UniTs - University of Trieste

Faculty of Scientific and Data Intensive Computing

Department of mathematics informatics and geosciences

# High Performance Computing

*Lecturer:*
**Prof. Stefano Cozzini**

*Authors:*

**Andrea Spinelli**
**Christian Faccio**

March 10, 2025

 github.com/Spina02,christianfaccio  ✉

andreaspinelli2002@gmail.com,christianfaccio@outlook.it

# Abstract

As a student of Scientific and Data Intensive Computing, I've created these notes while attending the **High Performance Computing** module of **High Performance and Cloude Computing** course. The course will introduce the fundamentals of High Performance Computing, exploring both its concepts and practical applications. The notes cover a wide range of topics, including:

- An overview of High Performance Computing and its importance in solving complex, real-world problems.

- The principles behind modern computer architectures and how they influence performance.

- Essential tools and techniques for parallel programming, alongside strategies to optimize code for advanced architectures.

- The evolution of computing facilities and how to effectively leverage them for large-scale computational challenges.

- Developing a proactive mindset, moving beyond the use of pre-packaged tools to a deeper understanding of the underlying systems.

This comprehensive approach not only equips you with technical skills but also fosters a critical perspective on technological innovation in computing.

While these notes were primarily created for my personal study, they may serve as a valuable resource for fellow students and professionals interested in High Performance Computing.

# Contents

# 1 Introduction

## 1.1 Base Concepts

High Performance Computing (HPC), also known as supercomputing, refers to computing systems with extremely high computatiional power that are able to solve hugely complex and demanding problems. [**europaHighPerformance**]

Often, high precision and accuracy are required in scientific and engineering simulations, which can be achieved by increasing the computational power of the system. This is where HPC comes into play, as it allows for the execution of large-scale **simulations** of complex problems in a reasonable amount of time. Simulations have become the key method for researching and developing innovative solutions in both scientific and engineering fields. They are especially prominent in leading domains such as the aerospace industry and astrophysics, where they enable the investigation and resolution of highly complex problems. However, the increasing reliance on simulation also introduces significant challenges related to complexity, scalability, and data management, which in turn impact the supporting IT infrastructure.

As scientific inquiry progresses along what is known as the Inference Spiral of System Science, the complexity of models intensifies and the influx of new data enriches these systems with additional insights. Consequently, this dynamic evolution necessitates ever increasing computational power to efficiently handle the enhanced simulations and data management challenges.
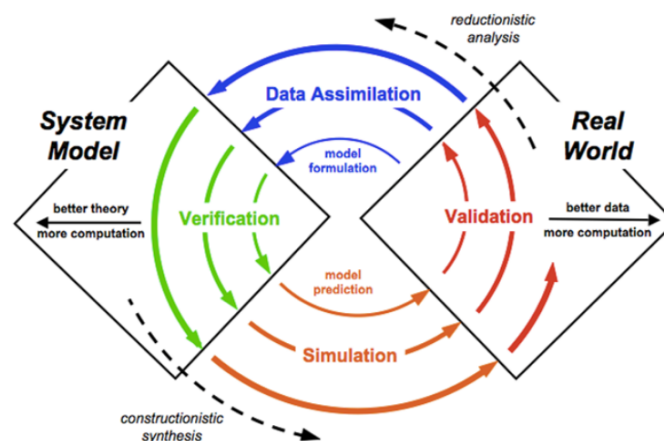


Figure 1.1: Research and Development

> 👁 **Observation**:
>
> In today's world, larger and larger amounts of data are constantly being generated, from 33 zettabytes globally in 2018 to an expected 181 zettabytes in 2025. This exponential growth is driving a shift towards data-intensive applications, making HPC indispensable for processing and analyzing these vast datasets efficiently. Consequently, HPC is key to unlocking valuable insights that benefit citizens, businesses, researchers, and public

### 1.1.1   What is High Performance Computing?

High Performance Computing (HPC) involves using powerful servers, clusters, and supercomputers, along with specialized software, tools, components, storage, and services, to solve computationally intensive scientific, engineering, or analytical tasks.

HPC is used by scientists and engineers both in research and in production across industry, government and academia.

