

# Center of Excellence in Supercomputing for AI & Big Data

## Supercomputing Architectures and Development Methodologies for AI and Big Data Applications

by: Tassadaq Hussain

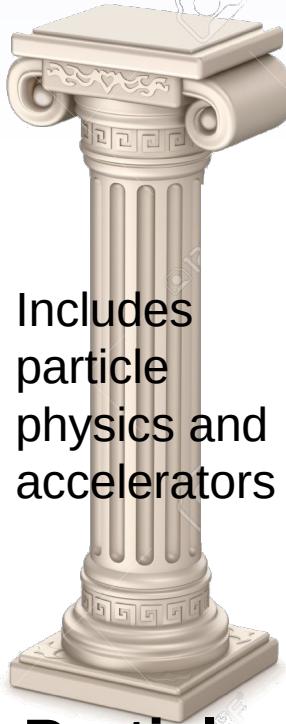
**Professor Department of Electrical Engineering  
Namal University Mianwali**

**Collaborations:**  
Barcelona Supercomputing Center, Spain  
European Network on High Performance and Embedded Architecture and Compilation  
Pakistan Supercomputing Center

- **Importance of HPC**
- Types of Processing System
- HPC System Architecture
- Supercomputing Classes and Infrastructure
- Cluster Software Stacks
- Programming Models
- Demonstration (Supercomputing System)

# Pillars of Science

Fermi National  
Accelerator Laboratory



Includes  
particle  
physics and  
accelerators

**Particle  
Physics**



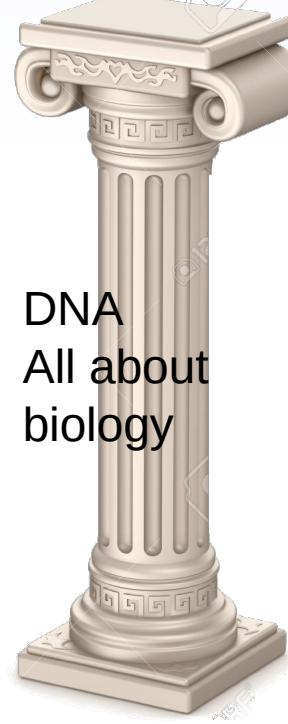
Includes all of  
cosmology,  
astrophysics

**Cosmology**



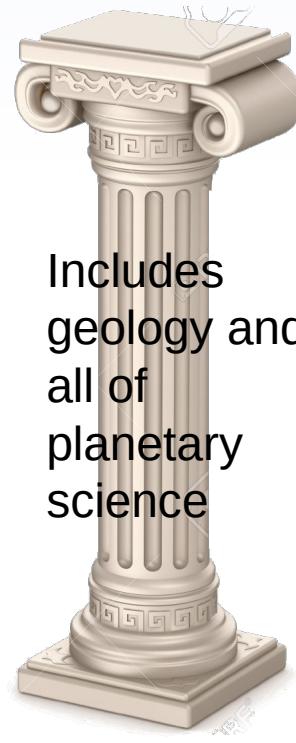
Data,  
Connectivity,  
AI &  
Processing

**Computing**



DNA  
All about  
biology

**Biology**



Includes  
geology and  
all of  
planetary  
science

**Space**

**QUARKS**

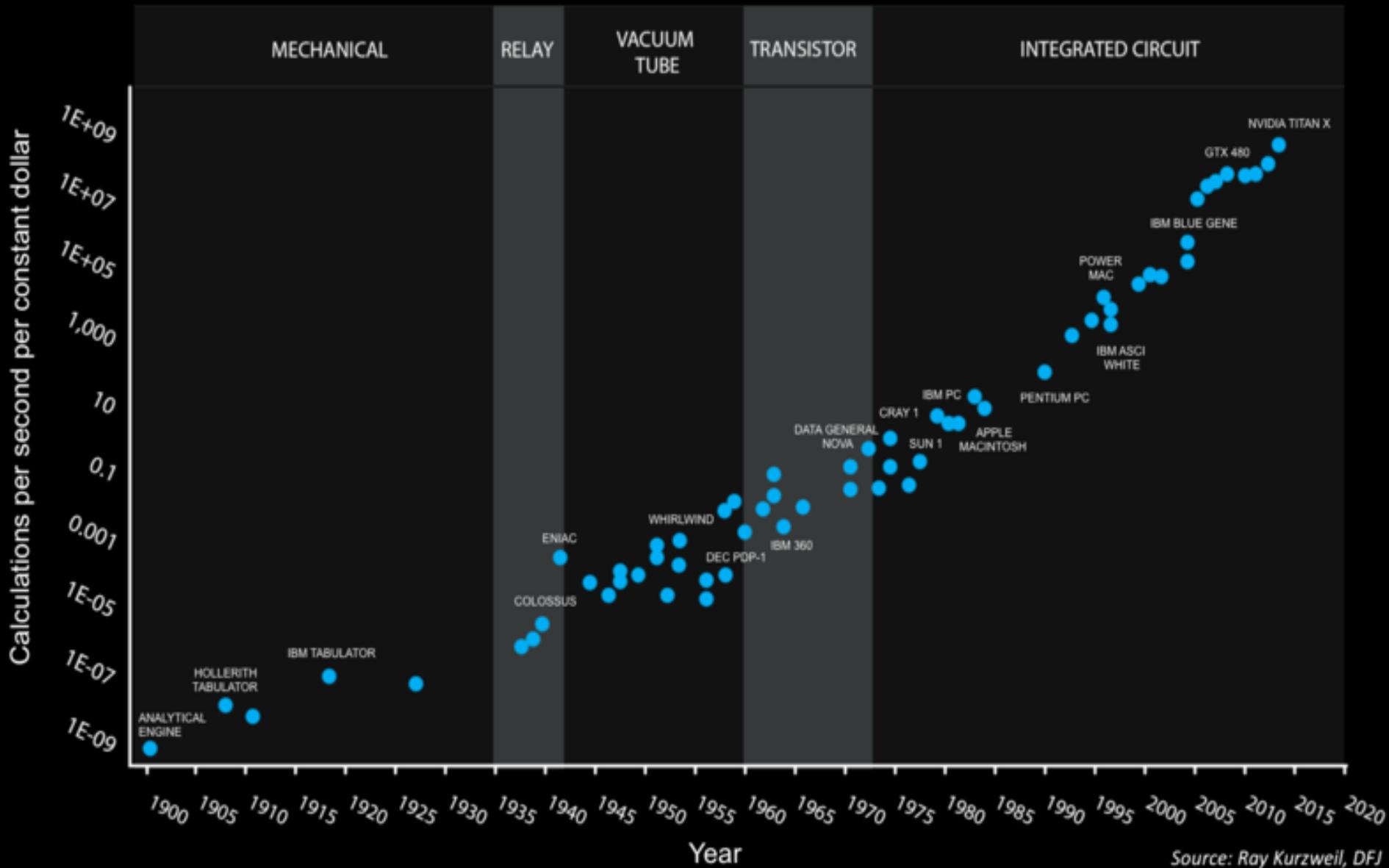
**BIG BANG**

**MicroElectronics**

**DNA**

**SPACE**

# 120 Years of Moore's Law



Source: Ray Kurzweil, DFI



Life Science



Earth Science



Social Science

Science

175 ZByte @2025

80%  
Data-Sciences

Data

100 ExaFLOPS  
@2020

87.04 B\$  
234.6 B\$ @2025

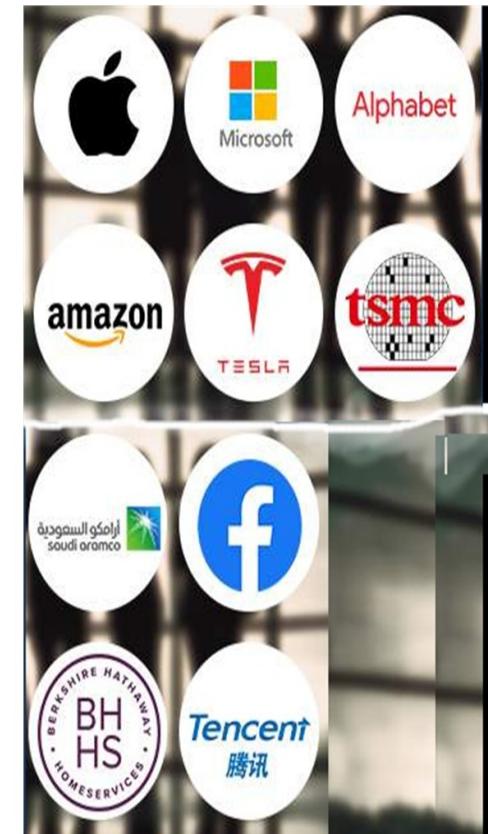
AI

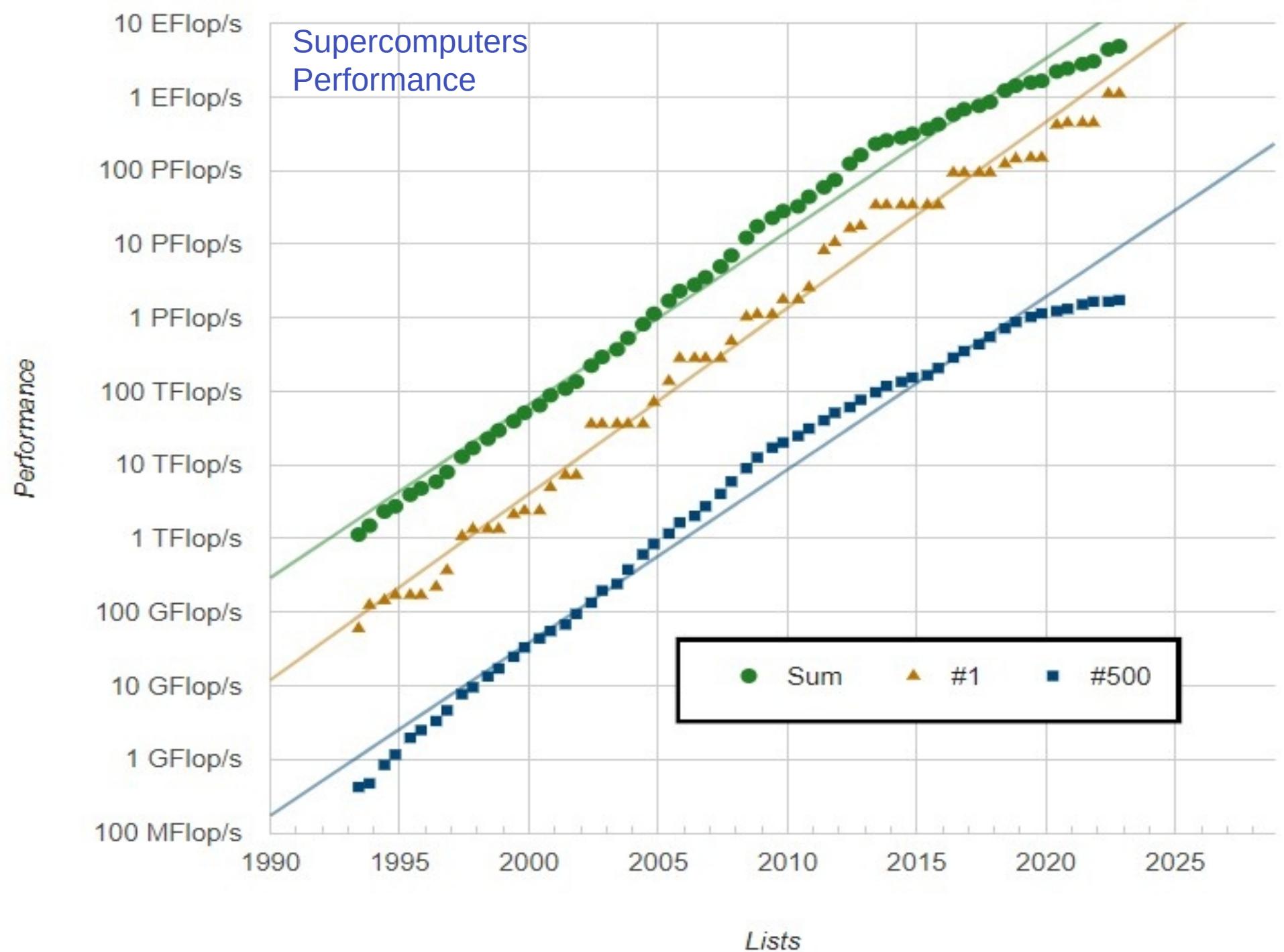
Top500 List  
8 PetaFLOPS  
@2022

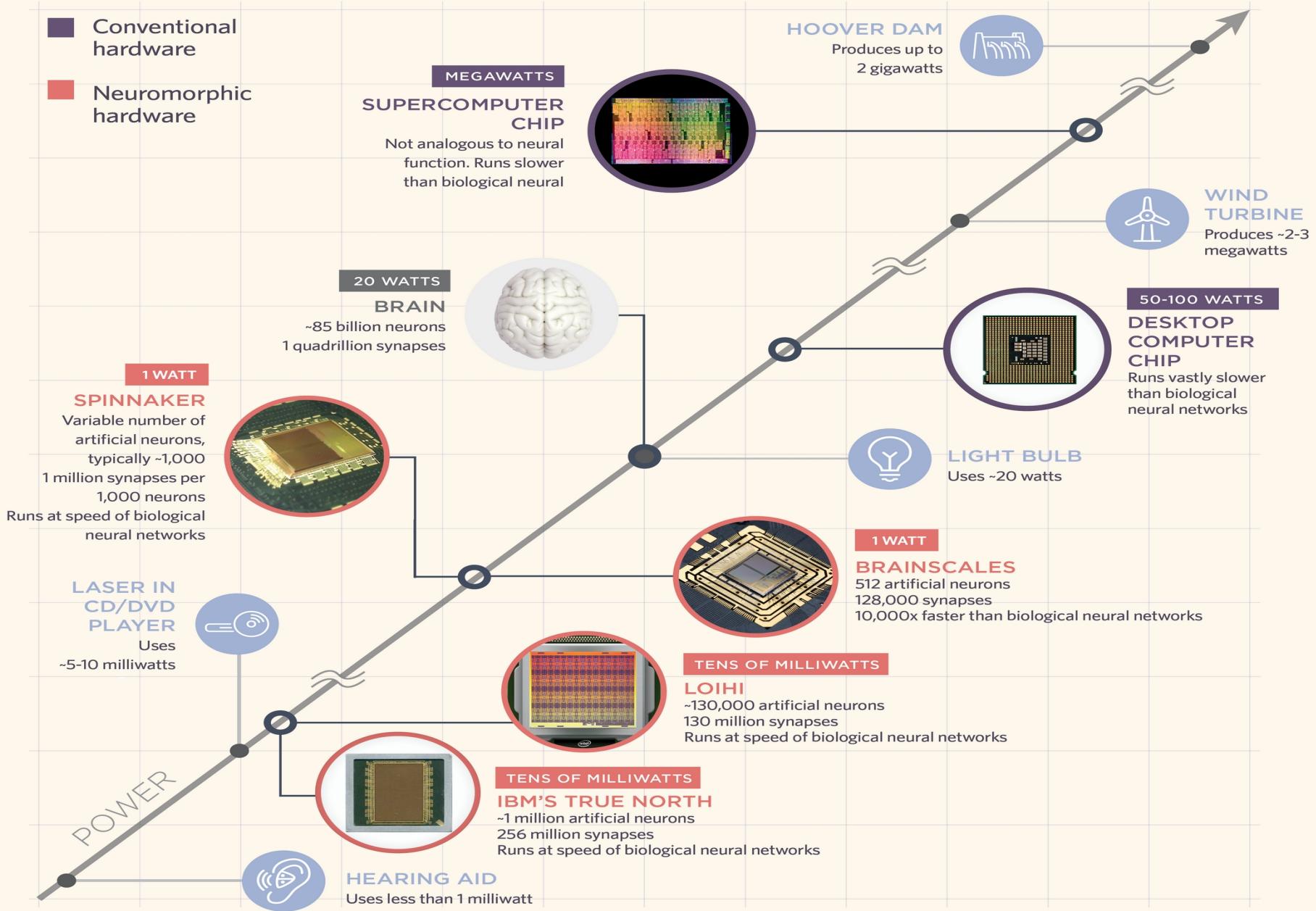
uProcessor  
100 B\$ @2020  
30% Cell Phone  
20% Embedded  
App  
50 Servers, PCs etc.

Computing

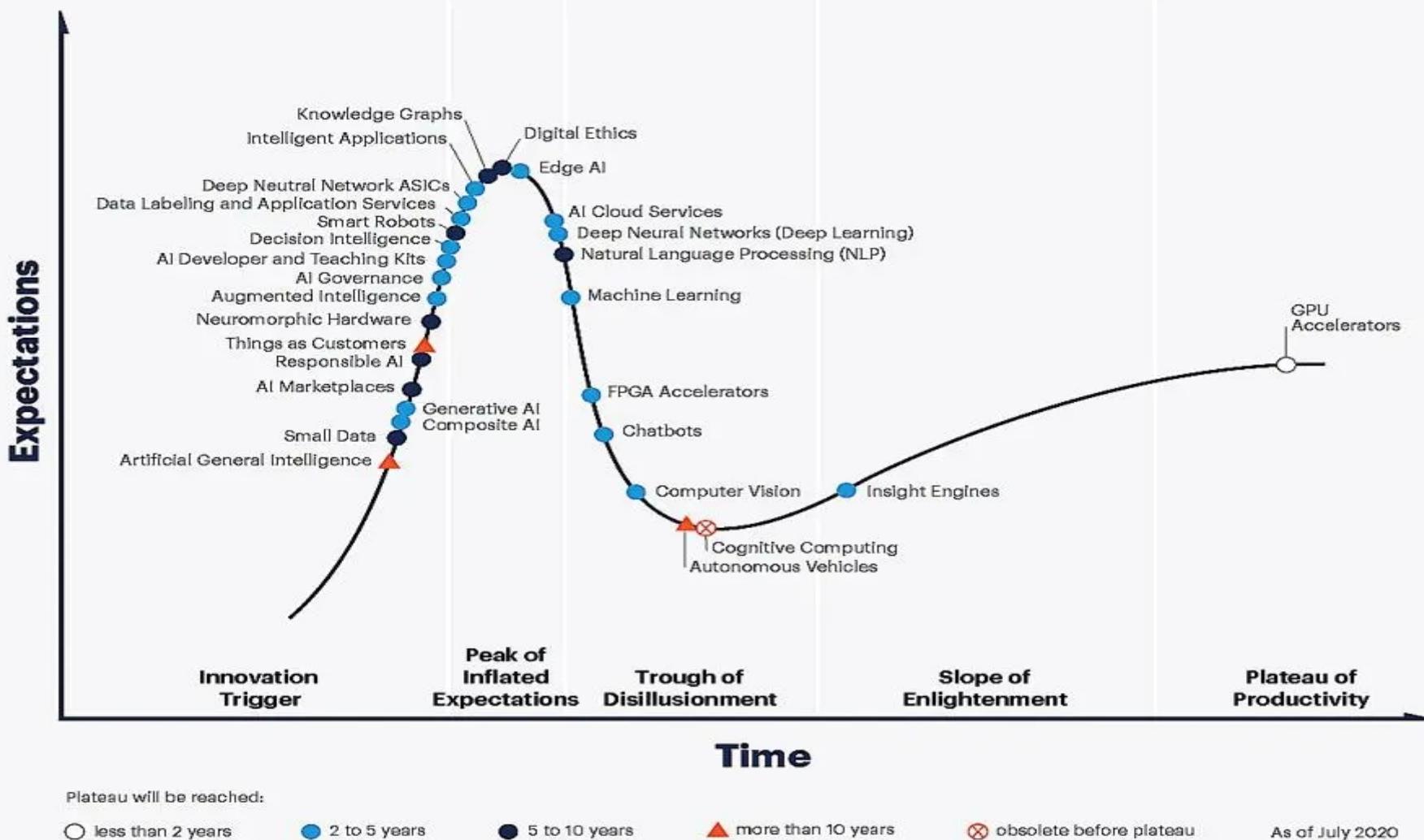
Digital Industrial Age  
5.5 Trillion \$ Revenue@2021







# Hype Cycle for Artificial Intelligence, 2020



[gartner.com/SmarterWithGartner](http://gartner.com/SmarterWithGartner)

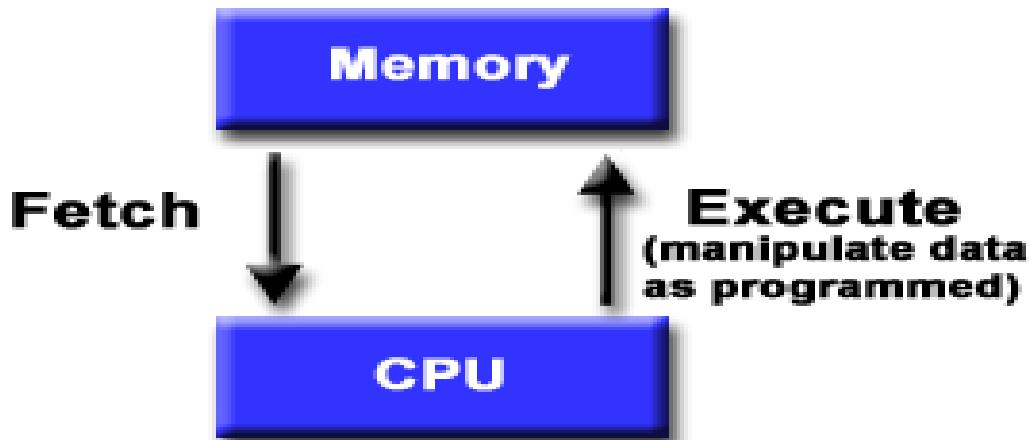
Source: Gartner  
© 2020 Gartner, Inc. and/or its affiliates. All rights reserved. Gartner and Hype Cycle are registered trademarks of Gartner, Inc. and its affiliates in the U.S.

