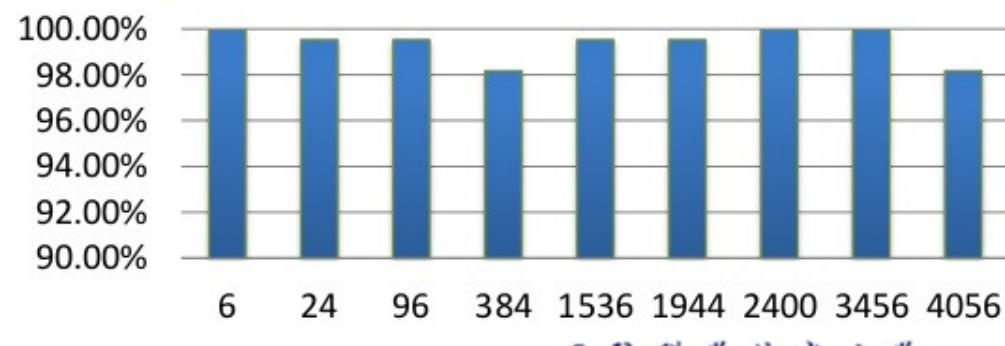
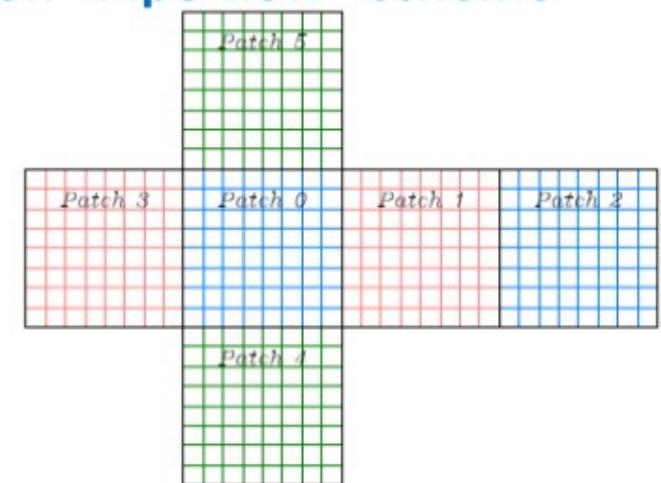
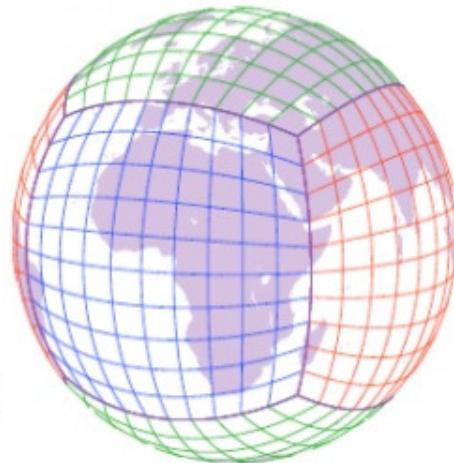


# Application

- Application of a global shallow water model: algorithms
  - ◆ Hierarchical data partition & communication on cubed-sphere
  - ◆ Balanced partition between CPU/MIC inside each node
  - ◆ Communication hiding algorithm based on “Pipe-flow” scheme

- Nearly ideal weak scaling on the Tianhe-2

- ◆ Using up to 4,056 nodes (97,344 CPU cores + 693,576 MIC cores)
  - ◆ # of unknowns for the largest run: 200 billion



# Course texts

- Course materials partly taken from the following texts.
  - But all topics covered by lecture slides.
- *Introduction to Parallel Computing*. Grama, Karypis, Kumar, Gupta. Pearson, 2003.
- *An Introduction to Parallel Programming*. Peter Pacheco. Morgan Kaufmann 2011.
- *Programming Massively Parallel Processors*. Kirk, Hwu. Morgan Kaufmann 2016.
- *CUDA by Example*. Sanders, Kandrot. Addison-Wesley 2010.