Key elements of the HPC ecosystem include:

- **Hardware:** High-performance servers, clusters, and supercomputers.
- **Software:** Specialized tools and applications designed to optimize complex computations.
- **Applications:** Scientific, engineering, and analytical tasks that leverage high computational power.

**People in HPC**

Human capital is by far the most important aspect in the HPC landscape. Two crucial roles include HPC providers, who plan, install, and manage the resources, and HPC users, who leverage these resources to their fullest potential. The mixing and interplaying of these roles not only enhances individual competence but also drives overall advancements in high-performance computing.

### 1.1.2   Performance and metrics

**Performance** in the realm of high-performance computing is a multifaceted concept that extends far beyond a mere measure of speed. While terms such as "how fast" something operates are often used to describe performance, they tend to be vague. Many factors contribute to the overall performance of a system, and the interpretation of these factors can vary depending on the specific context and objectives of the computational task. Performance, therefore, remains a complex and central issue in the field of HPC, as it involves more than just the raw computational speed.

The discussion often extends to the idea that the "P" in HPC might stand for more than just performance. A growing sentiment among professionals in the field suggests that high performance should be complemented by high productivity. This broader view recognizes that the true efficiency of a computing system is not only determined by its ability to perform tasks quickly but also by the ease and speed with which applications can be developed and maintained. In other words, while raw performance is critical, the overall productivity of a system—combining the system's speed with the programmer's effort—plays an equally important role.

To further clarify the distinction, consider that performance can be seen as a measure of how effectively a system executes tasks, whereas productivity is the outcome achieved relative to the effort invested in developing the application. For instance, if a code optimization leads to a system that runs twice as fast but requires an extensive period of development—say, six months of work—the benefits of the improvement must be weighed against the increased effort required. This example underlines the importance of balancing performance gains with the associated development costs.

Ultimately, the challenge lies in understanding and optimizing both aspects. A successful HPC system is one that not only achieves high computational throughput but also enhances the productivity of the developers who create and refine the applications. This balance is essential for advancing the capabilities of high-performance computing in both research and production environments.

**Number Crunching on CPU**

When evaluating the performance of a high-performance computing (HPC) system, one of the most fundamental metrics is the rate at which floating point operations are executed. This rate is typically expressed in millions (Mflops) or billions (Gflops) of operations per second. In essence, it quantifies how many calculations, such as additions and multiplications, the system is capable of performing every second.

To estimate this capability, we rely on the concept of theoretical peak performance. This value is computed by considering the system's clock rate, the number of floating point operations that can be executed in a single clock cycle, and the total number of processing cores available. Under ideal conditions, the theoretical peak performance can be expressed as follows:

$$\text{FLOPS} = \text{clock\_rate} \times \text{Number\_of\_FP\_operations} \times \text{Number\_of\_cores}$$

This formula provides an upper bound on the computational power of the system. However, it is important to note that this is a best-case scenario estimate and does not always reflect the performance achievable in real applications.

**Sustained (Peak) Performance**

While the theoretical peak performance offers insight into the maximum potential of an HPC system, the actual performance observed during real-world operations is better captured by the sustained (or peak) performance. In practice, several factors such as memory bandwidth limitations, communication latencies, and input/output overhead can prevent a system from reaching its theoretical maximum.

Sustained performance refers to the effective throughput that an HPC system attains when executing actual workloads. Since it is challenging to exactly measure the number of floating point operations performed by every application, standardized benchmarks are commonly used to assess this performance. One widely recognized benchmark is the HPL Linpack test, which forms the basis for the TOP500 list of supercomputers. This benchmark emphasizes the importance of sustained performance, as it reflects the system's efficiency and reliability under realistic operational conditions.

Understanding both the theoretical and sustained performance metrics is crucial. While the former provides an idealized estimate of a system's capabilities, the latter offers a more practical perspective, thereby guiding decisions on system improvements and resource allocation in high-performance computing environments.