**Gartner**

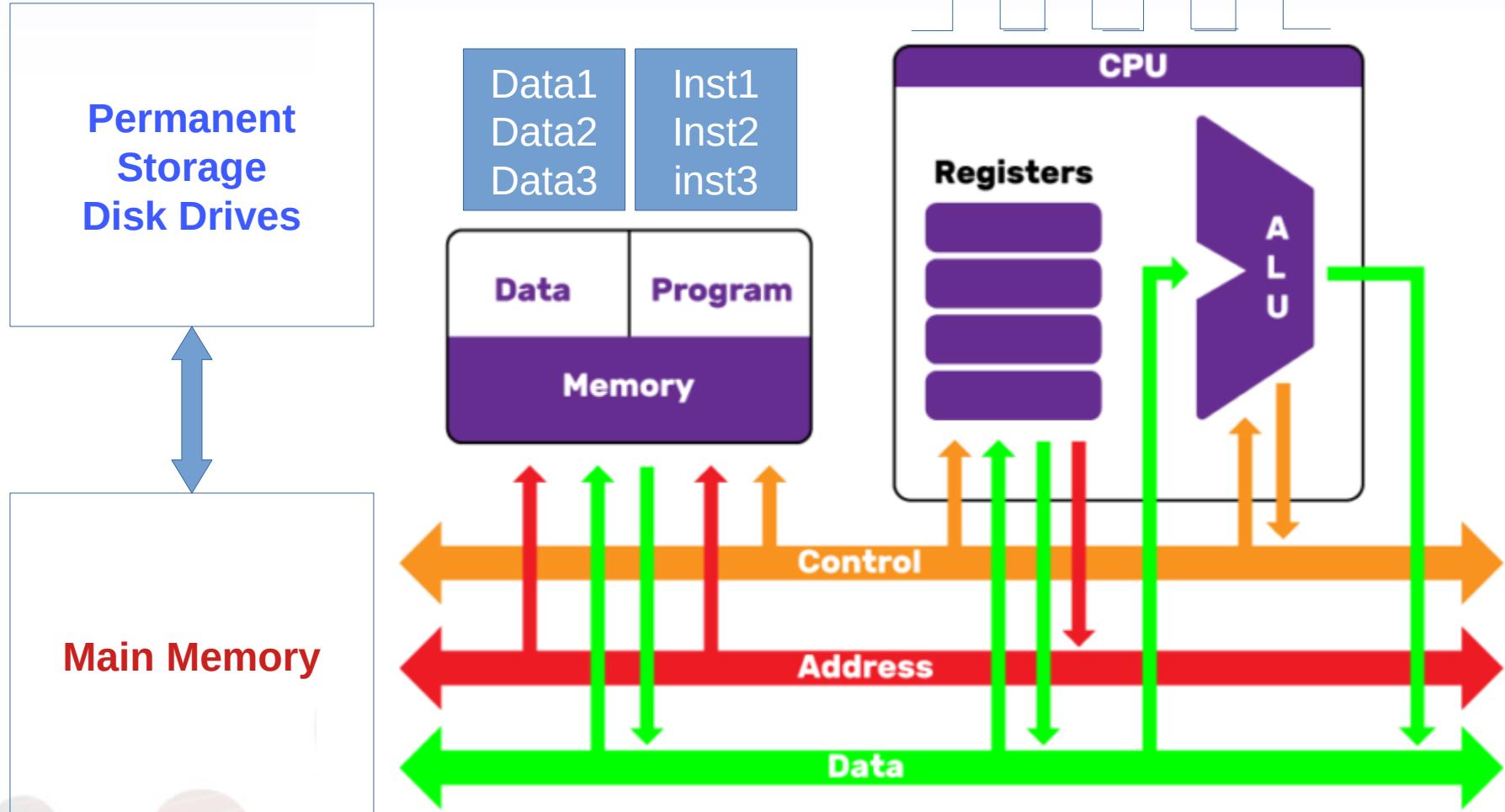
- Importance of HPC
- **Types of Processing System**
- HPC System Architecture
- Supercomputing Classes and Infrastructure
- Cluster Software Stacks
- Programming Models
- Demonstration (Supercomputing System)

# Basic Processor Architecture

- A central processing unit (CPU) gets instructions and/or data from memory, decodes the instructions and then **sequentially** performs them.
- Memory is used to store both program and data instructions
  - Program instructions are coded data which tell the computer to do something
  - Data is simply information to be used by the program



# Information and Computer



# Processors for Supercomputers

## **Microprocessor development directions:**

- Increasing of clock frequency and speed instruction stream processing
- Processing of large collection of data in single processor instruction - SIMD
- Control path multiplication – multi threading

## **Vector processors**

- NEC SX-6
- Cray (Cray X1)

## **RISC processors**

- MIPS
- IBM Power4
- Alpha

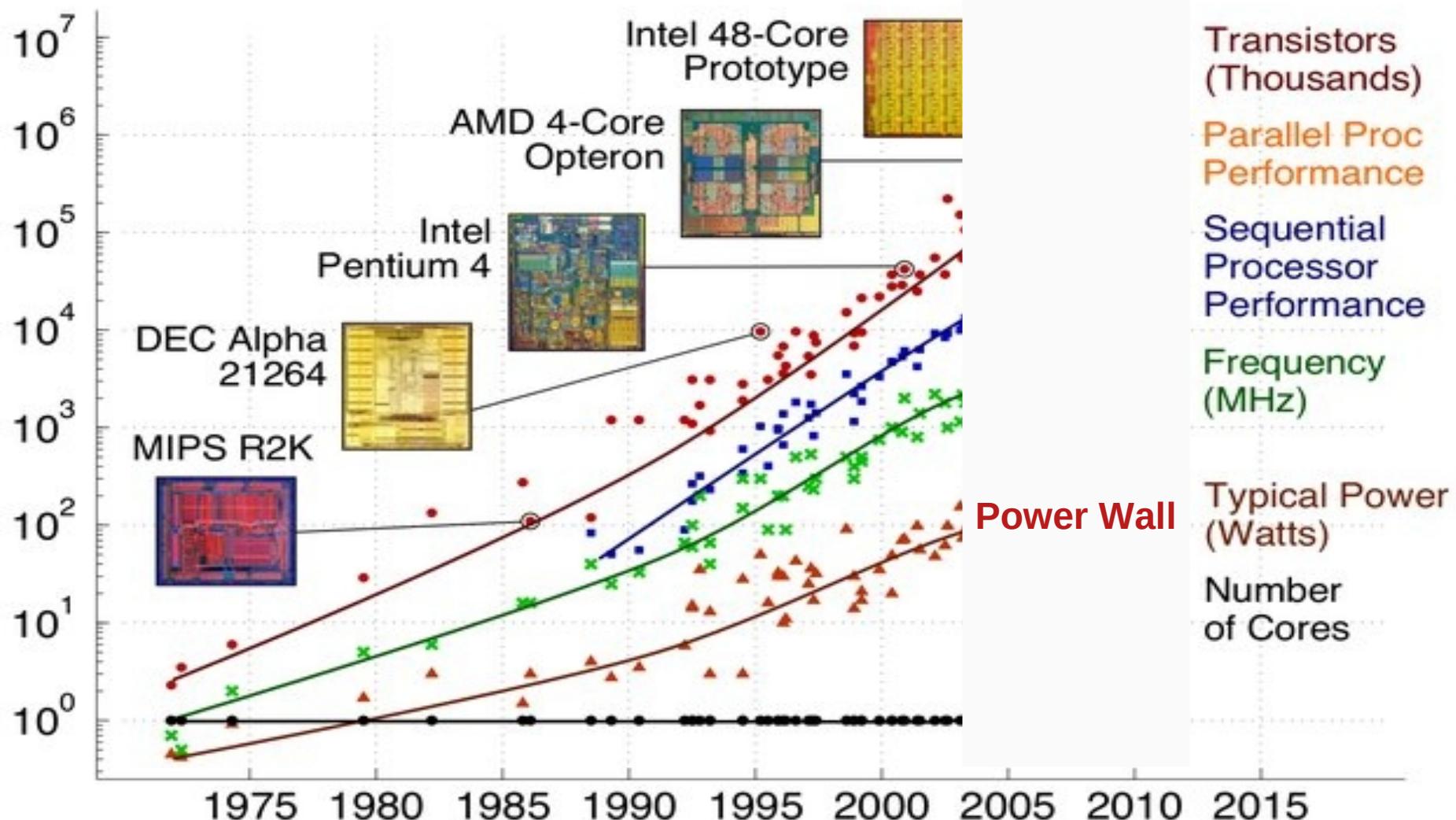
## **CISC processors**

- IA32
- AMD x86-64

## **VLIW processors**

- IA64

# Performance: Transistor Count, Frequency, Power and Multi Processor



- Importance of HPC
- Types of Processing System
- **HPC System Architecture (Parallel Processors)**
- Supercomputing Classes and Infrastructure
- Cluster Software Stacks
- Programming Models
- Demonstration (Supercomputing System)

# HPC System

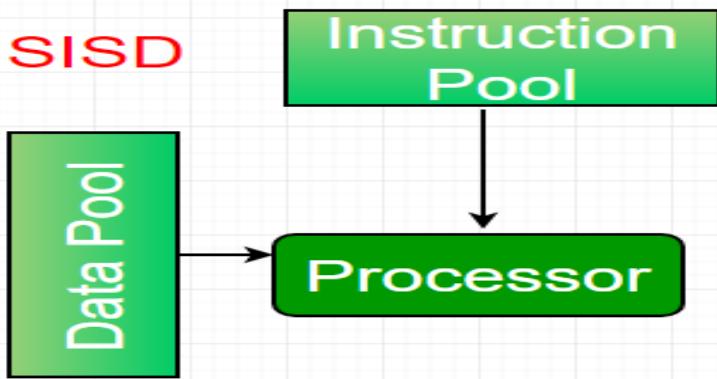
- Multi/Many Processor System Architecture
- Computer Architecture wrt Memory Arrangement
- Shared Storage Solutions
- Network Configuration

# Flynn's Classical Taxonomy

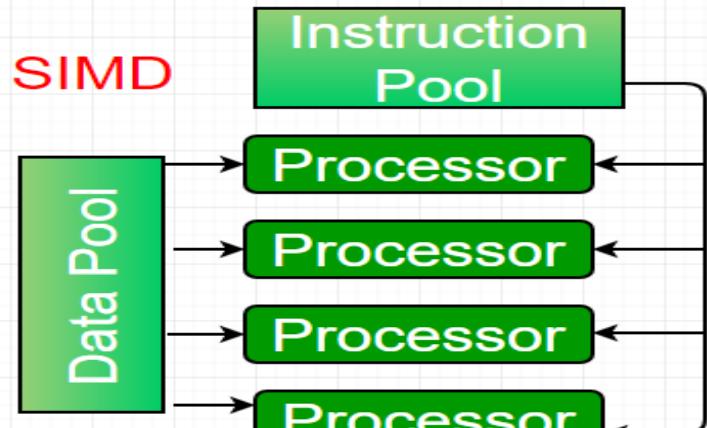
- There are different ways to classify parallel computers. One of the more widely used classifications, in use since 1966, is called Flynn's Taxonomy.
- Flynn's taxonomy distinguishes multi-processor computer architectures according to how they can be classified along the two independent dimensions of ***Instruction*** and ***Data***.
- Each of these dimensions can have only one of two possible states: ***Single*** or ***Multiple***.

# Processor Architectures

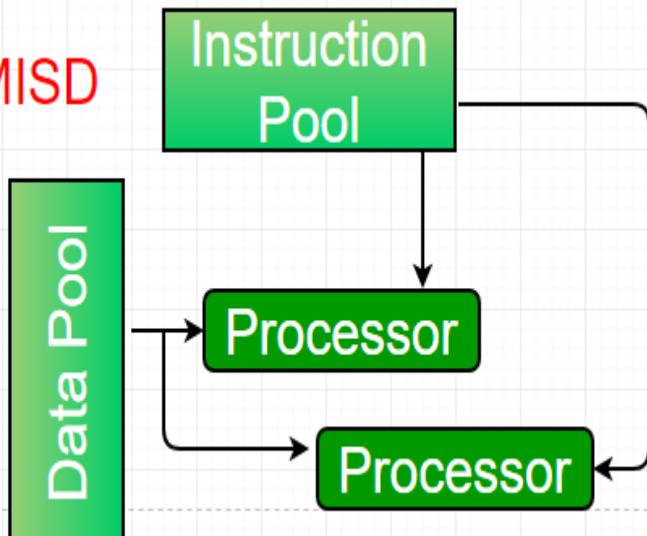
SISD



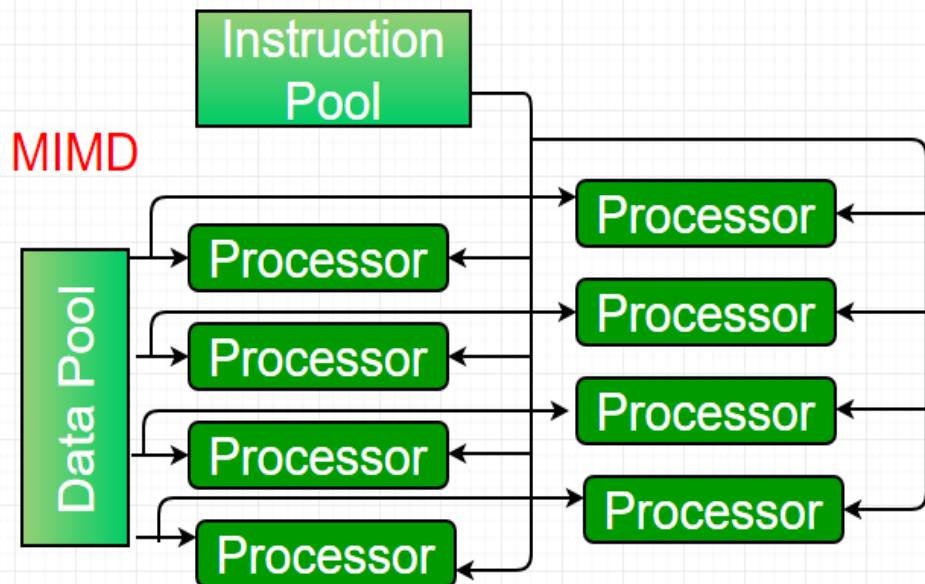
SIMD



MISD



MIMD

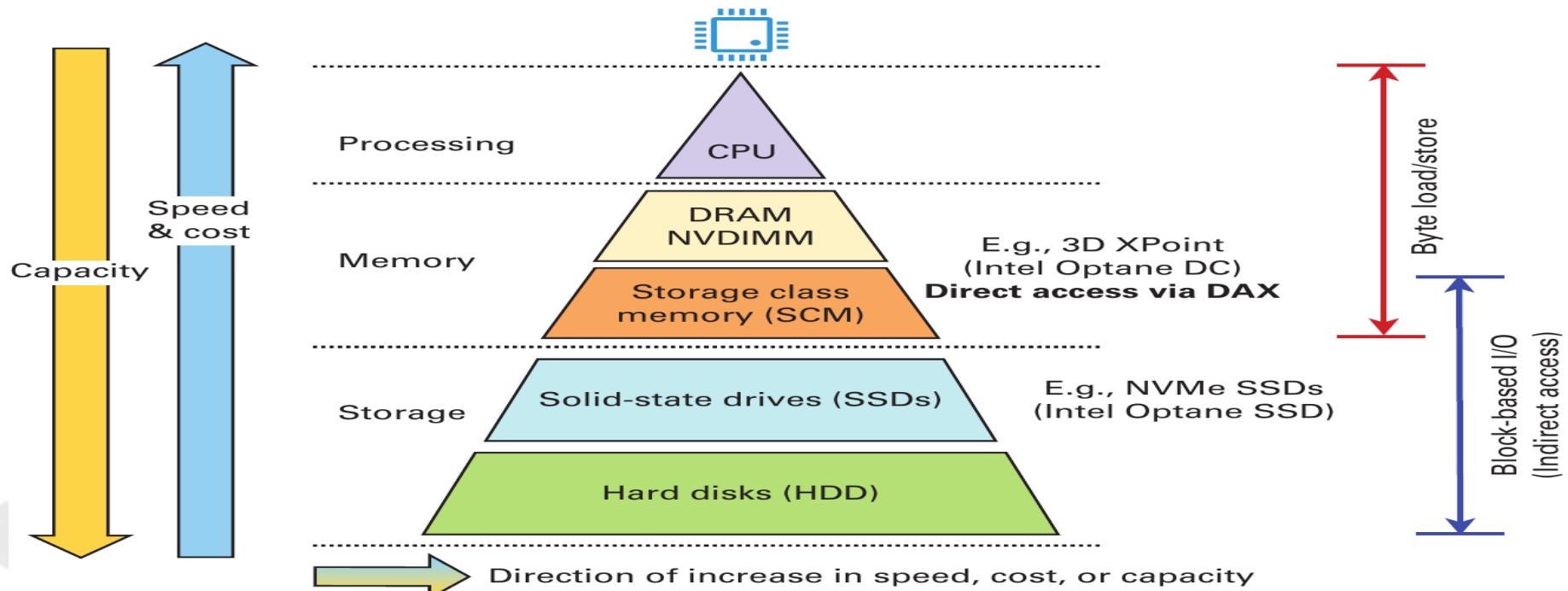


# Computer Architecture wrt Memory Arrangement

Shared-memory

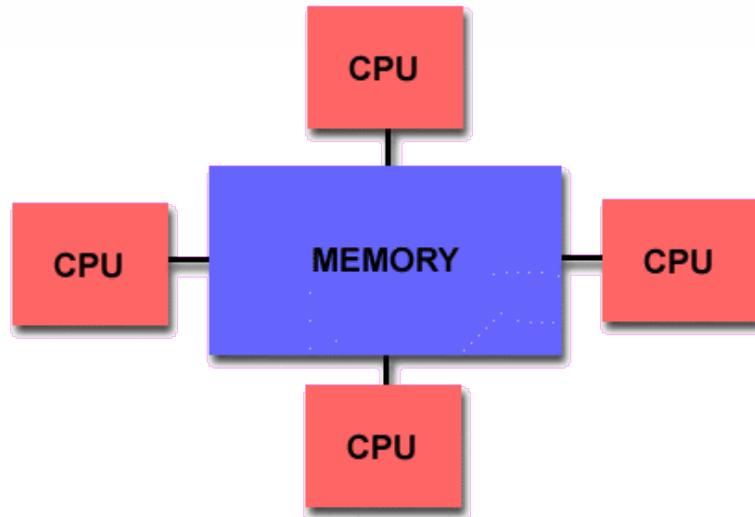
Distributed-memory

Distributed Shared-memory



# Shared Memory

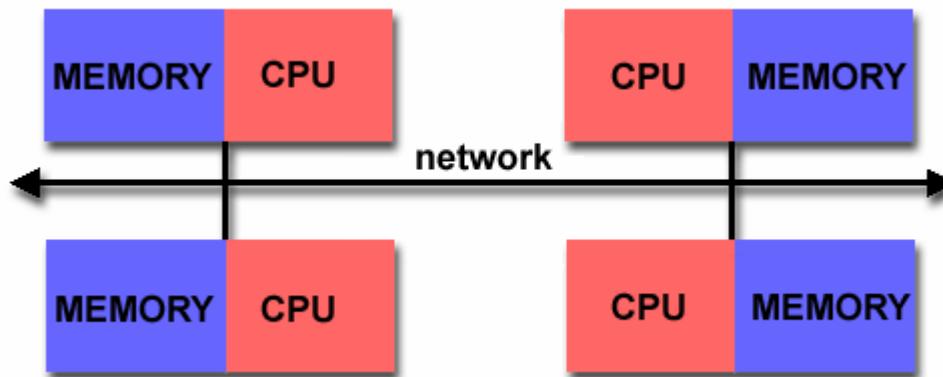
- Shared memory parallel computers vary widely, but generally have in common the ability for all processors to access all memory as global address space.