## 1.1.3 Supercomputers and TOP500

Supercomputers are the most powerful and advanced computing systems available today. They are designed to handle the most complex and demanding computational tasks, such as weather forecasting, climate modeling, and nuclear simulations. Supercomputers are typically used in scientific research, engineering, and other fields that require massive computational power. The TOP500 list (www.top500.org) is a ranking of the world's most powerful supercomputers. It is published twice a year and provides a snapshot of the current state of high-performance computing. The list ranks supercomputers based on their performance on the Linpack benchmark, which measures the speed at which a system can solve a dense system of linear equations. The TOP500 list is widely regarded as the most authoritative ranking of supercomputers and is used by researchers, industry professionals, and policymakers to track trends in high-performance computing.

The **HPL Linpack benchmark** measures how fast a system can solve a large dense system of linear equations. It estimates floating-point performance by stressing both CPU and memory

| Performance | |
| --- | --- |
| Linpack Performance (Rmax) | 1,742.00 PFlop/s |
| Theoretical Peak (Rpeak) | 2,746.38 PFlop/s |
| Nmax | 25,446,528 |
| **Power Consumption** | |
| Power: | 29,580.98 kW |
| Power Measurement Level: | 2 |
| **Software** | |
| Operating System: | TOSS |
| Compiler: | g++ 12.2.1 and hipcc 6.2.0 |
| Math Library: | AMD rocBLAS 6.0.2 and Intel MKL 2016 |
| MPI: | HPE Cray MPI |

Figure 1.2: Best Supercomputers in the World at the moment (March 2025)

resources, making it a standard method for comparing supercomputers. HPL results are reported in floating-point operations per second (FLOPS), reflecting sustained system throughput under a realistic workload. This test underpins the widely recognized TOP500 list, where higher HPL scores signify more capable machines in real-world, data-intensive scenarios.

For each machine the following numbers are reported using HPL:

- **Rmax:** The maximum performance achieved by the system.
- **Rpeak:** The theoretical peak performance of the system.
- **Efficiency:** The ratio of Rmax to Rpeak, indicating how effectively the system utilizes its computational resources.

What about the AI workload? Well, floating point in TOP500 is double precision, but for AI works we need single or half precision. So, the TOP500 list is not the best for AI workload. Moreover, data movement is not considered in TOP500, but it is very important for AI workload.

<div align="right">

# 2
# Hardware

</div>

In the classical Von Neumann architecture there is only one processing unit (CPU) that processes instructions, the ALU is responsible for math and logic operations and the register stores data.
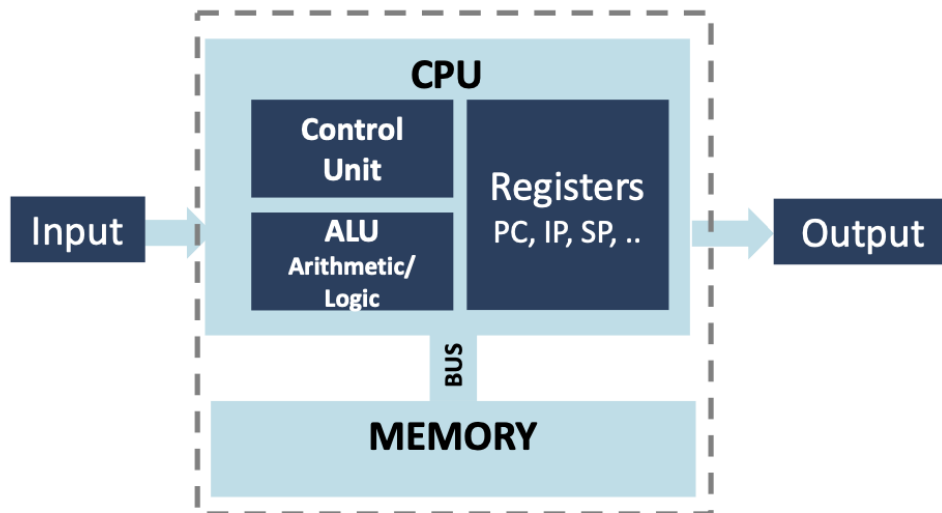


Figure 2.1: Von Neumann architecture

One instruction is executed at a time, the CPU fetches the instruction from the memory, decodes it and executes it. The CPU can access the memory to read or write data. Accessing any location in the mempory has always the same cost, this is called the **uniform memory access** (UMA).

## 2.1 Moore's Law

> **Definition**: *Moore's Law*
>
> It states that the number of transistors in a dense integrated circuit doubles about every two years.

How can we go from Moore's Law to processor performance? Through Dennard Scaling:

"Power density stays constant as transistors get smaller"

Intuitively,

- **Smaller transistors** → shorter propagation delay → faster frequency
- **Smaller transistors** → smaller capacitance → lower voltage

$$Power \propto Capacitance \times Voltage^2 \times Frequency$$

**But**. . . even with smaller transistors, we cannot continue reducing power, there are then two options:
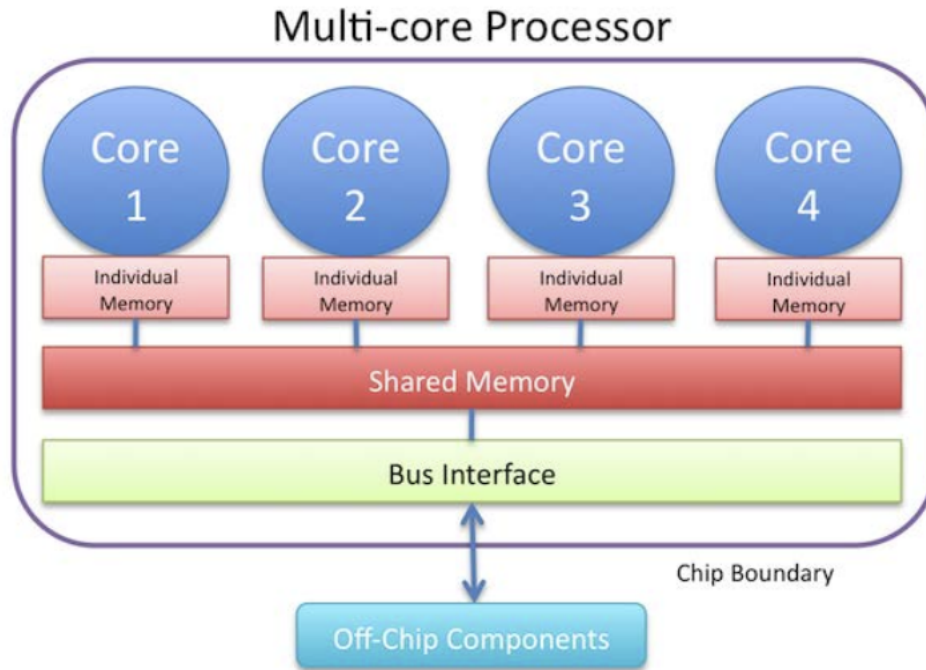
---

Figure 2.2: Multicore processors

- **Increase power**
- **Stop frequency scaling**

From 2006, single-core performance stopped increasing, so the only way to increase performance is to use more cores. The first solution is to write efficient software to make the efficient use of hardware resources. The second is to use specialized architectural solutions, like GPUs, FPGAs, etc.

Today, **CPUs are multicore processors**, lowering clock frequency because of power and heat dissipation, but packing more computing cores onto a chip. These cores will share some resources (memory, network, disk, ...) but are still capable of independent calculations.

## 2.2 Parallel Computers

Flynn Taxonomy (1966) is used to classify parallel computers based on the number of instruction streams and data streams.



Figure 2.3: Flynn Taxonomy

---

| | HW level | SW level |
|---|---|---|
| SISD | A Von Neumann CPU | no parallelism at all |
| MISD | On a superscalar CPU, different ports executing different *read* on the same data | • ILP on same data;<br>• Multiple tasks or threads operating on the same data |
| SIMD | Any vector-capable hardware, the vector registers on a core, a GPU, a vector processor, an FPGA, ... | data parallelism through vector instructions and operations |
| MIMD | Every multi-core/processor system; on a superscalar CPUs, different ports executing different ops on different data | • ILP on different data;<br>• Multiple tasks or threads executing different code on different data. |