- Multiple processors can operate independently but share the same memory resources.
- Changes in a memory location effected by one processor are visible to all other processors.
- Shared memory machines can be divided into two main classes based upon memory access times: **UMA** and **NUMA**.

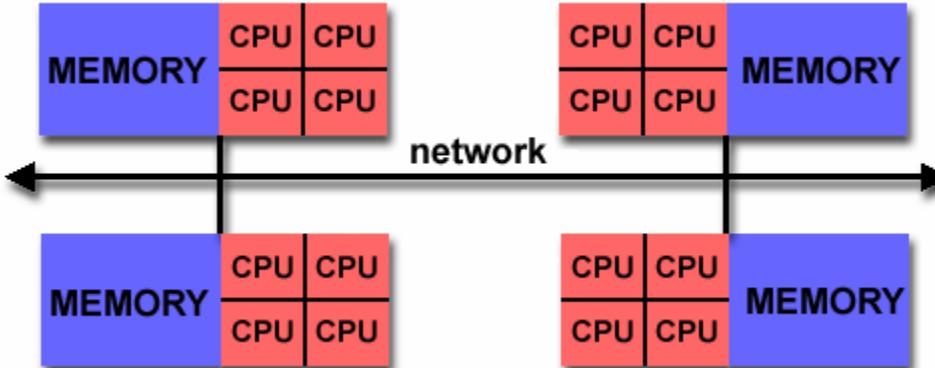
# Distributed Memory

- Like shared memory systems, distributed memory systems vary widely but share a common characteristic. Distributed memory systems require a communication network to connect inter-processor memory.
- Processors have their own local memory. Memory addresses in one processor do not map to another processor, so there is no concept of global address space across all processors.
- Because each processor has its own local memory, it operates independently. Changes it makes to its local memory have no effect on the memory of other processors. Hence, the concept of cache coherency does not apply.
- When a processor needs access to data in another processor, it is usually the task of the programmer to explicitly define how and when data is communicated. Synchronization between tasks is likewise the programmer's responsibility.
- The network "fabric" used for data transfer varies widely, though it can be as simple as Ethernet.



# Hybrid Distributed-Shared Memory

- The largest and fastest computers in the world today employ both shared and distributed memory architectures.



- The shared memory component is usually a cache coherent SMP machine. Processors on a given SMP can address that machine's memory as global.
- The distributed memory component is the networking of multiple SMPs. SMPs know only about their own memory - not the memory on another SMP. Therefore, network communications are required to move data from one SMP to another.
- Current trends seem to indicate that this type of memory architecture will continue to prevail and increase at the high end of computing for the foreseeable future.
- Advantages and Disadvantages: whatever is common to both shared and distributed memory architectures.

# Types: Memory Access

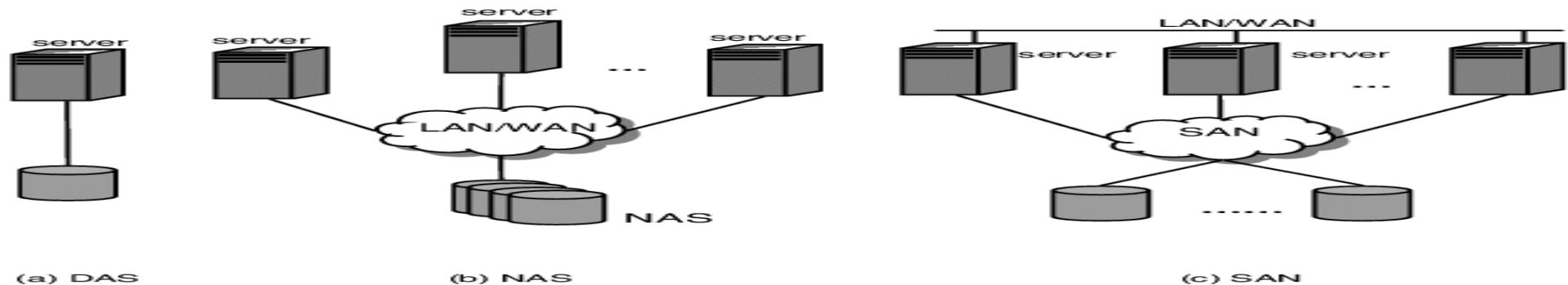
Unified Memory Access (UMA)

Non Unified Memory Access (NUMA)

- Cache-Coherent NUMA (ccNUMA)
- Distributed Memory (DM) / MPI
- In-Memory Computing
- NUMA Hybrid Architectures

# Shared Disk Storage

- DAS (Direct Attached Storage)
- NAS (Network Attached Storage)
- SAN (Storage Area Network)
- Parallel File System Storage
- Object Storage:
- Tiered Storage: SSD and HDD
- Ceph : Object, Blocks, File Storage



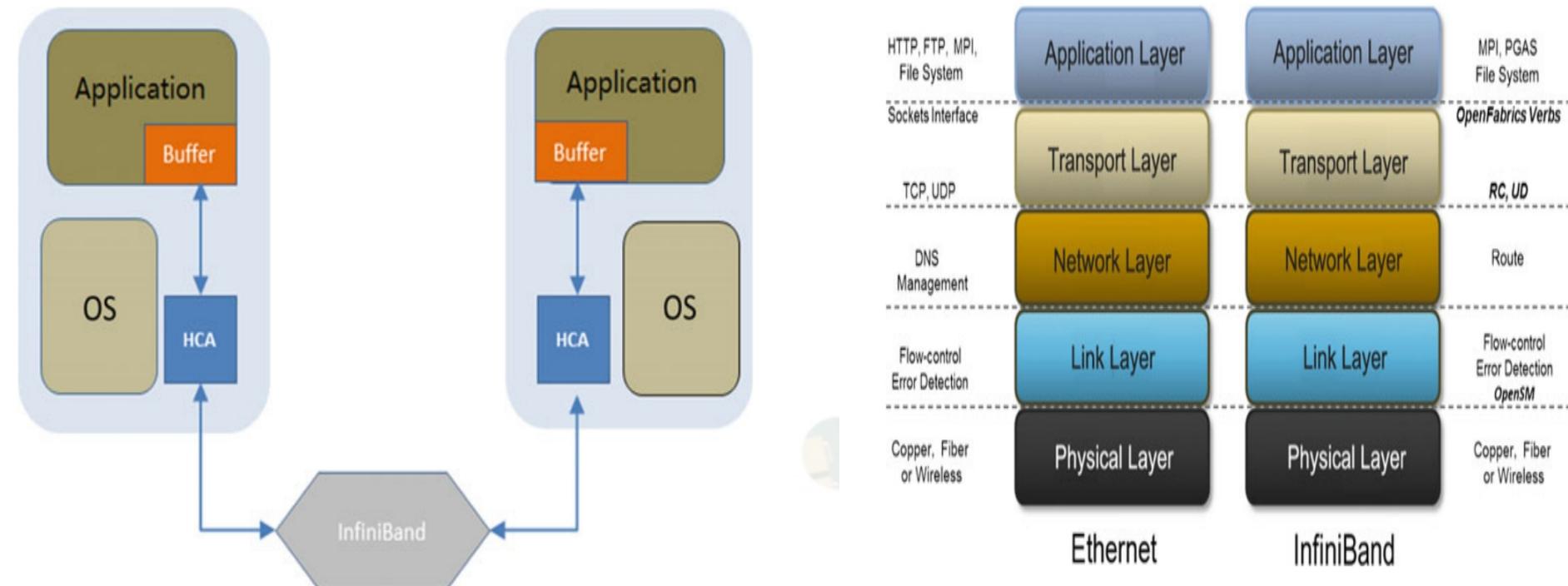
# Comparison: SAN, NAS, PFS

Criteria	Storage Area Network (SAN)	Network Attached Storage (NAS)	Parallel File System
Access Type	Block-level access	File-level access	File-level access with parallel capabilities
Performance	High-performance, low-latency	Varies based on NAS solution	High-performance, scalable
Workload Flexibility	Well-suited for demanding workloads	Suitable for general workloads	Designed for high-performance computing
Scalability	Scalable, can handle growth	Varies based on NAS solution	Highly scalable and optimized for large clusters
Complexity	Can be complex to set up and manage	Relatively simple to set up and manage	Can be complex to optimize for specific workloads
Isolation	Can provide isolated storage access	Shared access across clients	Shared access, optimized for parallelism
Hardware Requirement	Requires dedicated SAN hardware (FC)	NAS devices with Ethernet connectivity	Requires specialized parallel file system software
Cost	Higher cost due to dedicated infrastructure	Generally cost-effective	Costs can vary based on solution and scale
Management	May require specialized expertise	Generally easier to manage	Requires expertise in parallel file system management
Protocols Supported	Typically Fibre Channel, iSCSI, etc.	NFS, SMB, other file protocols	NFS, Lustre, GPFS, BeeGFS, etc.
Use Cases	High-performance applications	General-purpose file sharing	High-performance computing, data-intensive tasks

# Network Communication Protocols

InfiniBand is a high-speed networking technology primarily designed for low-latency, high-bandwidth communication within data centers and high-performance computing (HPC) environments.

InfiniBand excels in short to medium-range communication within a data center or between closely located computing clusters.



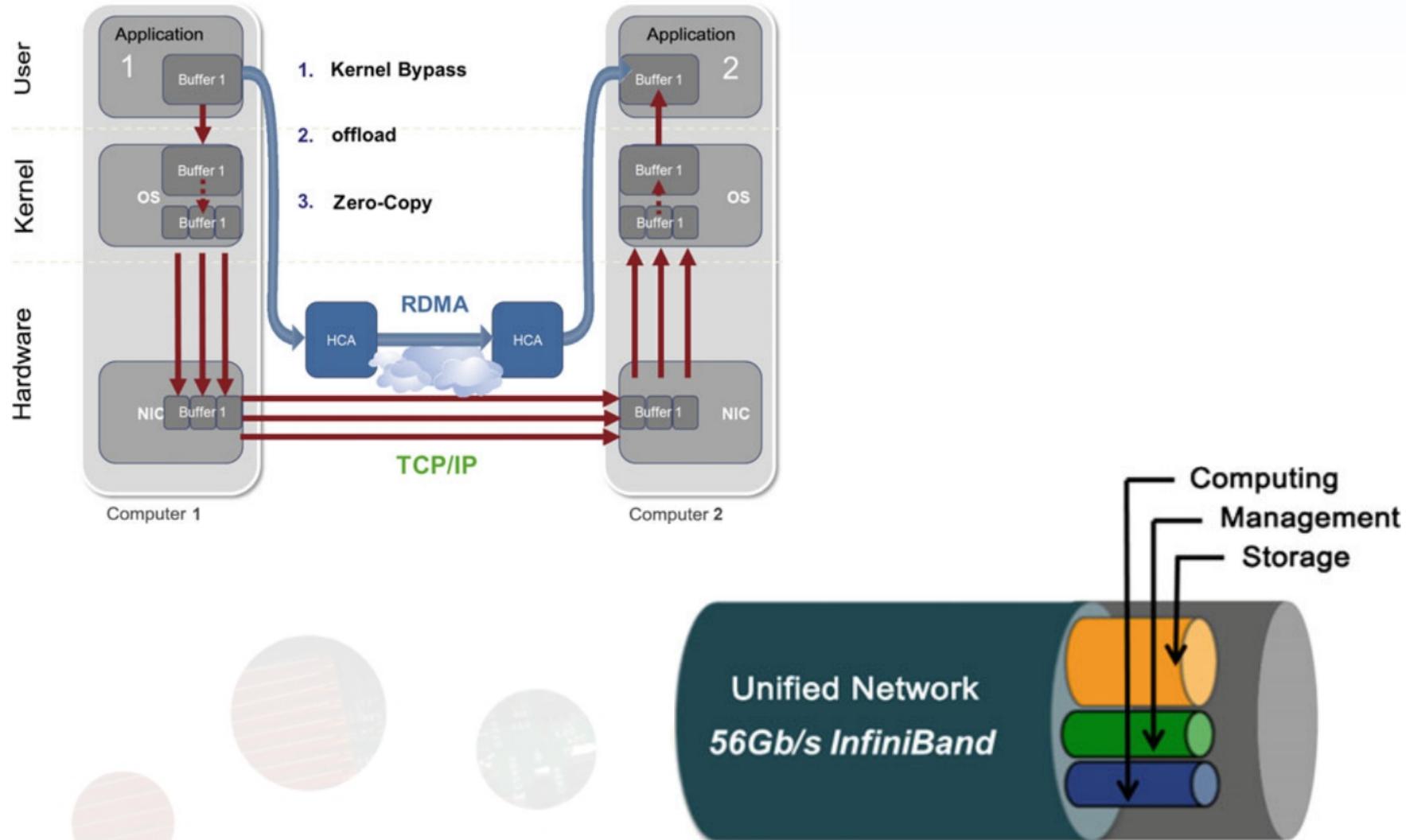
# RDMA and HCA

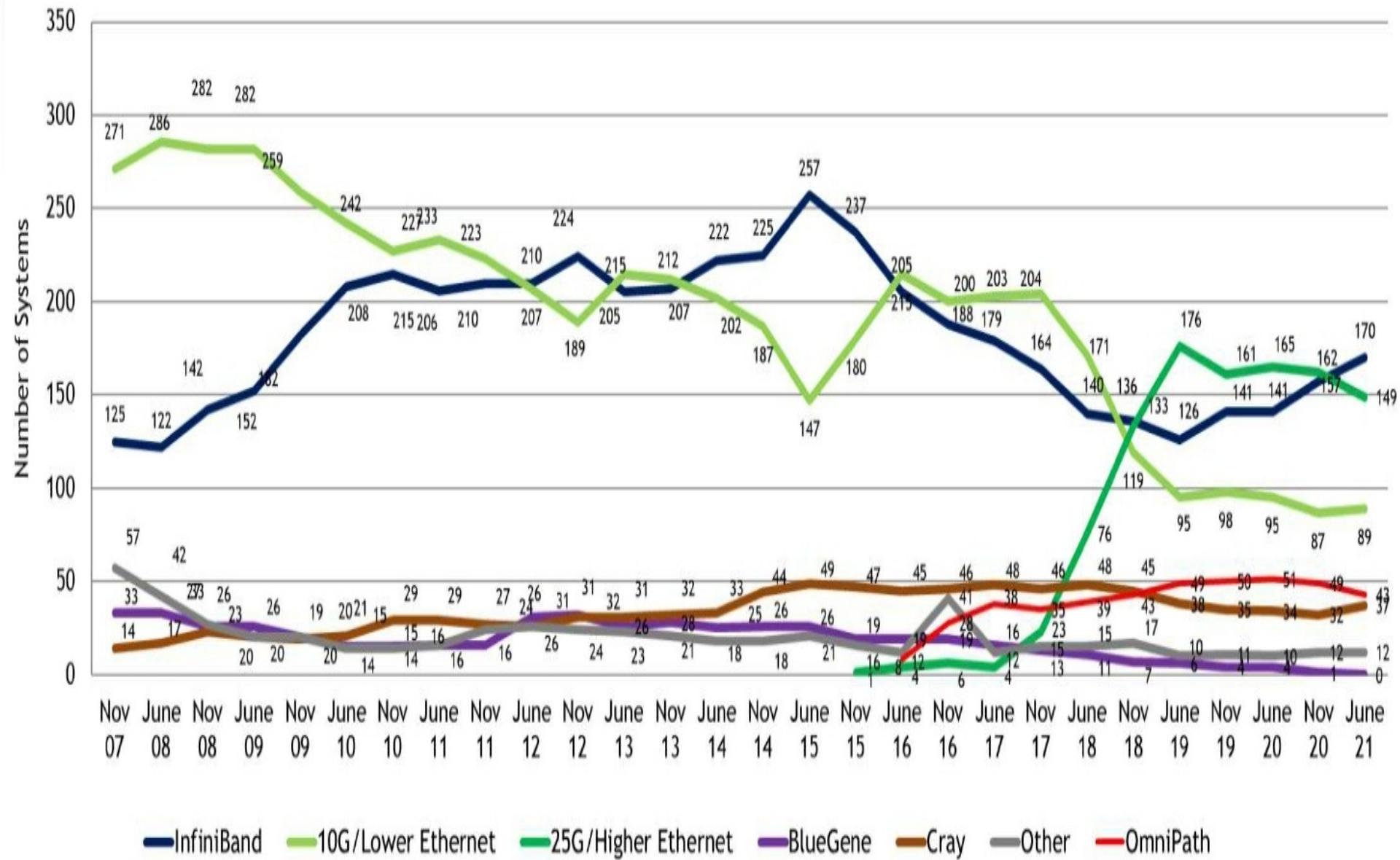
Remote Direct Memory Access (RDMA) and Host Channel Adapter (HCA) allows computers of Cluster to access each other's memory directly without involving the host CPU.

It executes **data offloading** and **data transfer tasks** from the host to dedicated network adapters that support **RDMA**, such as **InfiniBand or RoCE (RDMA over Converged Ethernet)**.

This bypasses the need for the usual network stack processing, resulting in reduced latency and improved data throughput.

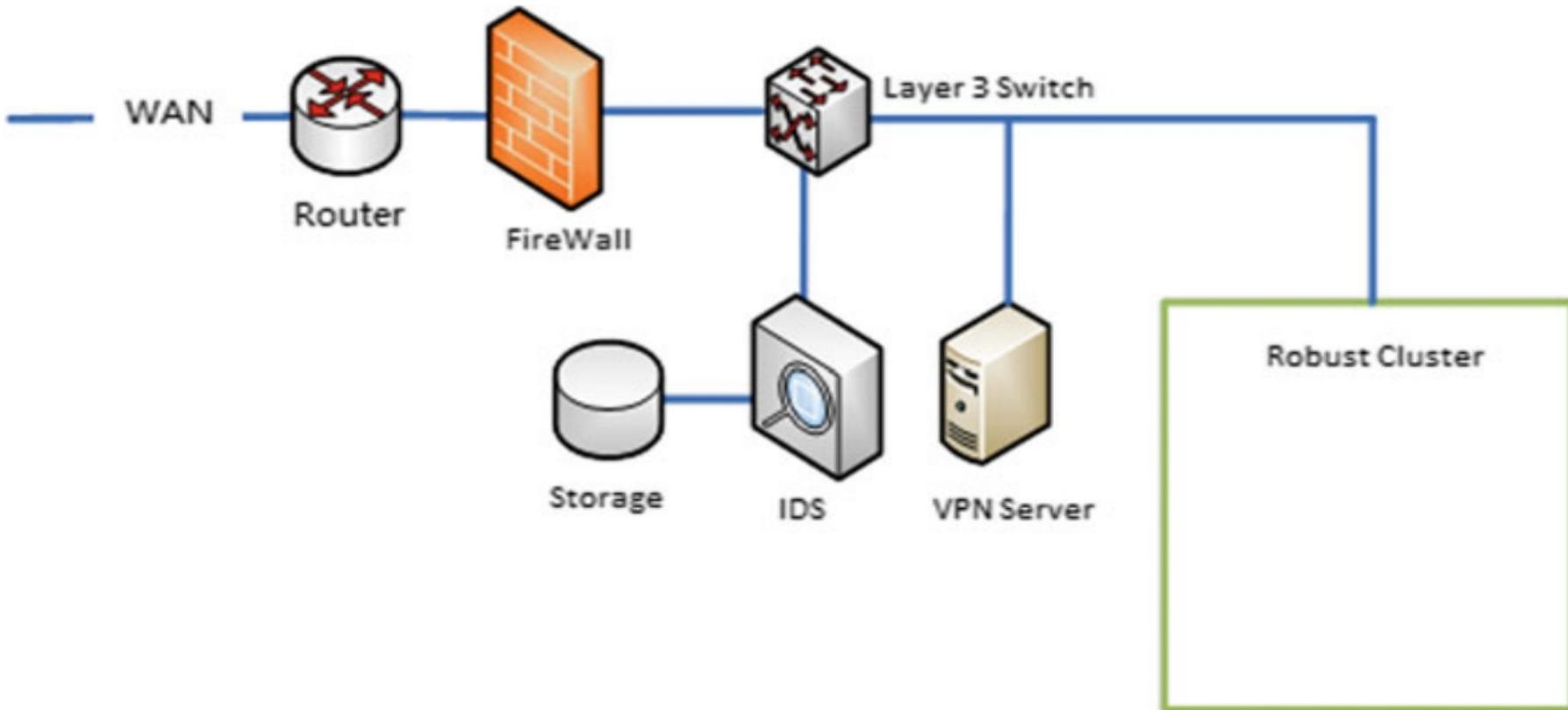
# TCP/IP and Remote Direct Memory Access





# Network Security

Soft Firewall: Zentyal, pfSense, IPFire, SmoothWall, ClearOS, and Iptables.



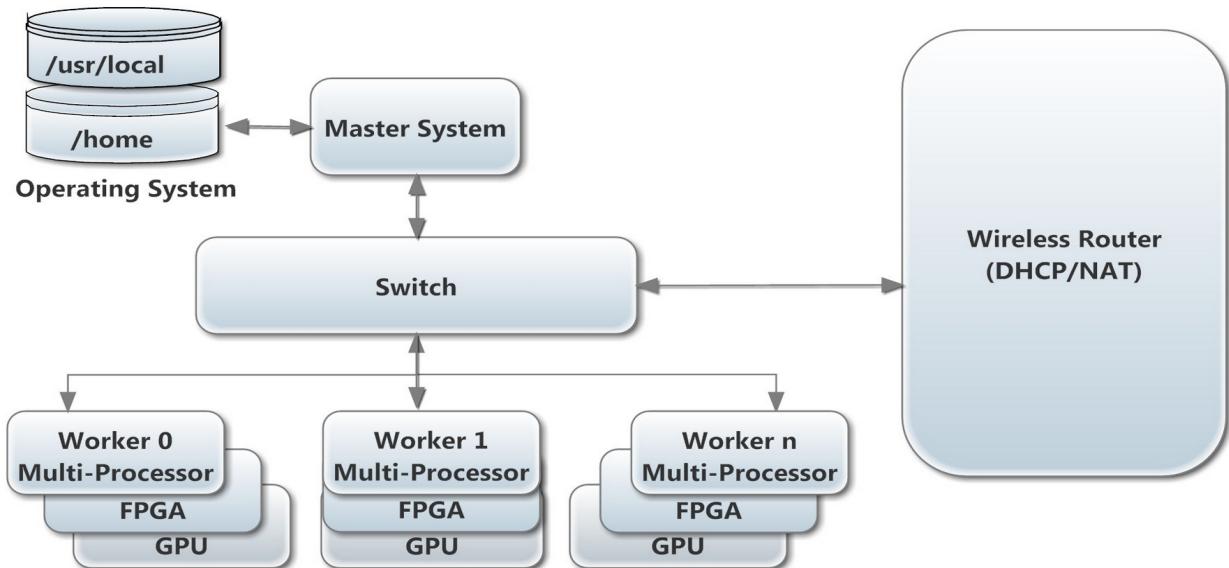
- Importance of HPC
- Types of Processing System
- HPC System Architecture
- **Supercomputing Classes and Infrastructure**
- Cluster Software Stacks
- Programming Models
- Demonstration (Supercomputing System)

# What is a Supercomputer?

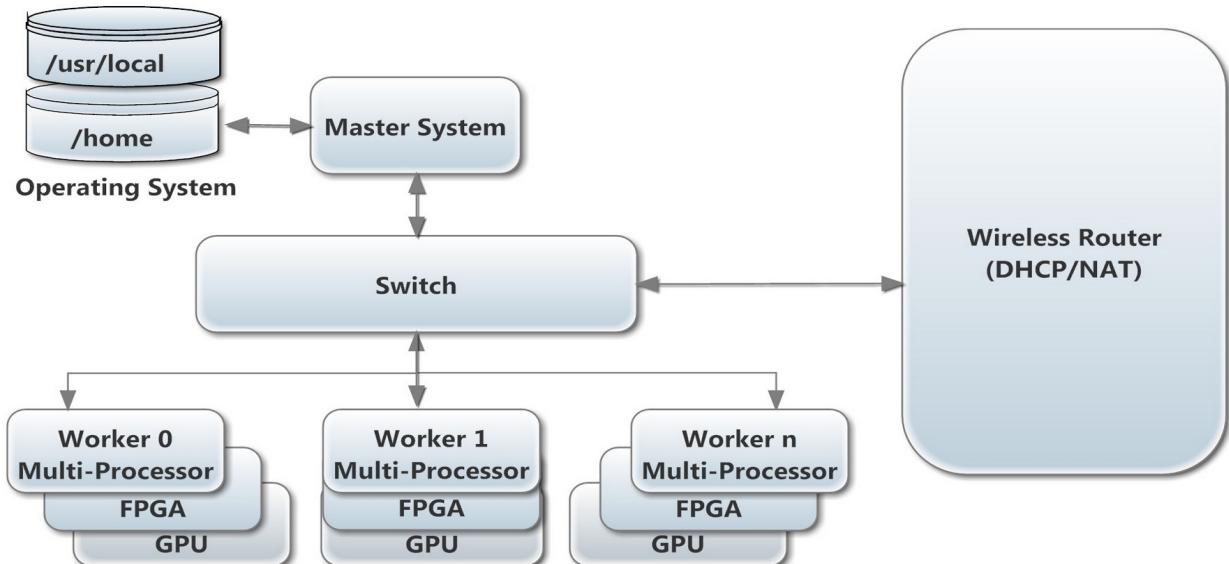
*“A supercomputer is a device for turning compute-bound problems into I/O-bound problem”* - Seymour Cray

A supercomputer is a computer system that leads the world in terms of processing capacity, particularly speed of calculations, at the time of its introduction. (Wikipedia)

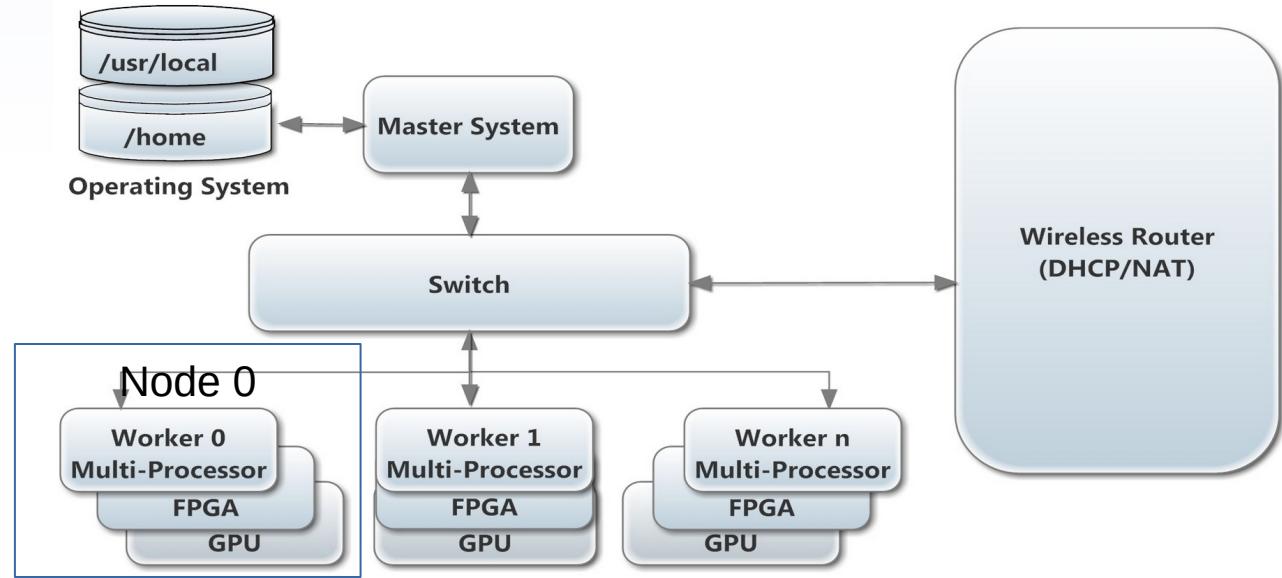
# Basic Arch of Supercomputer



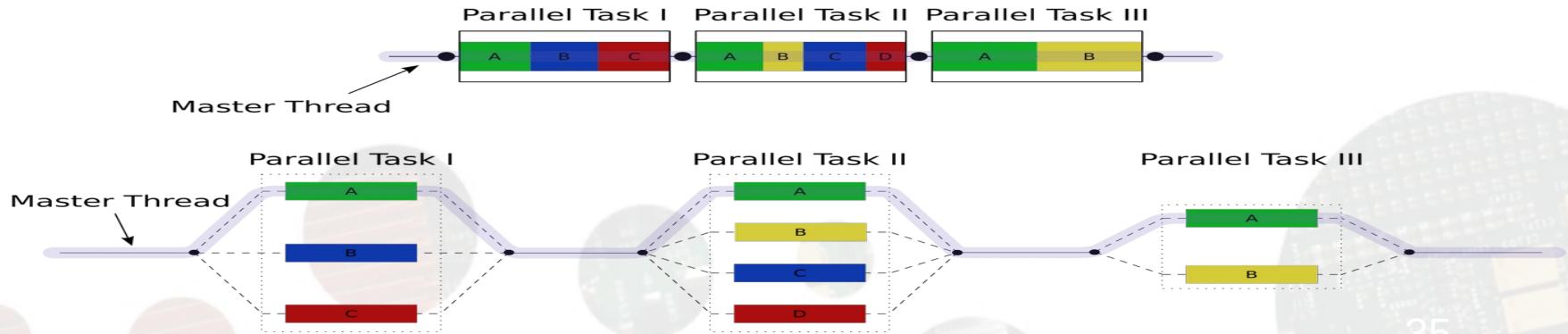
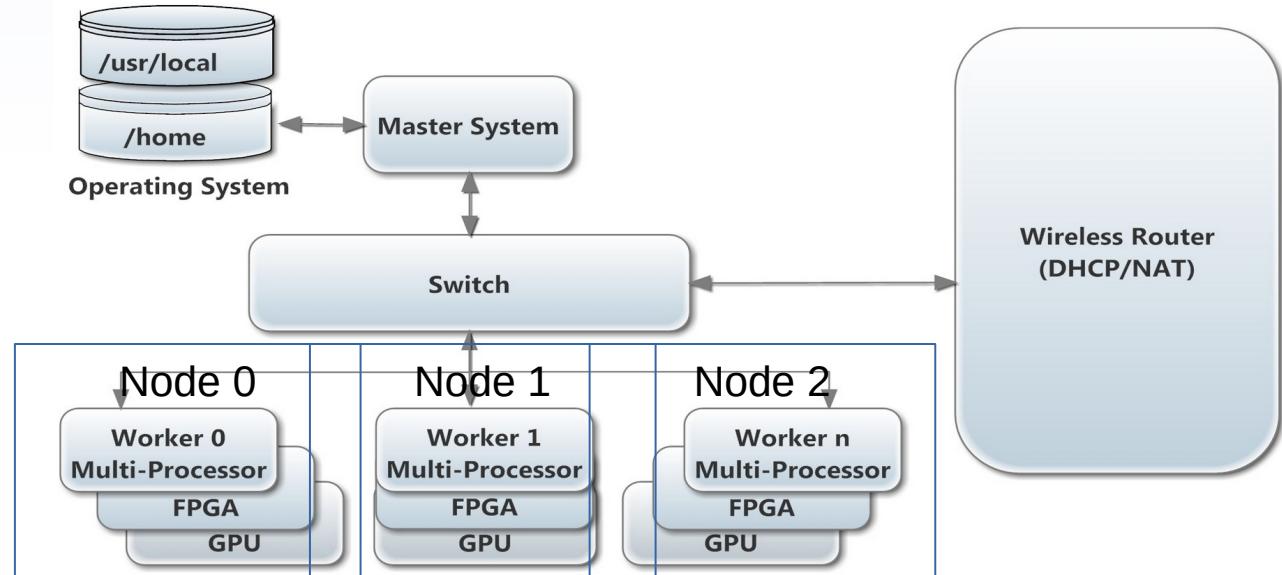
# Basic Arch of Supercomputing



# Basic Arch of Supercomputing

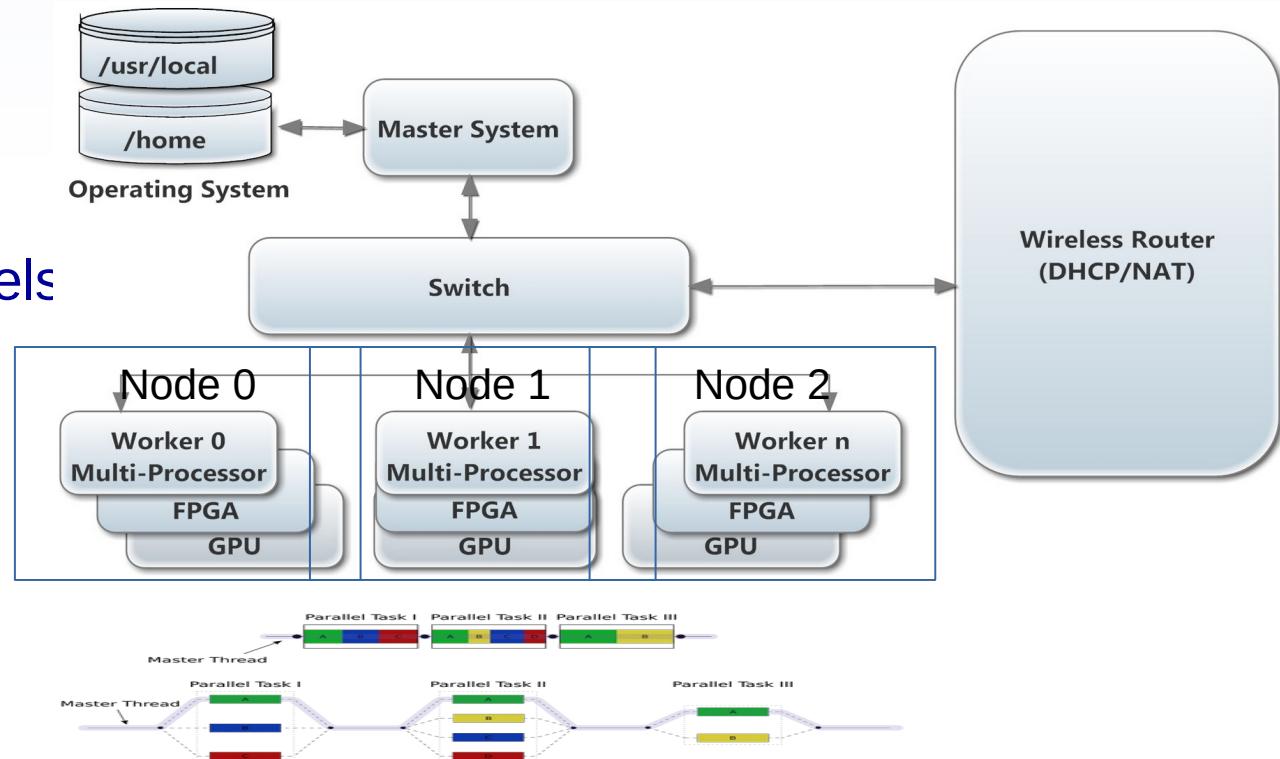


# Basic Introduction of Supercomputing



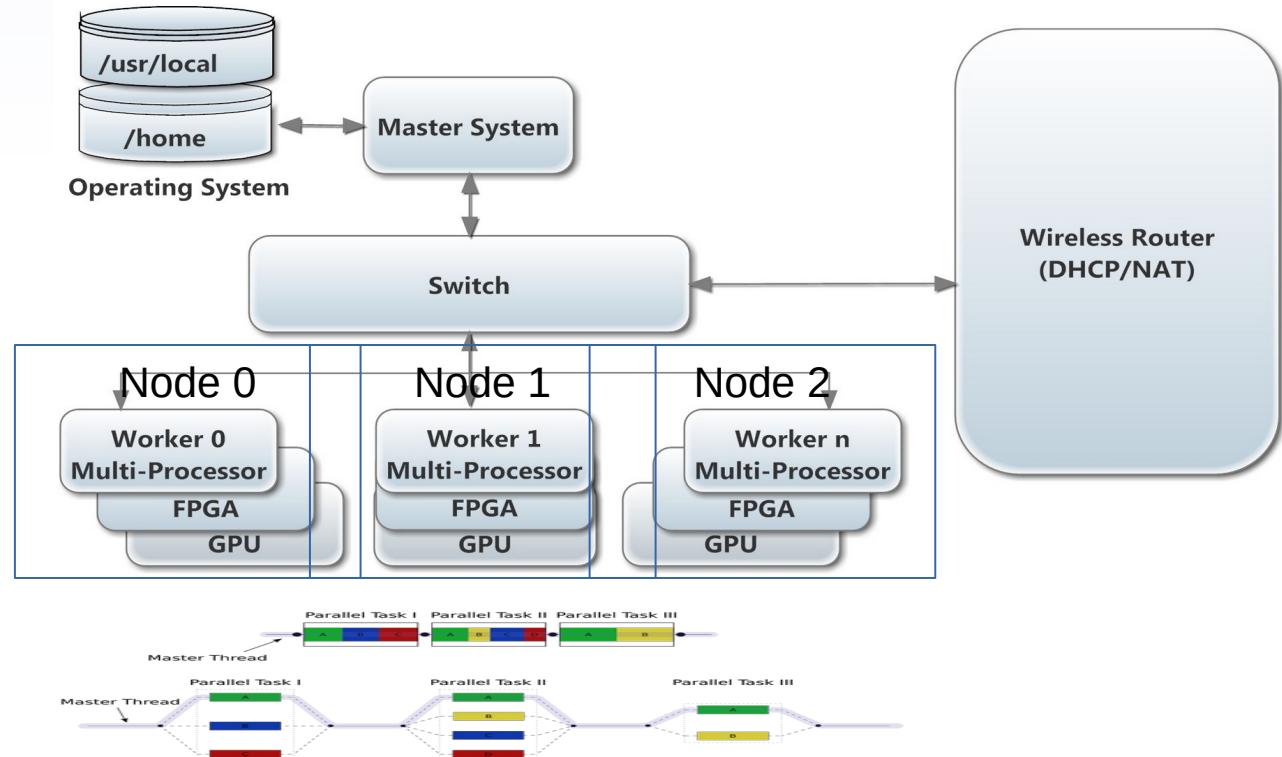
# Constituents

- Processing
- Software
- Programming Models
- Storage
- Network
- Accessibility
- Power