Figure 2.4: Parallel computers

The essential components of a HPC cluster are:
- **Compute nodes**: the actual computers that perform the calculations
- **Interconnect**: the network that connects the compute nodes
- **Storage**: the disk space where data is stored
- **Software**: the operating system and the software stack that runs on the cluster



Figure 2.5: HPC cluster

A core is the smallest unit of computing, having one or more threads and is responsible for executing instructions.
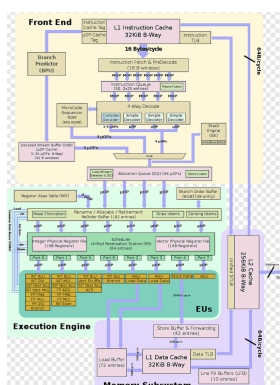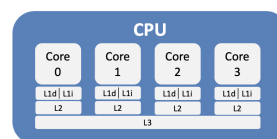


Figure 2.7: Cache hierarchy
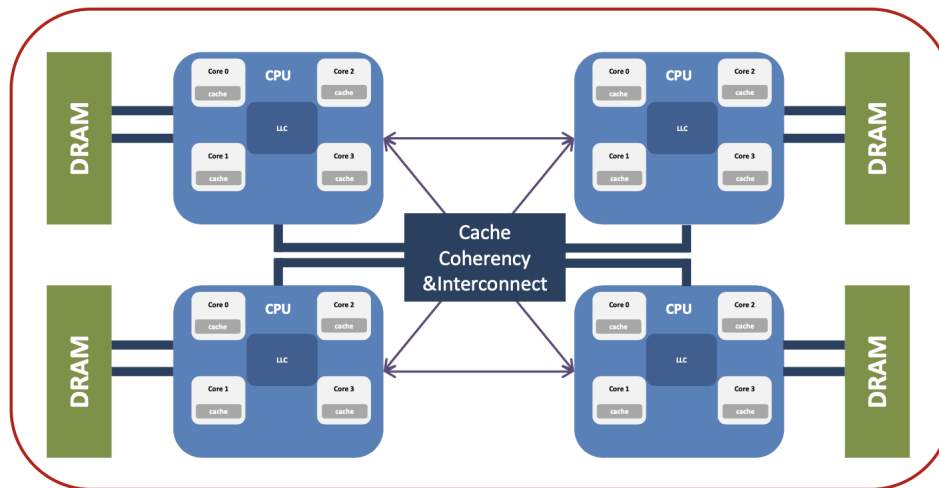
Cache hierarchy can have different topologies.



Figure 2.6: Core

Figure 2.8: Node topology

> 💡 **Tip**:
>
> **Cache hierarchy** has been invented cause accessing memory (moving data) takes almost 100x times computing (doing operations). The cache is a small memory that stores the most frequently accessed data. The cache is faster than the main memory but smaller. The cache is divided into levels, the first level is the fastest but the smallest, the second level is slower but bigger, and so on.

In a cluster there are different types of networks:

- **High-speed network**: used for communication between nodes (parallel computation, low latency/high bandwidth, infiniband or ethernet as examples)
- **I/O NETWORK**: used for communication with storage (I/O requests, latency not fundamental/good bandwidth, NFS, Lustre, GPFS as examples)
- **In band Management Network**: used for cluster management, monitoring, and control, LRMS (Load Resource Management System) as examples
- **Out of band Management Network**: used for cluster management, remote control of nodes and any other device, IPMI (Intelligent Platform Management Interface) as examples

What about **memory**?

It is fundamental and on a supercomputer there is a hybrid approach as for the memory placement. The memory on a single node can be accessed directly by all the cores on that node (**shared-memory**). When many nodes are used at a time, a process cannot directly access the memory on a different node, it needs to issue a request for that. That is named **distributed-memory**.