# Basic Introduction of Supercomputing

- Performance
- Programmability
- Portability
- Scalability
- Accessibility
- Usability
- Power
- Cost



FLOPS = Floating Point Operation Per Second  
KFLOPS =  $10^3$   
= One Thousand Computation Per Second  
=  $12.00 * 1212.222 * 21212 + 232323 \dots$

MFLOPS =  $10^6$  Million Computation Per Second  
GFLOPS =  $10^9$  Billion  
TFLOPS =  $10^{12}$  Trillion  
PFLOPS =  $10^{15}$  Quadrillion  
EFLOPS =  $10^{18}$  Quintillion  
ZFLOPS =  $10^{21}$  Sextillion

# History of Supercomputers

1945-50 - Manchester Mark I

1950-55 - MIT Whirlwind

1955-60 - IBM 7090 - **210 KFLOPS**

1960-65 - CDC 6600 - **10.24 MFLOPS**

1965-70 - CDC 7600 - **32.27 MFLOPS**

1970-75 - CDC Cyber **76 MFLOPS**

1975-80 - Cray-1 - **160 MFLOPS**

1980-85 - Cray X-MP - **500 MFLOPS**

1985-90 - Cray Y-MP - **1.3 GFLOPS**

1990-95 - Fujitsu Numerical Wind Tunnel - **236 GFLOPS**

1995-00 - Intel ASCI Red - **2.150 TFLOPS**

2000-02 - IBM ASCI White, SP Power3 375 MHz - **7.226 TFLOPS**

2002-03 - NEC Earth Simulator - **35 TFLOPS**

# Supercomputing Generations

Gen I: 1970s. SIMD Array of Processors.

- USAs developed ILLIAC-IV processor array, consisting of an 8 x 8 array of 64 processors. 106 MFLOPS

Gen II: Cray-1 Pipelined Vector Machines 1976, 12 pipelined arithmetic units, 160 MFLOPS

Gen III: Shared-Memory multi-processor systems (MIMD Arch) 40 TFLOPS

Gen IV: Massive Parallel Processor having 10 to thousands of processors

Gen V: Cluster based (CPU+GPU+MIC)

# Classes of Supercomputers

## General Purpose

**vector processing machines** - the same operation carried out on a large amount of data simultaneously

**Tightly connected cluster computers (NUMA)** - communication oriented architectures engineered from ground up, based on high speed interconnects and large number of processors

**Commodity clusters** - collection of large number of commodity PCs (COTS) interconnected by high-bandwidth low-latency network

**Special Purpose:** high performance computing devices with a hardware architecture dedicated to solve a single problem (equipped with custom ASICs or FPGA chips)

# Commodity Cluster

A cluster is a collection of connected, independent computers working in unison to solve a problem COTS technology nodes are interconnected by Ethernet LAN, Myrinet, QsNet ELAN etc. computation can be performed by using popular programming toolkits and frameworks: OpenMP, MPI clusters require dedicated management software.

# Components of a Supercomputing System

- Compute Nodes
- Login/Access Nodes
- I/O Nodes
- Network Nodes (Switches/Routers)
- Management Nodes
- Storage Nodes
- Service Nodes

## **Compute Nodes:**

These nodes perform the actual computational tasks and simulations.

Equipped with powerful CPUs and GPUs for parallel processing.

The number of compute nodes can vary based on the expected workload and desired level of parallelism.

## **Login/Access Nodes:**

Provide remote access for users to connect to the cluster.

Users log in to these nodes to submit jobs, manage files, and interact with the cluster.

Not meant for computation but facilitate user interaction.

## **Management Nodes:**

Host cluster management software, schedulers, and monitoring tools.

Manage resource allocation, job scheduling, system health monitoring, and overall cluster administration.

## **Storage Nodes:**

Manage the storage infrastructure of the cluster.

Provide access to high-performance storage systems, such as parallel file systems or distributed storage solutions.

Store large datasets and simulation results. crucial for parallel processing and data sharing.

## **I/O Nodes:**

Handle data input/output operations efficiently.

Optimize data access and transfer between compute nodes and storage systems.

Reduce I/O bottlenecks and latency for data-intensive applications.

## **Network Nodes (Switches/Routers):**

High-speed interconnects (e.g., InfiniBand, Ethernet with RDMA) enable fast communication between compute nodes.

Low-latency, high-bandwidth communication is

## **Service Nodes:**

Host essential services for cluster operation and user interaction.

Services may include Light Weight Directory Access Protocol authentication, DNS, NTP (Network Time Protocol), and other services needed for system stability.

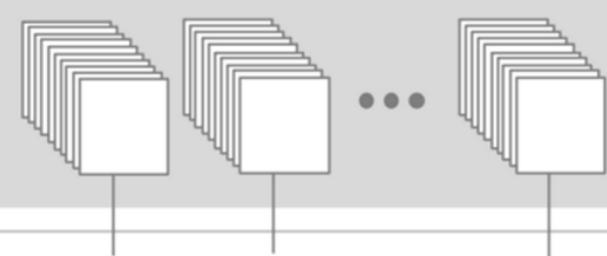
## **Compute Node Interconnect:**

High-performance network interconnects between compute nodes facilitate communication for parallel processing and distributed computing.

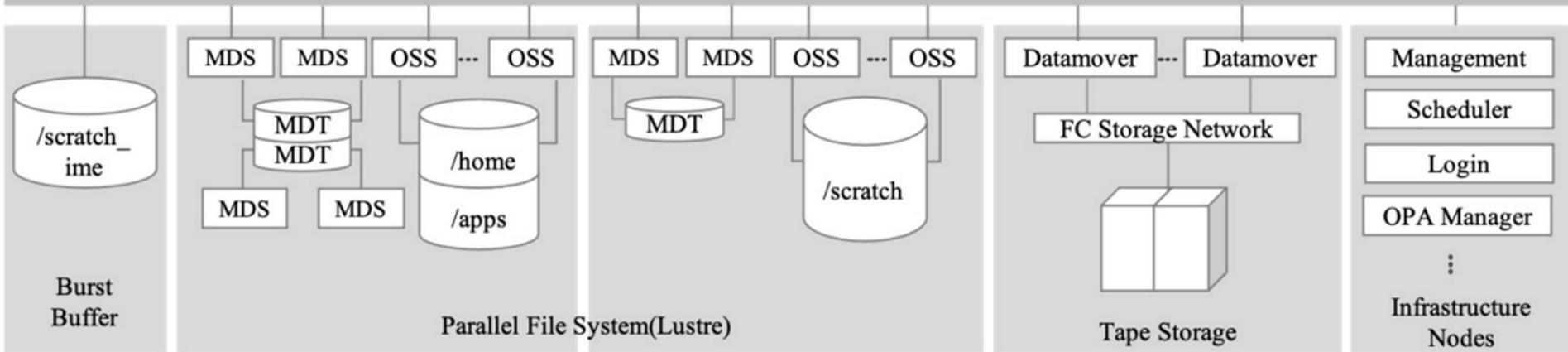
Compute Nodes-KNL(8,305)



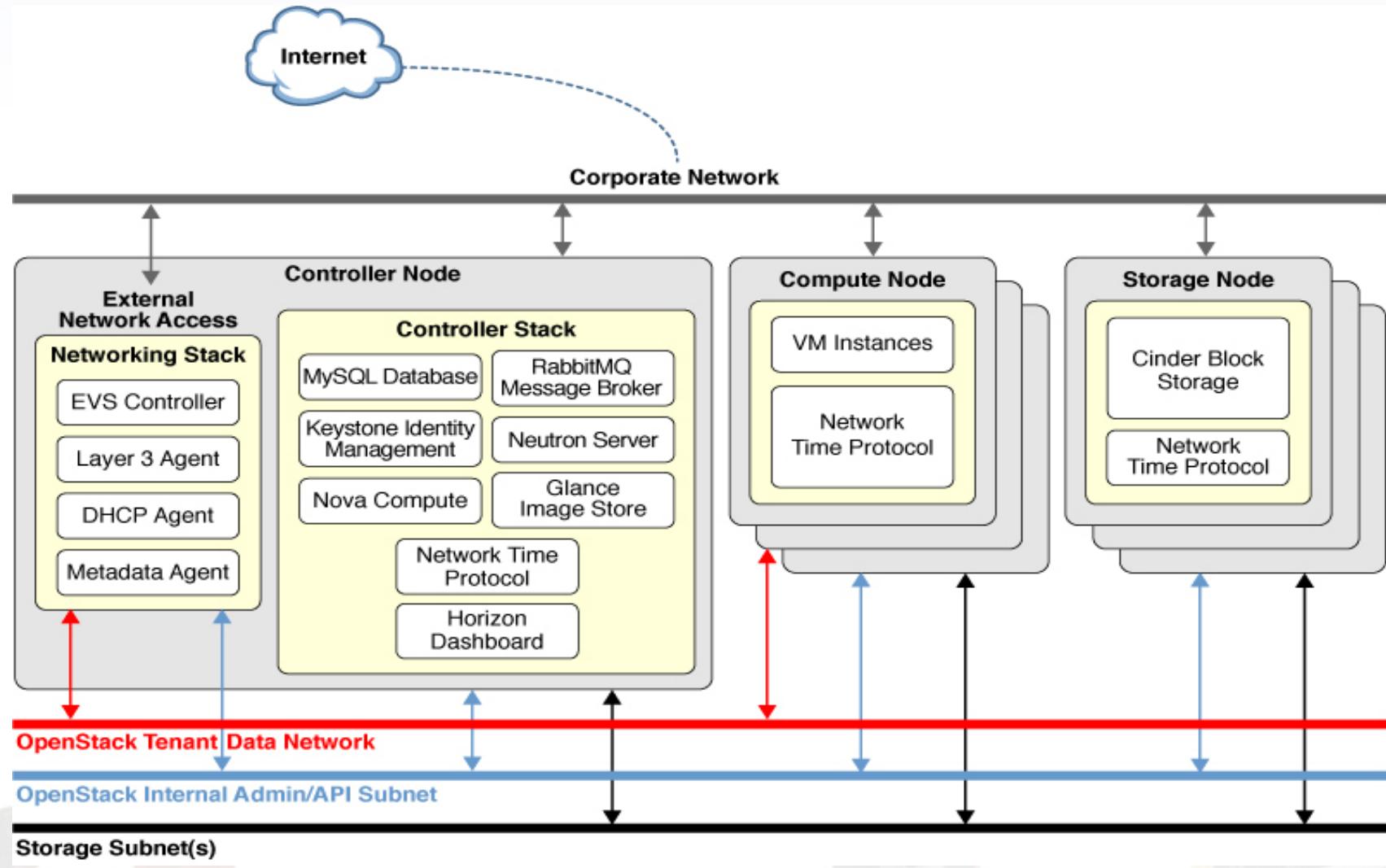
Compute Nodes-SKL(132)



Omni Path High-Speed Interconnect(OPA)



# OpenStack



# Performance Measurement

2 login node and 52 compute nodes, each of them:

2 x IBM Power9 8335-GTH @ 2.4GHz (3.0GHz on turbo, 20 cores and 4 threads/core, total 160 threads per node)

$$- \mathbf{160 \times 2.4 \text{ G FLOPS} = 384 \text{ GFLOPS}}$$

512GB of main memory distributed in 16 dimms x 32GB @ 2666MHz

2 x SSD 1.9TB as local storage and 2 x 3.2TB NVME

4 x GPU NVIDIA V100 (Volta) with 16GB HBM2.

$$- \mathbf{4 \times 14.13 \text{ TFLOPS} = 56.52}$$

GPFS via one fiber link 10 Gbit

Supercomputing System Performance =  $( 56.52 + .384 ) \times 52$

**2.959 PFLOPS**

# Performance Bottleneck

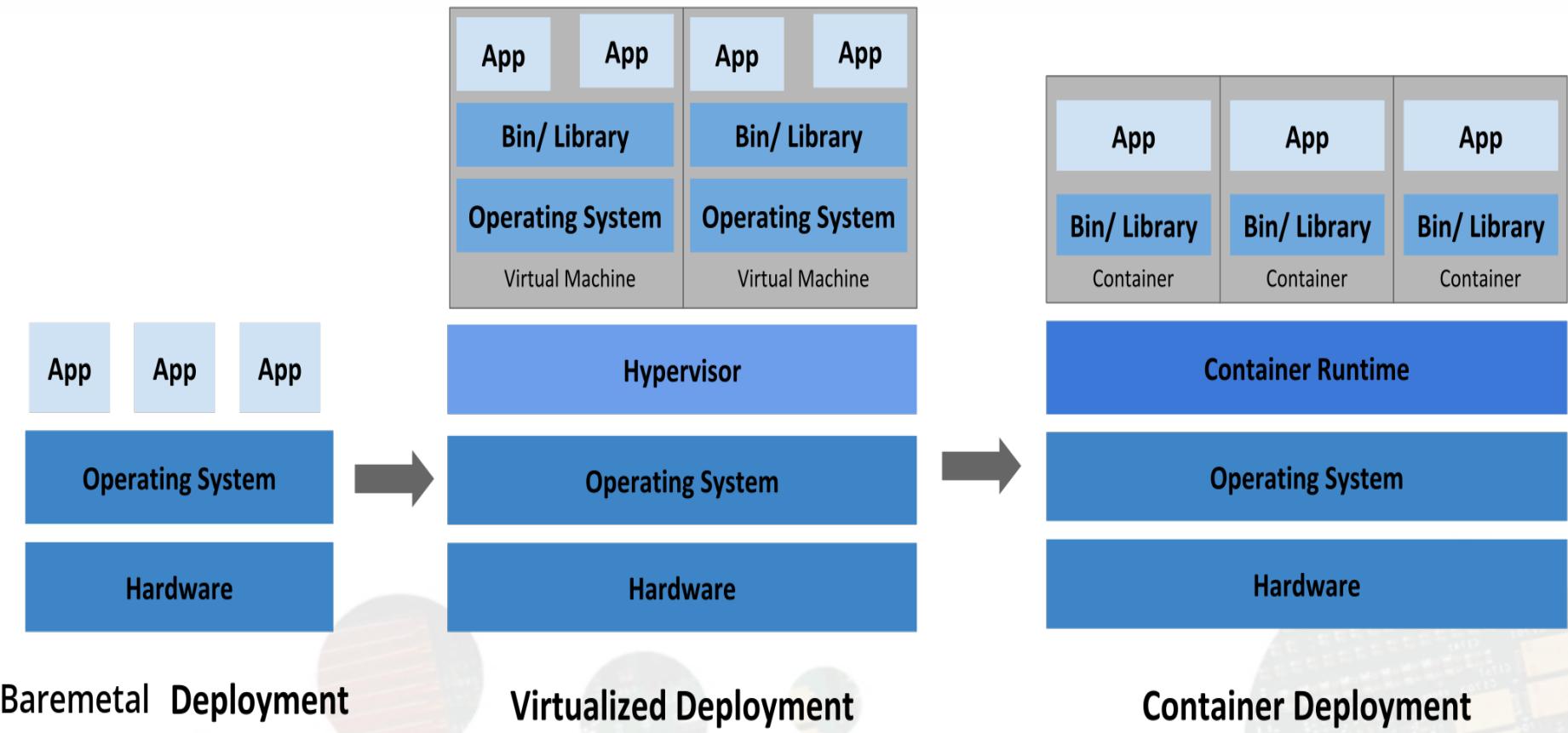
Main Memory Delay

Network (I/O) Delay

Disk Storage Delay

- Importance of HPC
- Types of Processing System
- HPC System Architecture
- Supercomputing Classes and Infrastructure
- **Cluster Software Stacks**
- Programming Models
- Applications

# Software Stack Structures and Usage



# Cluster System Software Stack Configuration Tools

## Linux Distributions

- Red Hat Enterprise Linux (RHEL) and CentOS
- SUSE Linux Enterprise Server (SLES)
- Ubuntu Server

## HPC-Optimized Linux Distributions

- CentOS Stream for HPC
- OpenHPC

## Specialized Linux Variants:

- Rocky Linux
- Scientific Linux

## Environment Modules System

- Singularity containers in an HPC environment

## Cluster Management Tools

- Automate the deployment and management of nodes

# High-Performance Computing (HPC) Bare Metal

Slurm: A popular open-source workload manager and job scheduler for managing and optimizing HPC clusters.

Distributed computation environments are used for software development (OpenMP, MPI). OpenMPI: A widely used message-passing interface for developing parallel and distributed computing applications.

Intel oneAPI: A unified software stack from Intel for developing high-performance applications that can leverage CPU, GPU, and FPGA architectures.

GNU Compiler Collection (GCC): A set of compilers and libraries often used for building high-performance applications in various programming languages.

CUDA Toolkit: NVIDIA's software platform for developing applications that leverage NVIDIA GPUs for high-performance computing.

High-Performance Math Function Library: BLAS (Basic Linear Algebra Subprograms), LINPACK, LAPACK (Linear Algebra PACKage) Linux

# Kubernetes

Docker: A platform for developing, shipping, and running applications using containerization.

Kubernetes: An open-source container orchestration platform for automating the deployment, scaling, and management of containerized applications.

Helm: A package manager for Kubernetes that simplifies the deployment and management of applications.

Prometheus: An open-source monitoring and alerting toolkit designed for cloud-native environments.

Fluentd: An open-source data collector for unifying log management across different services.

# OpenStack Platform

Nova: The compute service in OpenStack, responsible for managing virtual machines and instances.

Neutron: The networking service that provides networking resources and connectivity for OpenStack deployments.

Cinder: The block storage service that enables the creation and management of block storage volumes.

Swift: The object storage service designed to store and retrieve large amounts of unstructured data.

Horizon: The web-based dashboard interface for managing and monitoring OpenStack resources.

Heat: The orchestration service that allows users to define and automate the deployment of resources using templates.

# HPC , AI and BigData Software Stack

Deep Learning and BigData Environment	Frameworks	Caffe, Caffe2, Caffe-MPI, Chainer, Microsoft CNTK, Keras, MXNet, Tensorflow, Theano, PyTorch Apache Hadoop, Apache Spark, Apache Flink, Apache Storm, Apache			
	Libraries	cnDNN, NCCL, cuBLAS, Apache Hive, Apache Pig, Apache HBase, Spark SQL, MLlib, GraphX			
	User Access	<b>NVIDIA DIGITS, Apache Zeppelin</b>			
Programming Environment	Development & Performance Tools	Intel Parallel Studio XS Cluster Edition	PGI Cluster Development Kit	GNU Toolchain	NVIDIA CUDA
	Scientific and Communication Libraries	Intel MPI	MVAPICH2, MVAPICH	IBM Spectrum LSF	Open MPI
	Debuggers	Intel IDB	PGI PGDBG	GNU GDB	
Schedulers, File Systems and Management	Resource Management/Job Scheduling	Adaptive Computative Moab, Maui TORQUE	SLURM	Altair PBS Professional	IBM Spectrum LSF
	File Systems	Lustre	NFS	GPFS	Local (ext3, ext4, XFS)
	Cluster Management	Beowulf, xCat, OpenHPC, Rocks, Bright Cluster Manager for HPC including support for NVIDIA Data Center GPU Manager			
Operating Systems and Drivers	Drivers & Network Mgmt.	Accelerator Software Stack and Drivers		OFED, OPA	
	Operating Systems	Linux (RHEL, CentOS, SUSE Enterprise, Ubuntu, etc.)			

		<b>Supercomputing(SC)</b>	<b>Deep Learning(DL)</b>	<b>Big Data Computing (BDC)</b>
Apps...  Middleware & MGMT  System SW  Hardware	Boundary Interaction Services	<b>In_house, commercial &amp; OSS applications</b> [e.g. Paraview], <b>Remote desktop</b> [e.g. Virtual Network Computing (VNC)], <b>Secure Sockets Layer</b> [e.g. SSL certificates]	<b>Framework_dependent applications</b> [e.g. NLP, voice, image], <b>Web mechanisms</b> [e.g. Google & Amazon Web Services], <b>Secure Sockets Layer</b> [e.g. SSL certificates]	<b>Framework_dependent applications</b> [e.g. 2/3/4-D], <b>Secure Sockets Layer</b> [e.g. SSL certificates]
	Processing Services	<b>Domain Specific frameworks</b> [e.g. PETSc], <b>Batch processing of large tightly coordinated parallel jobs</b> [100s - 10000s of processes communicating frequently with each other]	<b>DNN training &amp; inference frameworks</b> [e.g. Caffe, Tensorflow, Theano, Neon, Torch], <b>DNN numerical libraries</b> [e.g. dense LA]	<b>Machine Learning</b> (traditional) [e.g. Mahout, Scikit-learn, BigDL], <b>Analytics/Statistics</b> [e.g. Python, ROOT, R, Matlab, SAS, SPSS, SciPy], <b>Iterative</b> [e.g. Apache Hama], <b>Interactive</b> [e.g. Dremel, Drill, Tez, Impala, Shark, Presto, BlinkDB, Spark], <b>Batch / Map Reduce</b> [e.g. MapReduce, YARN, Sqoop, Spark], <b>Real-time / streaming</b> [e.g. Flink, YARN, Druid, Pinot, Storm, Samza, Spark]
	Model / Infromation Management Services	<b>Data Storage:</b> Parallel File Systems [e.g. Lustre, GPFS, BeeGFS, PanFS, PVFS], <b>I/O libraries</b> [e.g. HDF5, PnetCDF, ADIOS]	<b>Data Storage</b> [e.g. HDFS, Hbase, Amazon S3, GlusterFS, Cassandra, MongoDB, Hana, Vora]	<b>Serialization</b> [e.g. Avro], <b>Meta data</b> [e.g. HCatalog], <b>Data Ingestion &amp; Integration</b> [e.g. Flume, Sqoop, Apache NiFi, Elastic Logstash, Kafka, Talend, Pentaho], <b>Data Storage</b> [e.g. HDFS, Hbase, Amazon S3, GlusterFS, Cassandra, MongoDB, Hana, Vora], <b>Cluster Mgmt</b> [e.g. YARN, MESO]
	Communication Services	<b>Messaging &amp; Coordination</b> [e.g. MPI/PGAS, direct fabric access], <b>Threading</b> [e.g. OpenMP, task-based models]	<b>Messaging &amp; Coordination</b> [e.g. Machine Learning Scaling library(MLSL)]	<b>Messaging</b> [e.g. Apache Kafka (streaming)]
	Workflow / Task Services	<b>Conventional compiled languages</b> [e.g. C/C++/Fortran], <b>Scripting languages</b> [e.g. Python, Julia]	<b>Scripting languages</b> [e.g. Python]	<b>Workflow &amp; Scheduling</b> [e.g. Oozie], <b>Scripting languages</b> [e.g. Leras, Mocha, Pig, JAQL, Python, Java, Scala]
	System Management & Security Services	<b>Domain numerical libraries</b> [e.g. PETSc, SCALAPACK, BLAS, FFTW,...], <b>Performance &amp; debugging</b> [e.g. DDT, Vampire], <b>Accelerator APIs</b> [e.g. CUDA, OpenCL, OpenACC] <b>Data Protection</b> [e.g. System SSS, OS/PFS file access control] <b>Batch scheduling</b> [e.g. SLURM], <b>Cluster management</b> [e.g. OpenHPC], <b>Container Virtualization</b> [e.g. Docker], <b>Operating System</b> [e.g. Linux OS Variant]	<b>Batching for training</b> [built into DL frameworks], <b>Reduced precision</b> [e.g. interference engines], <b>Load distribution layer</b> [e.g. Round robin/load balancing for interference], <b>Accelerator APIs</b> [e.g. CUDA, OpenCL], <b>Hardware Optimization Libraries</b> [e.g. cuDNN, MKL-DNN, etc.]  <b>Virtualisation</b> [e.g Dockers, Kubernetes, VMware, Xen, KVM, HyperX], <b>Operating System</b> [e.g. Linux (RedHat, Ubuntu, etc.), Windows]	<b>Distributed Coordination</b> [e.g. ZooKeeper, Chubby, Paxos], <b>Provisioning, Managing &amp; Monitoring</b> [e.g. Ambari, Whirr, BigTop, Chukwa], <b>SVM systems</b> [e.g. Google Sofia, libSVM, svm-py,...], <b>Hardware Optimization Libraries</b> [e.g. DAAL, DPDK, MKL, etc.],  <b>Vitualization</b> [e.g. Dockers, Kubernetes, VPware, Xen, KVM, HyperX], <b>Operating System</b> [e.g. Linux (PedHat, Ubuntu, etc.), Windows]
	Infrastructure	<b>Local storage</b> [e.g. Storage & I/O nodes, NAS] <b>Servers</b> [e.g. CPU & Memory [Gen Purpose CPU nodes, GPUs, FPGAs]] <b>Network</b> [e.g. Infiniband & OPA fabrics]	<b>Local storage</b> [e.g. Local storage or NAS/SAN] <b>Services</b> [e.g. CPU & Memory [Gen Purpose CPU + GPU/FPGA, TPU]] <b>Network</b> [e.g. Ethernet]	<b>Local Storage</b> [e.g. Direct attached Storage] <b>Services</b> [e.g. CPU & Memory, [Gen Purpose CPU hyper-convergent nodes]] <b>Netwrok</b> [e.g. Ethernet fabrics]

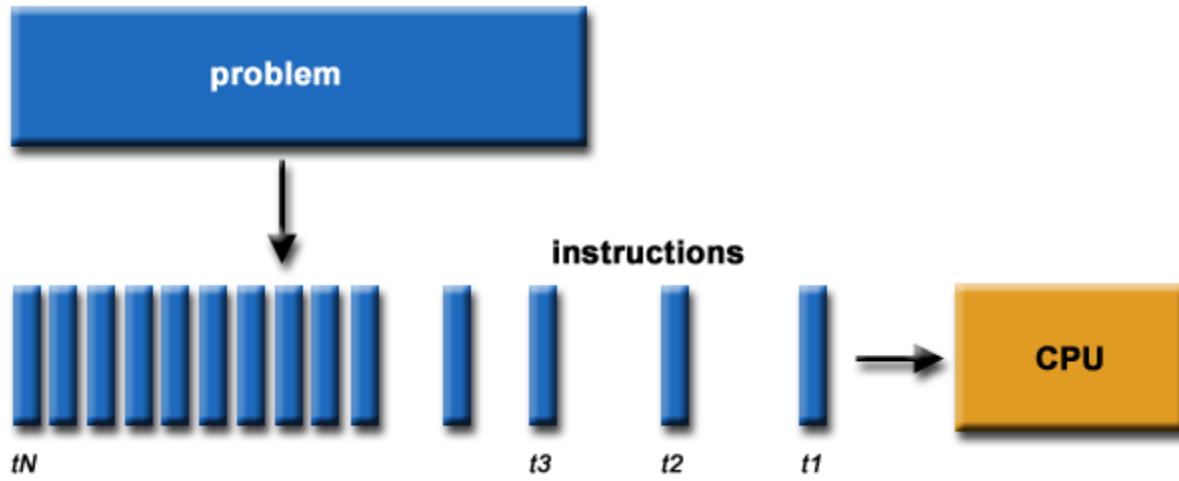
# BigData Framework: Infiniband Support

- Big data Hadoop (inherent support for RDMA);
- Oracle, IBM DB2, Microsoft SQL server database;
- Block storage (support RDMA storage protocol iSER and SRP);
- Parallel file system (GPFS, Lustre, Ceph, etc.); and
- Enterprise applications, such as finance, futures, and animations.
- High-Performance Computing (HPC) Clusters:
- Data Analytics and Big Data Platforms:
- AI and ML

- Importance of HPC
- Types of Processing System
- HPC System Architecture
- Supercomputing Classes and Infrastructure
- Cluster Software Stacks
- **Parallel Programming Models**
- Demonstration (Supercomputing System)

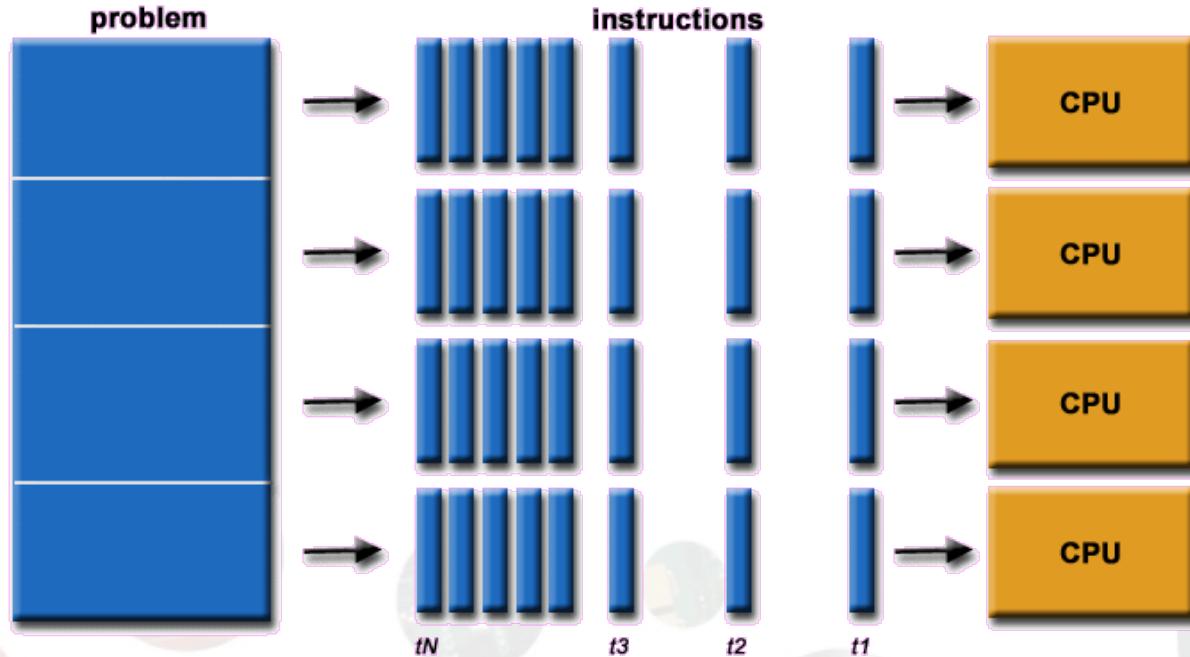
# Parallel Programming ?

- Traditionally, software has been written for ***serial*** computation:
  - To be run on a single computer having a single Central Processing Unit (CPU);
  - A problem is broken into a discrete series of instructions.
  - Instructions are executed one after another.
  - Only one instruction may execute at any moment in time.



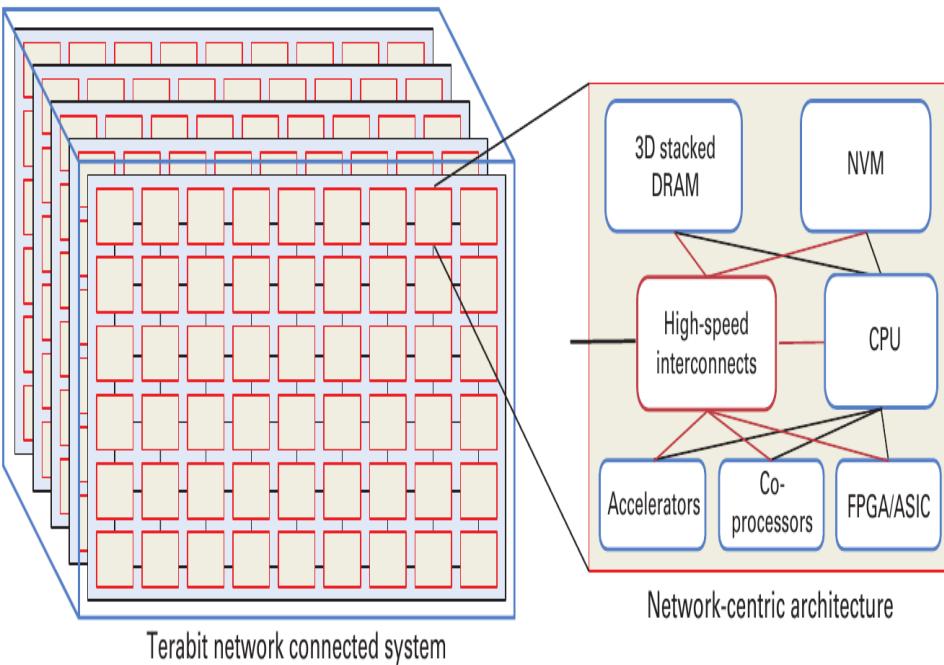
# Parallel Computing

- In the simplest sense, ***parallel computing*** is the simultaneous use of multiple compute resources to solve a computational problem.
  - To be run using multiple CPUs
  - A problem is broken into discrete parts that can be solved concurrently
  - Each part is further broken down to a series of instructions
- Instructions from each part execute simultaneously on different CPUs



# Parallel Computing: Resources

- The compute resources can include:
  - A single computer with multiple processors;
  - A single computer with (multiple) processor(s) and some specialized computer resources (GPU, FPGA ...)
  - An arbitrary number of computers connected by a network;
  - A combination of both.



192 **Core** : single-precision cores

64 **DP Unit** : double -precision cores

32 **LD/ST** : load/store units

32 **SFU** : Special Function Units

# Parallel Computing: The Computational Problem

- The computational problem usually demonstrates characteristics such as the ability to be:
  - Broken apart into discrete pieces of work that can be solved simultaneously;
  - Execute multiple program instructions at any moment in time;
  - Solved in less time with multiple compute resources than with a single compute resource.
  - Amdahl's Law
    - Speedup =  $1 / [(1 - P) + (P / N)]$

# Parallel Computing: what for? (1)

- Traditionally, parallel computing has been considered to be "the high end of computing" and has been motivated by numerical simulations of complex systems and "Grand Challenge Problems" such as:
  - weather and climate
  - chemical and nuclear reactions
  - biological, human genome
  - geological, seismic activity
  - mechanical devices - from prosthetics to spacecraft
  - electronic circuits
  - manufacturing processes

# The future

- During the past 10 years, the trends indicated by ever faster networks, distributed systems, and multi-processor computer architectures (even at the desktop level) clearly show that ***parallelism is the future of computing.***
- It will be multi-forms, mixing general purpose solutions (PC) and very specialized solutions as IBM Cells, ClearSpeed, GPGPU from Nvidia etc.