## Shared Memory

- **Uniform Memory Access (UMA)**:
  Each processor has uniform access to memory. Also known as symmetric multiprocessors (SMP).
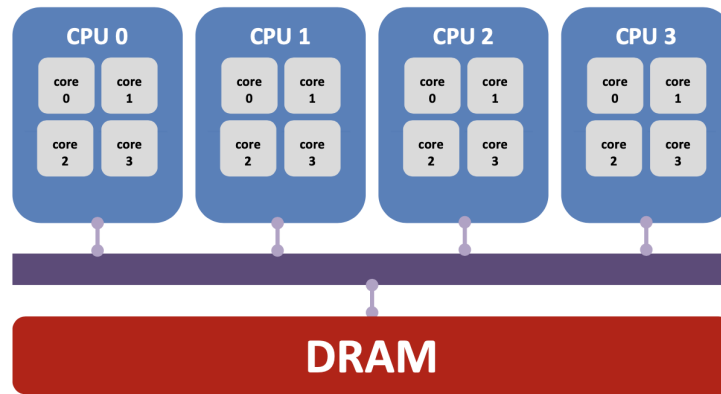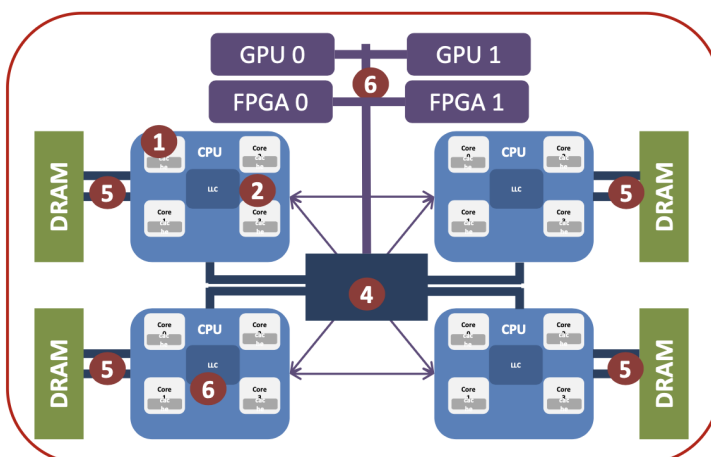
Figure 2.9: Shared memory - UMA

- **Non-Uniform Memory Access (NUMA)**: Time for memory access depends on location of data. Local access is faster on non-local access.

> ⚠ **Warning**: *Challenges for multicore*
>
> It aggravates the **Memory Wall problem**, where the memory access time is much slower than the CPU speed.
> - **Memory bandwidth**: the memory bandwidth is limited and shared among all the cores
> - **Memory latency**: the memory latency is high and can be a bottleneck
> - **Memory contention**: the memory contention can be a problem when multiple cores access the same memory location



1. ILP/SIMD units
2. Cores
3.
4. Socket/ccNuma domains
5. Inner cache levels
6. Multiple accelerators

Figure 2.10: Parallelism within a HPC node
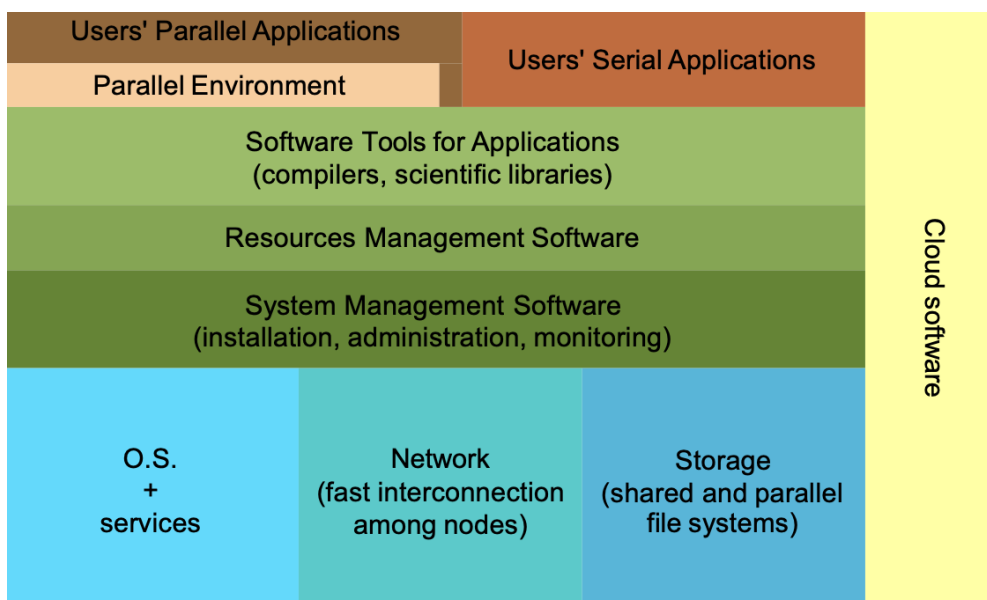
# HPC Software Stack