# How Parallel Processing?

## Instruction Level Parallelism

- Pipelining and Superscalar Execution

## Data Level Parallelism

- Vector Instruction or Specilizaed Accelerator

## Thread Level Parallelism

- Execute multiple function on different cores

## Task Level Parallelism

- Breaking multiple tasks (Memory, Input Output etc) into subtasks and execute using TLP

# Writing Parallel Application

A=10;

B=20;

C=A+B;

D=A\*B;

E=C\*A;

F=C\*B;

G=D\*A;

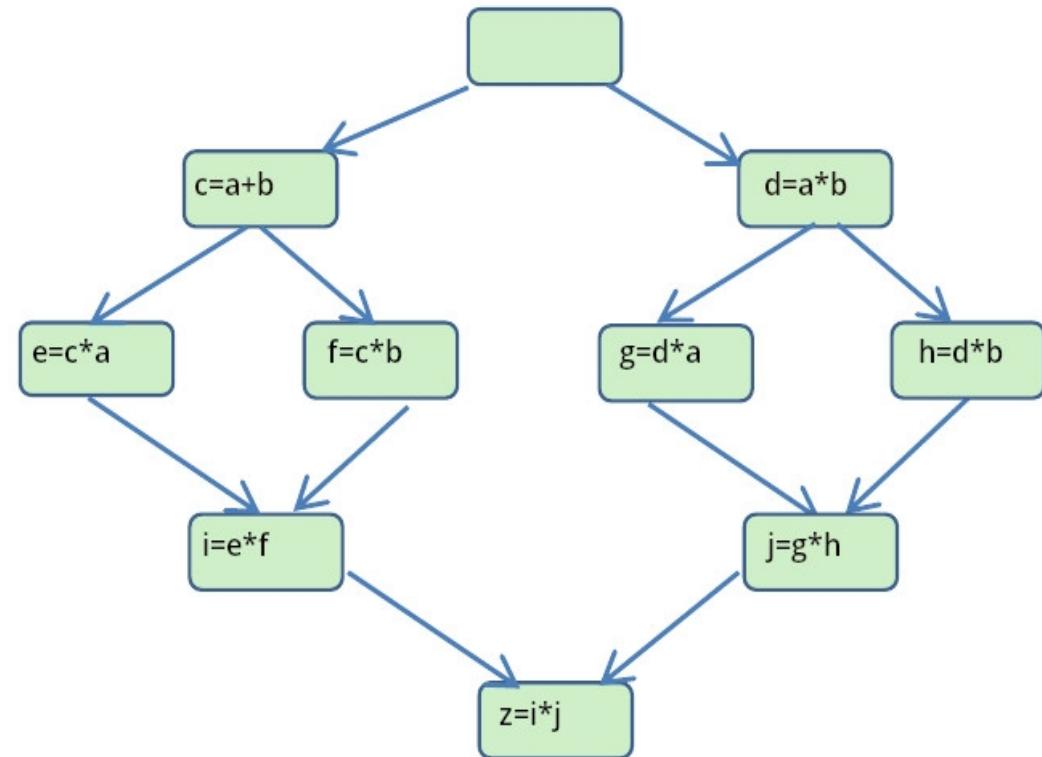
H=D\*B;

I=E\*F;

J=G\*H;

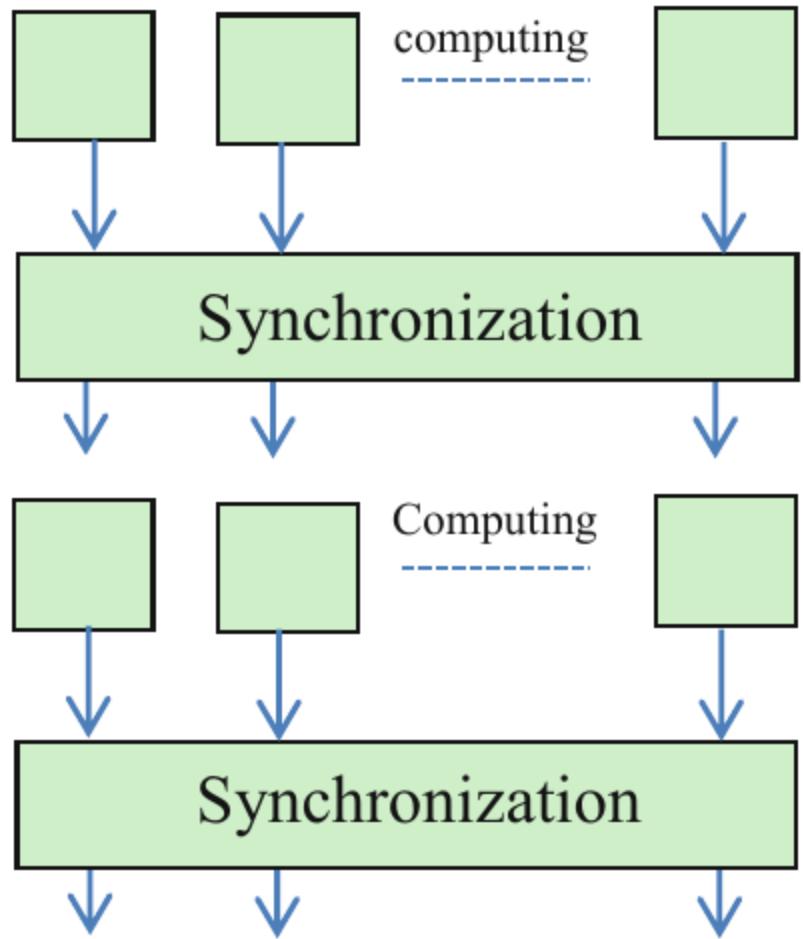
Z=I\*J;

Development and Application of Supercomputing

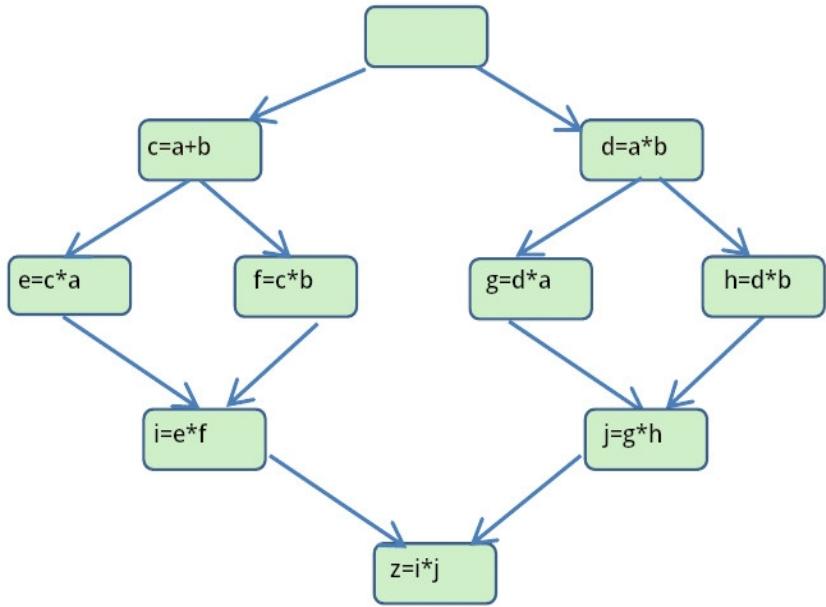


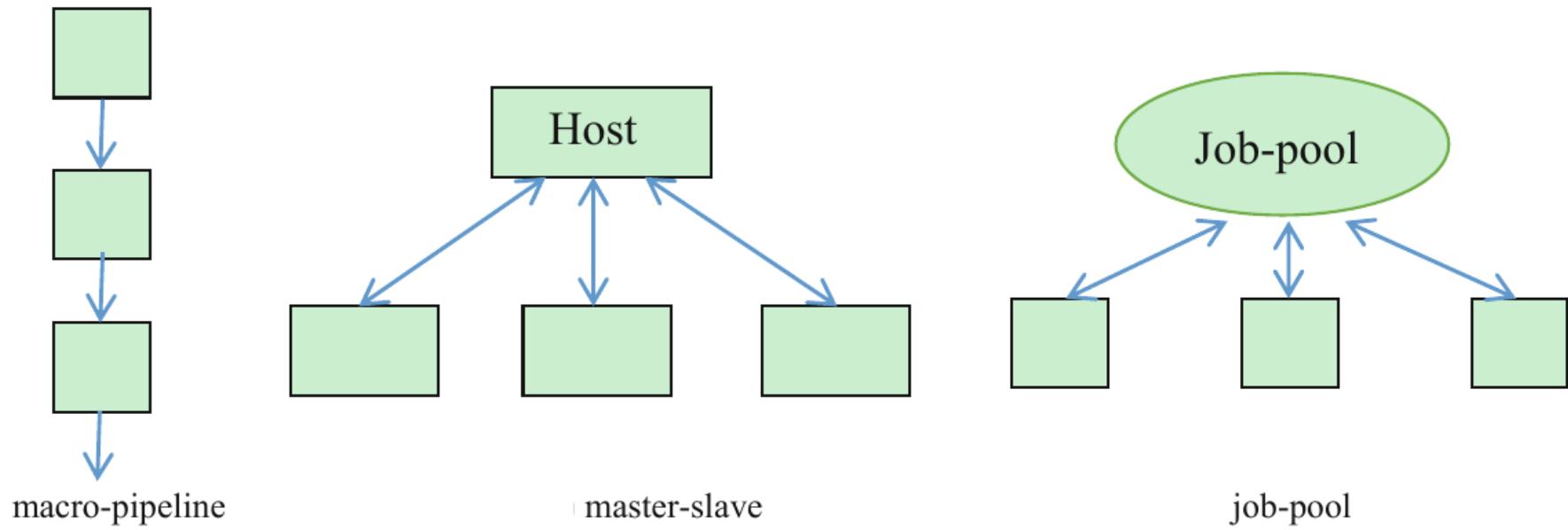
# Considerations for Parallel Programs

- Automatic vs. Manual Parallelization
- Understand the Problem and the Program
- Data Dependencies
- Partitioning
- Communications
- Synchronization
- Load Balancing
- Granularity
- I/O
- Limits and Costs of Parallel Programming
- Performance Analysis and Tuning



Development and Application of Supercomputing





# Frameworks for Parallel Programming

**Scientific Computing**

**Big Data Processing**

**AI and Machine Learning**

# Parallel Programming Models

Shared Memory: Multi-threading programming languages: such as POSIX Threads, Java Threads, OpenMP.

- OpenMP: Parallelize loops and sections of code using compiler directives.

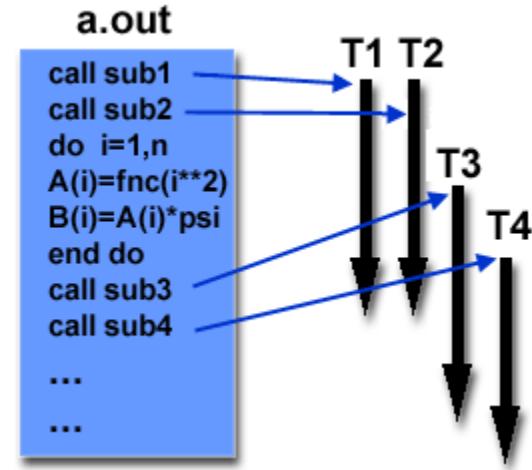
MPI (Message Passing Interface): A widely used standard for writing parallel programs that execute on distributed memory systems. It's commonly used in high-performance computing for scientific simulations and calculations.

OpenACC: OpenACC is a directive-based approach to parallel programming that focuses on using high-level directives to guide the compiler in parallelizing code for accelerators like GPUs.

CUDA: NVIDIA's parallel computing platform and programming model that enables developers to use GPUs for accelerating scientific computations.

BLAS (Basic Linear Algebra Subprograms)

# Threads Model



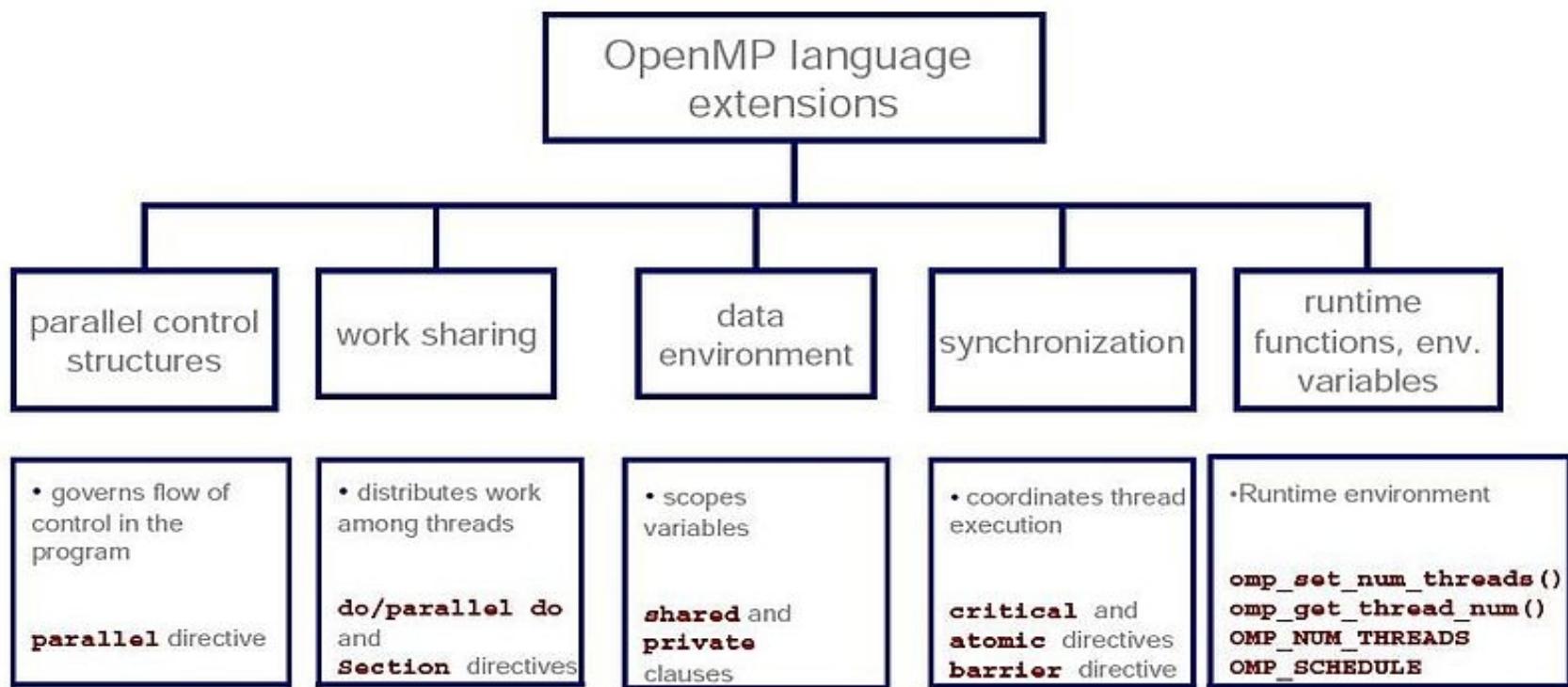
- In the threads model of parallel programming, a single process can have multiple, concurrent execution paths.
- Perhaps the most simple analogy that can be used to describe threads is the concept of a single program that includes a number of subroutines.
- Threads are commonly associated with shared memory architectures and operating systems.

# Shared Memory Programming Model: OpenMP

- **OpenMP**
  - Compiler directive based; can use serial code
  - Jointly defined and endorsed by a group of major computer hardware and software vendors. The OpenMP Fortran API was released October 28, 1997. The C/C++ API was released in late 1998.
  - Portable / multi-platform, including Unix and Windows NT platforms
  - Available in C/C++ and Fortran implementations
  - Can be very easy and simple to use - provides for "incremental parallelism"
- Microsoft has its own implementation for threads, which is not related to the UNIX POSIX standard or OpenMP.

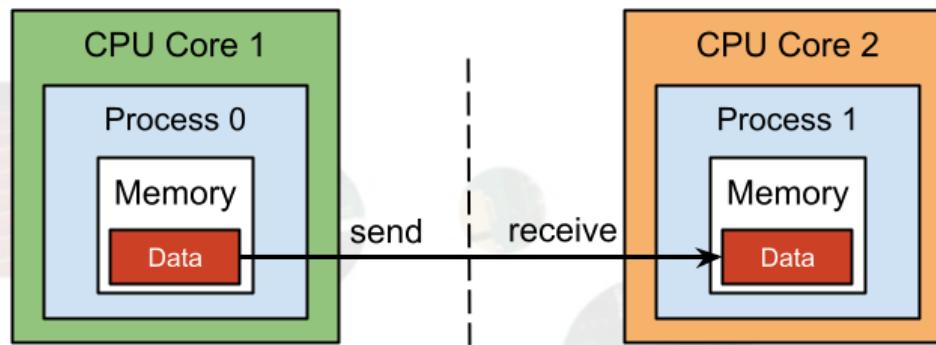
# Example:

## OpenMP Constructs



# Message Passing Model

- The message passing model demonstrates the following characteristics:
  - A set of tasks that use their own local memory during computation. Multiple tasks can reside on the same physical machine as well across an arbitrary number of machines.
  - Tasks exchange data through communications by sending and receiving messages.
  - Data transfer usually requires cooperative operations to be performed by each process. For example, a send operation must have a matching receive operation.



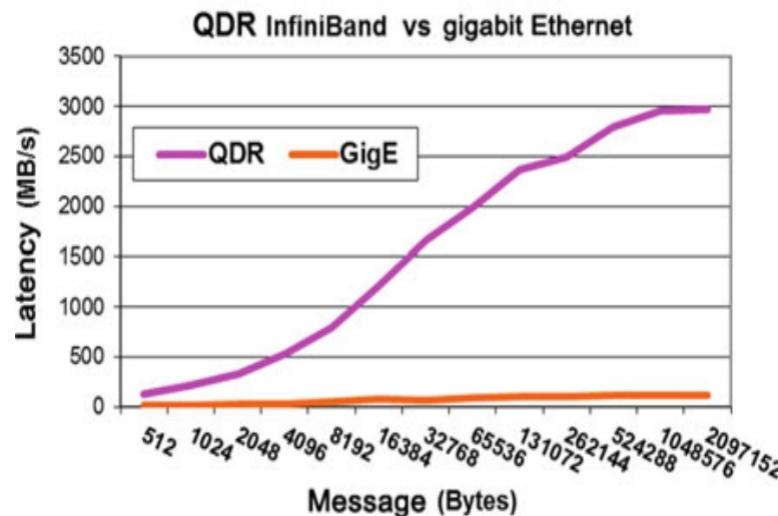
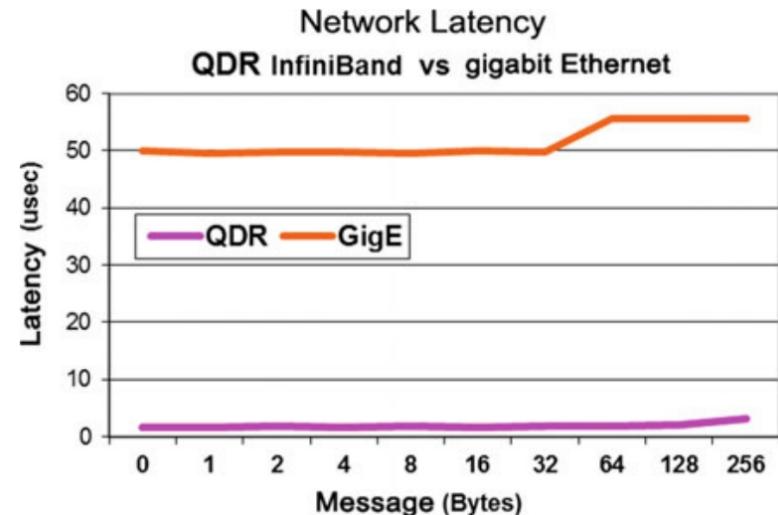
# Point-to-Point Communication

MPI defines more than 35 point-to-point communication methods, including

`MPI_Send`, `MPI_Recv`, `MPI_Sendrecv`,  
`MPI_Isend`, `MPI_Irecv`, `MPI_Probe`,  
`MPI_Iprobe`, `MPI_Test`, `MPI_Testall`,  
`MPI_Wait`, and `MPI_Waitall`.