Figure 3.1: What we need

We refer to **Cluster Middleware** as the software stack that is used to manage the cluster. It is composed of several layers, each one with its own purpose. The first layer is the **Operating System**, which is the base layer of the stack. The second layer is the **Resource Manager**, which is used to manage the resources of the cluster. The third layer is the **Job Scheduler**, which is used to schedule the jobs on the cluster.

- The **Administration Software** is used to manage the cluster. It is used to install, configure, and monitor the cluster. It is also used to manage the users and groups on the cluster.

- The **Resource management and scheduling software (LRMS)** takes care of having many users and distribute process, balances the load among different users and schedules multiple tasks.

## 3.1 Resource Management Problem

Once we have a pool of users and a pool of resources, then some software is needed to control the available resources, to decide which application to execute based on available resources and to execute applications.
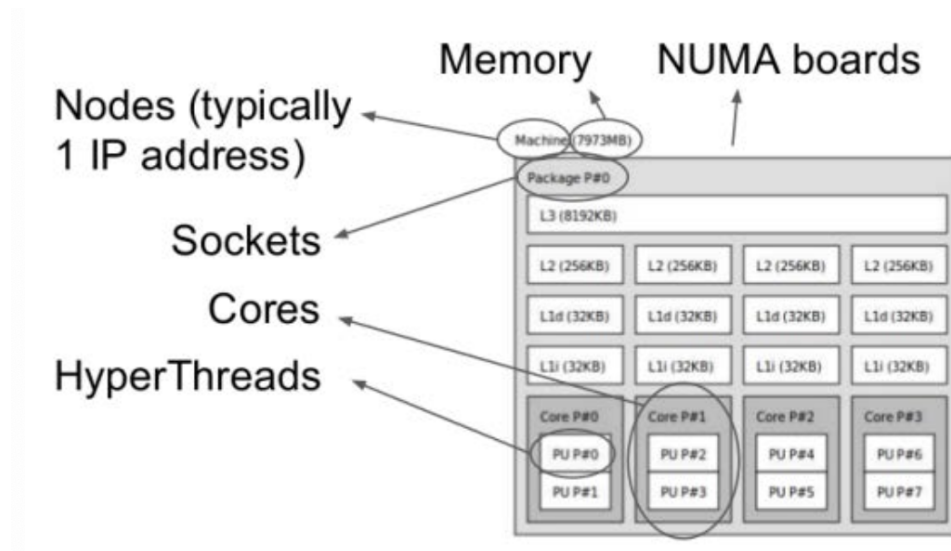
Figure 3.2: HPC resources

plus:
- network resources
- GPU/Accelerator
- Software resources

Let's now define some terms:
- **Batch Scheduler**: a software that manages the execution of jobs in a batch mode. It is used to schedule the jobs on the cluster.

> 👁 **Observation**: *Scheduling*
>
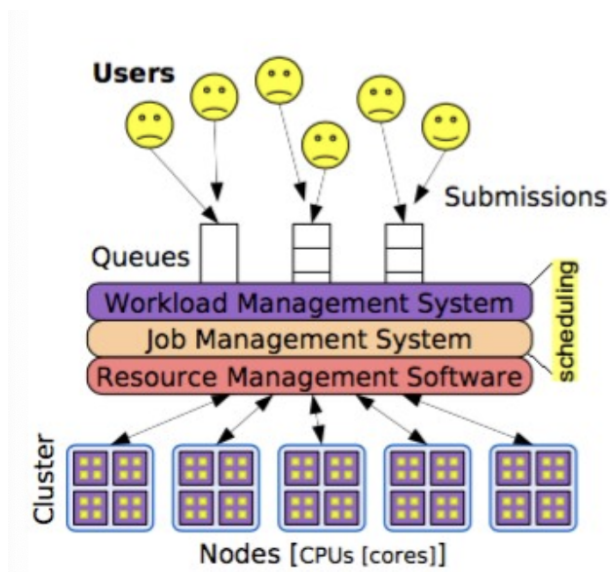> `Scheduling` is the process of assigning resources to jobs.