## Collective Communication

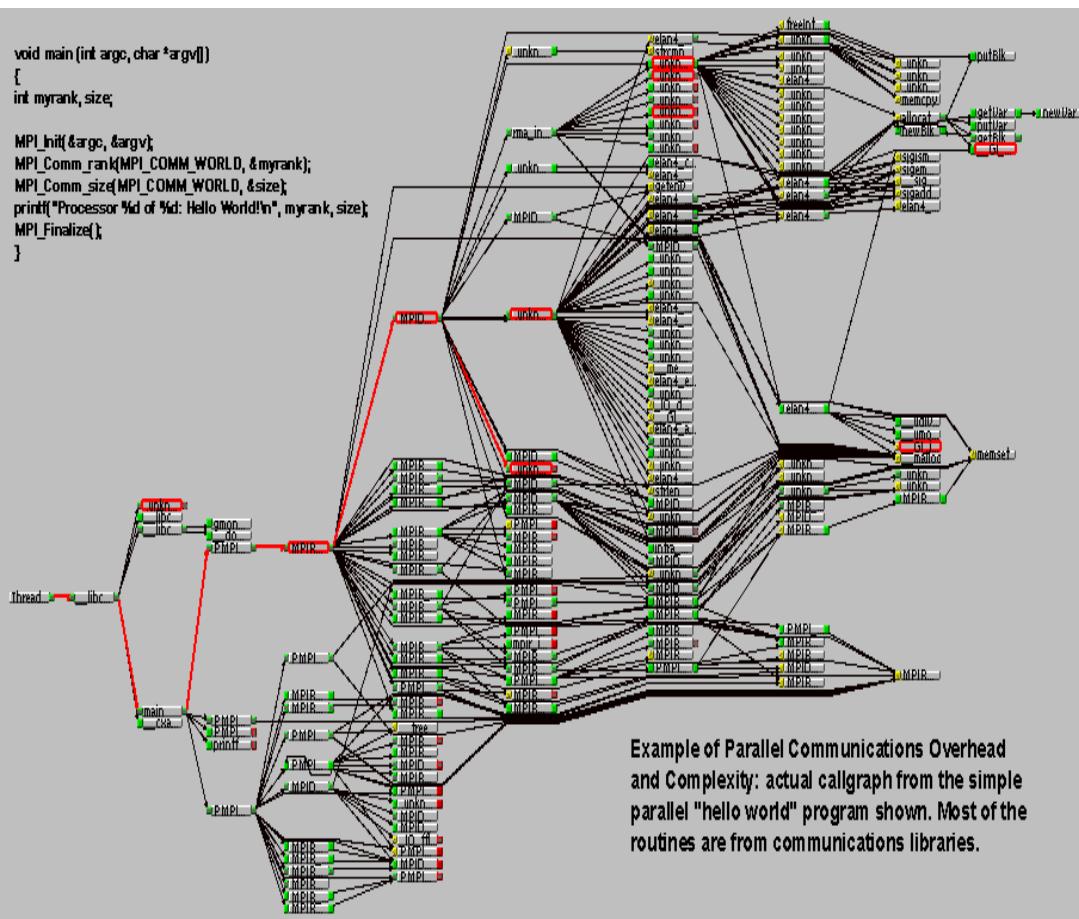
These include: `MPI_Allgather`,  
`MPI_Allgatherv`, `MPI_Allreduce`,  
`MPI_Alltoall`, `MPI_Alltoallv`,  
`MPI_Barrier`, `MPI_Bcast`, `MPI_Gather`,  
`MPI_Gatherv`, `MPI_Reduce`,  
`MPI_Scatter`, and `MPI_Scatterv`.



```

void main(int argc, char *argv[])
{
    int myrank, size;
    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);
    printf("Processor %d of %d: Hello World!\n", myrank, size);
    MPI_Finalize();
}

```



**MPI include file**

*Declarations, prototypes, etc.*

**Program Begins**

*Serial code*

**Initialize MPI environment**

*Parallel code begins*

**Do work & make message passing calls**

**Terminate MPI environment**

*Parallel code ends*

*Serial code*

**Program Ends**

# Hybryd

- In this model, any two or more parallel programming models are combined.
- Currently, a common example of a hybrid model is the combination of the message passing model (MPI) with either the threads model (POSIX threads) or the shared memory model (OpenMP). This hybrid model lends itself well to the increasingly common hardware environment of networked SMP machines.
- Another common example of a hybrid model is combining data parallel with message passing. As mentioned in the data parallel model section previously, data parallel implementations (F90, HPF) on distributed memory architectures actually use message passing to transmit data between tasks, transparently to the programmer.

# Big Data Processing

Apache Hadoop: An open-source framework for distributed storage and processing of large datasets across clusters of computers using the MapReduce programming model.

Apache Spark: A fast and general-purpose cluster computing system that provides in-memory data processing for big data analytics.

Apache Flink: A stream processing framework for real-time event-driven data processing with support for batch processing as well.

Apache Beam: A unified programming model for both batch and stream data processing, supporting multiple backends like Apache Spark, Google Cloud Dataflow, and more.

# AI and Machine Learning:

TensorFlow: An open-source deep learning framework developed by Google that supports both CPU and GPU acceleration for training neural networks.

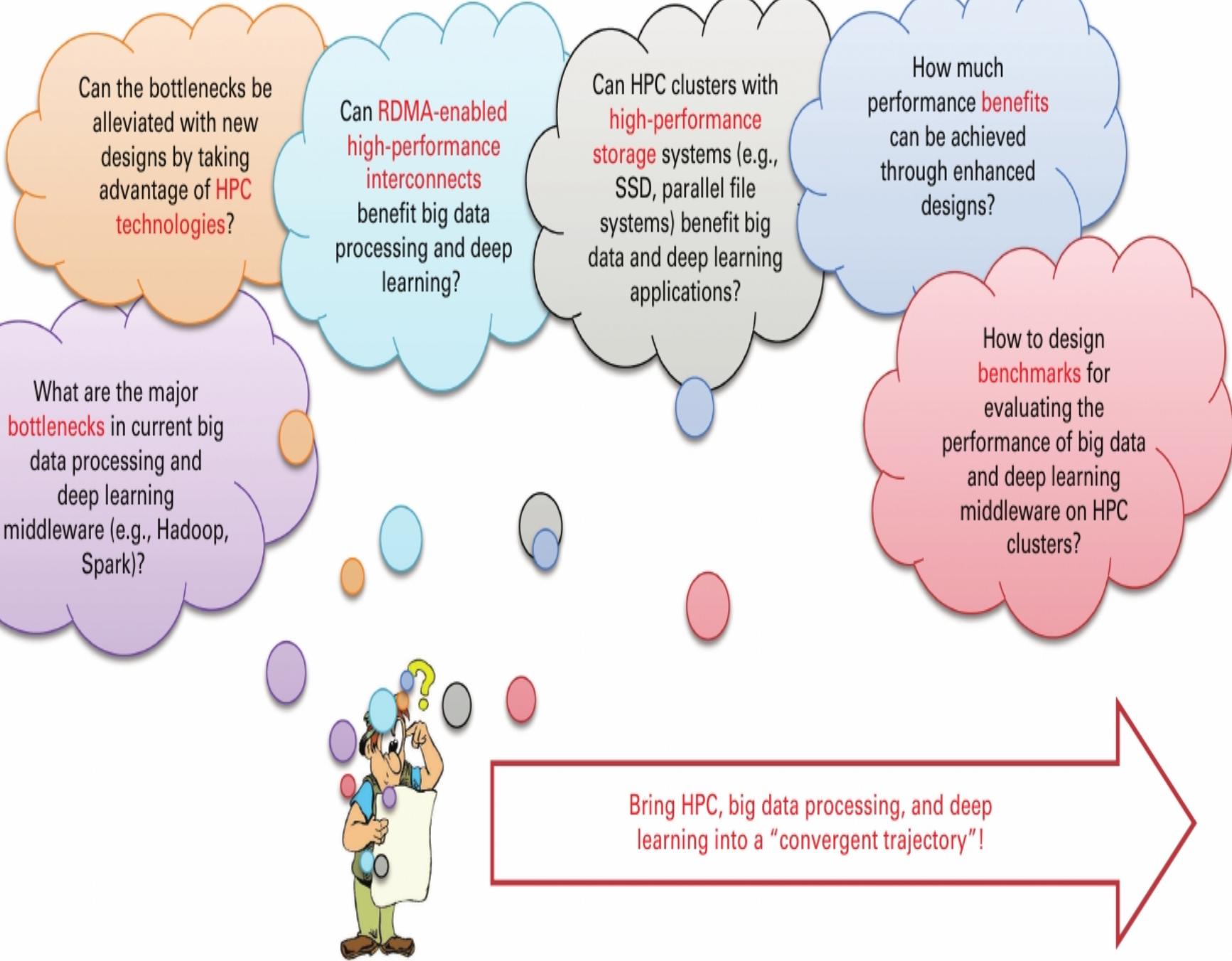
PyTorch: An open-source machine learning framework developed by Facebook's AI Research lab, known for its dynamic computation graph and ease of use.

Apache MXNet: A flexible and efficient deep learning framework that supports both symbolic and imperative programming paradigms.

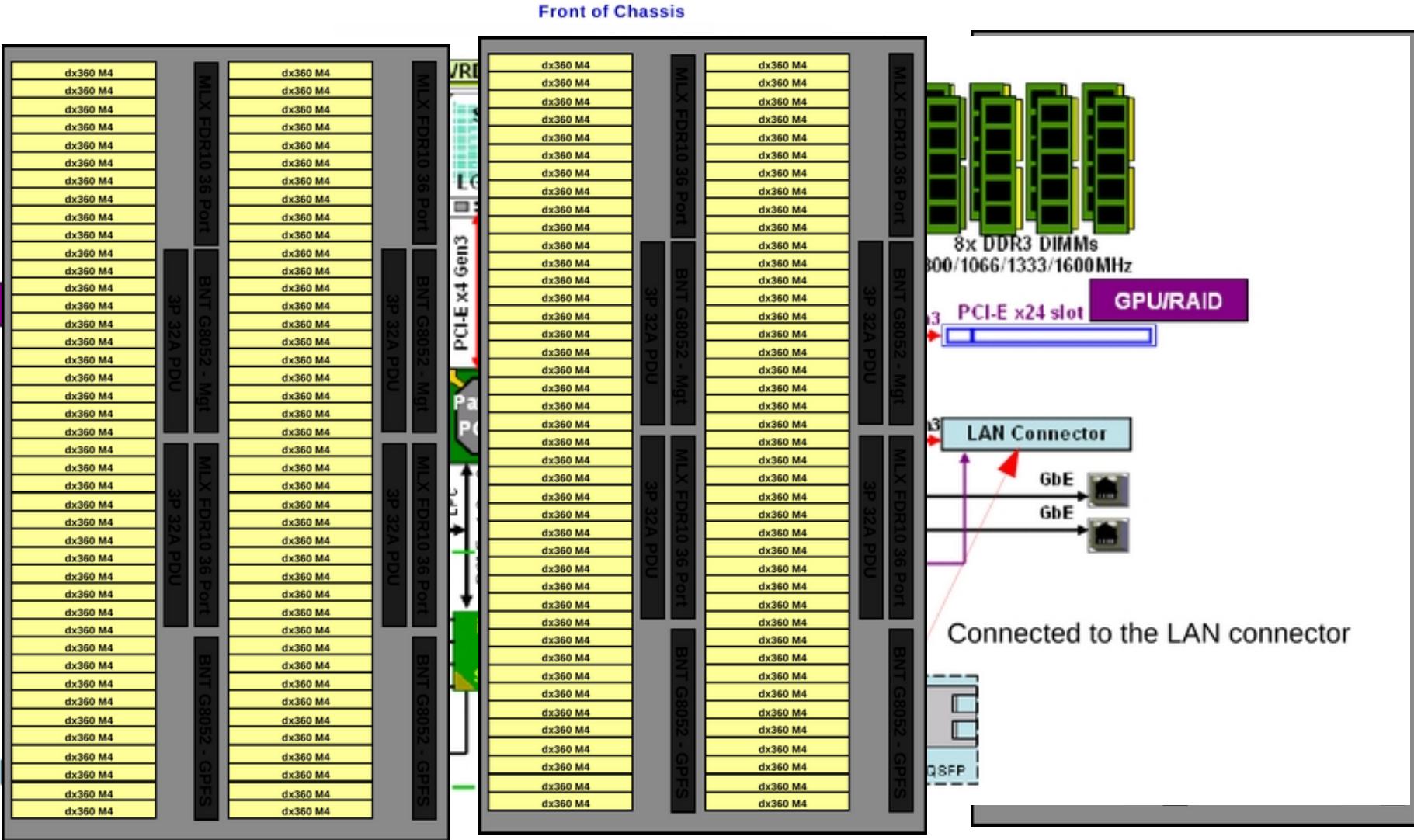
Horovod: A distributed deep learning framework designed to scale TensorFlow and PyTorch training across multiple GPUs and nodes.

Distributed TensorFlow: An extension of TensorFlow that enables distributed training of machine learning models across multiple devices and machines.

- Importance of HPC
- Types of Processing System
- HPC System Architecture
- Supercomputing Classes and Infrastructure
- Cluster Software Stacks
- Programming Models
- **Demonstration (Supercomputing System)**



# Bottom-up Design Approach



# Applications Services

Data Sciences

Health Science

Social Sciences

Agriculture

High Performance Computing

Web  
(IoT, VLSI Design)

# Development Frameworks and Libraries

Interactive

GCC

Python

OpenMP

MPI

CUDA

OpenACC

OpenCL

TensorFlow

Horovod

Hadoop

PowerAI

DeepSpeed

Spark

# Distributed System & Software Stack

OpenHPC, ROCKS

OpenShift, xCAT  
Nutanix Acropolis

Open-Stack  
Kubernetes

Linux Kernel: OpenPBS, PBS-Pro, SLURM, Ganglia , Open vSwitch, warewolf, Lustre, BeeGFS, Ceph, Mellanox OFED, IPoIB, OpenEth, Network Information Service, ACPI

Rolls, Singularity Image, Docker, Contrainer

# Hardware System

Intelligent RACK infrastructure  
PDU, PMS

Accelerators  
GPU/TPU/FPGAs

Multi-core  
CISC/SuperScalar

SAN/NAS,  
SSDs/NVMe

High-Speed Ethernet,  
Infiniband

# Challenges

How to Write Highly Scalable and Portable Parallel Programs?

How to Enable Automated Use of Effective Parallel Programming Techniques When Writing Parallel Programs?

How to Enforce the Use of Parallel Programming Design Principles During the Programming Process?

How to Employ Suitable Optimization Techniques?

How to Promote Interdisciplinary Collaboration?



## Barcelona Supercomputing Center

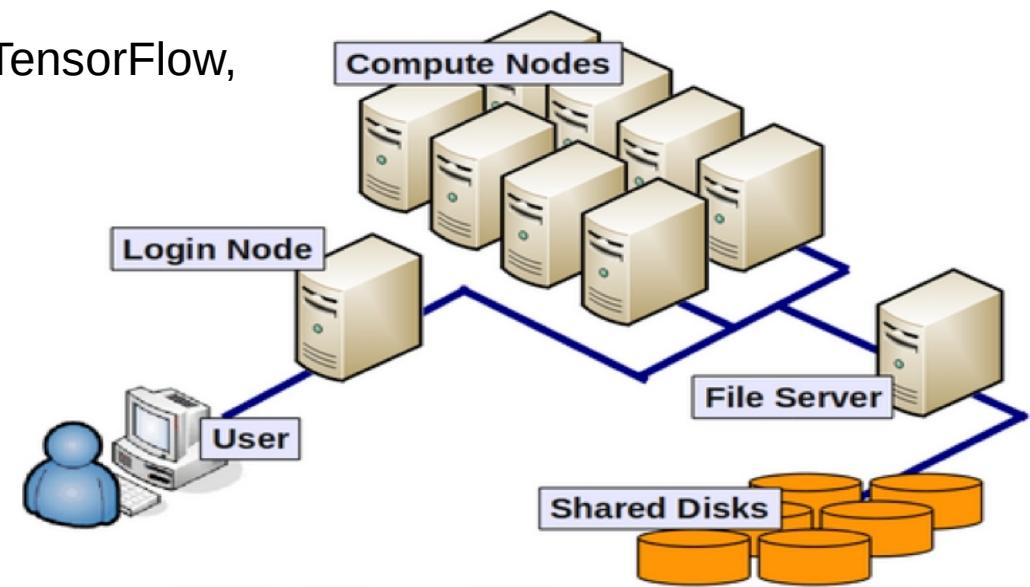
Linux Operating system,  
54 Servers

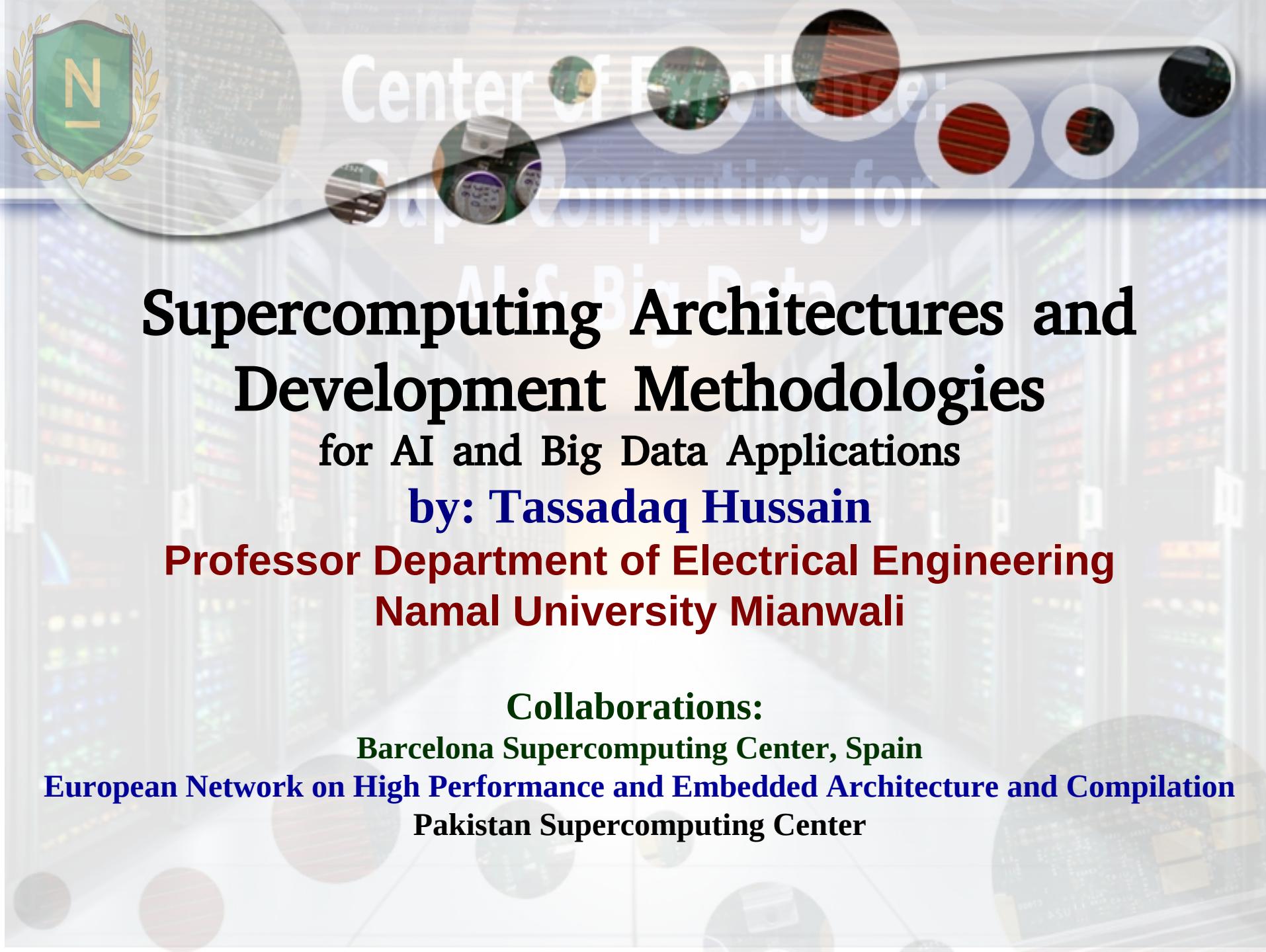
Two Power 9 processors, 4 NVIDIA GPUs with 512 GB of  
main memory

Interconnection InfiniBand is an industry-standard to  
interconnect servers that allows the local memory of one  
server to be accessed from remote servers speedily.

The software stack

**Horovod**, from Uber. Horovod Plugs into TensorFlow,  
PyTorch, and MXNet.





# **Supercomputing Architectures and Development Methodologies**

**for AI and Big Data Applications**

**by: Tassadaq Hussain**

**Professor Department of Electrical Engineering  
Namal University Mianwali**

## **Collaborations:**

**Barcelona Supercomputing Center, Spain**

**European Network on High Performance and Embedded Architecture and Compilation**

**Pakistan Supercomputing Center**