Figure 3.3: Batch scheduler

- **Resources Manager**: software that enable the jobs to connect the nodes and run.
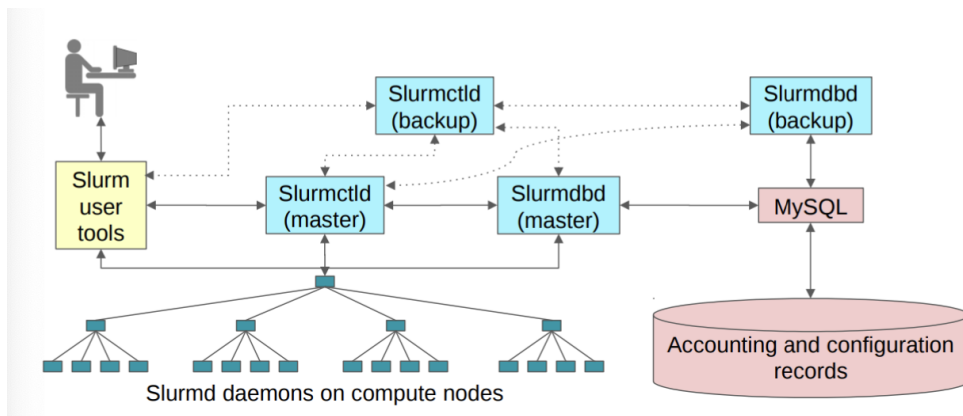
Figure 3.4: SLURM architecture

- **Node (aka Computing Node)**: computer used for its computational power.
- **Login/Master node**: it's through this node that the users will submit/launch/manage jobs.

> 📘 **Definition**: *SLURM*
>
> SLURM (Simple Linux Utility for Resource Management) is an open-source job scheduler for Linux and Unix-like kernels, used by many of the world's supercomputers and computer clusters.

- **Jobs**: Resource allocation requests.
- **Job Steps**: A job can be divided into multiple steps:
  - Typically an MPI and/or multi-threaded application program
  - Allocated resources from the job's allocation
  - A job that contains multiple job steps which can execute sequentially or concurrently
  - Lighter weight than jobs
- **Partitions**: Job queues with limits and access control
- **Qos**: Limits and policies

Table 3.1: SLURM Jobfile Parameters

| Directive | Meaning |
| --- | --- |
| `#!/bin/bash` | Shebang directive indicating the shell to use |
| `#SBATCH --job-name=<job_name>` | Assigns a name to the job |
| `#SBATCH --output=<file.out>` | Defines the output file |
| `#SBATCH --error=<file.err>` | Defines the error file |
| `#SBATCH --ntasks=<num_tasks>` | Specifies total number of tasks |
| `#SBATCH --cpus-per-task=<cpus>` | Sets desired CPUs per task |
| `#SBATCH --mem=<memory>` | Requests memory per node or per CPU option |
| `#SBATCH --time=<HH:MM:SS>` | Sets maximum runtime |
| `#SBATCH --partition=<partition>` | Selects partition (queue) |
| `#SBATCH --qos=<qos_class>` | Specifies quality of service class |

## 3.2  Scientific Software

This refers to the software above the Middleware:

- User's application (parallel and serial)
- Parallel Libraries and Tools
- Mathematical/Scientific Libraries
- I/O libraries
- Compilers

Here there is not so much standardization in HPC: every machine/app has a different software stack. Every code is optimized specifically for the hardware it runs on. This creates the so called **Dependency Hell**.

> **Definition**: *Dependency Hell*
>
> Dependency hell is a colloquial term for the frustration of some software users who have installed software packages which have dependencies on specific versions of other software packages.