

# Data Science and Advanced Programming — Lecture 12

## High Performance Computing

Simon Scheidegger  
Department of Economics, University of Lausanne, Switzerland

December 1st, 2025 | 12:30 - 16:00 | Internef 263

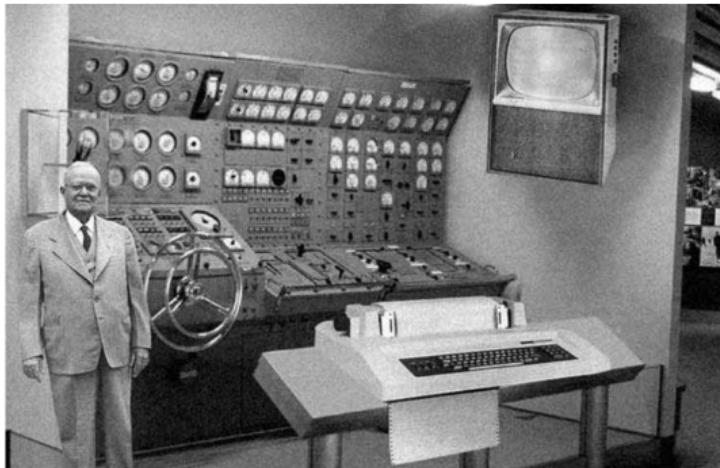
# Introduction to parallel and high-performance computing



# From making fire to flying rockets in 2 weeks



# From making fire to flying rockets in 2 weeks



Fortran: first appeared: 1957. For us, we start out here in early May 2023



By end of December, 2025, we are here...  
Since 2012: MPI, OpenMP, Fortran,  
C++, Cuda,...

# Scope of today's lecture

The purpose of the next 3 lectures is to give you an introduction to the main approaches to exploit modern **parallel** and **High-Performance Computing (HPC)** systems.

- ▶ Intended to provide only a very quick overview of the extensive and broad topic of parallel computing.
- ▶ Familiarize with the big picture, ideas and concepts.
- ▶ Familiarize with the main programming models.

**YOU SHOULD START TO THINK PARALLEL**

# Outline of this introductory lecture

1. Motivation — why should we use parallel programming.
2. Contemporary hardware.
3. Programming Models.

# Some Resources

## Full standard/API specification:

- ▶ <http://mpi-forum.org>
- ▶ <http://openmp.org>

## Tutorials:

- ▶ <https://computing.llnl.gov/tutorials/mpi/>
- ▶ <https://computing.llnl.gov/tutorials/openMP/>
- ▶ <http://cse-lab.ethz.ch/index.php/teaching/42-teaching/classes/577-hpcsei>

## Books:

- ▶ “Introduction to High Performance Computing for Scientists and Engineers”  
Georg Hager, Gerhard Wellein
- ▶ “Using Advanced MPI: Modern Features of the Message-Passing Interface”  
(Scientific and Engineering Computation) MIT Press, 2014 Gropp, W.; Höfler, T.; Thakur, R. & Lusk, E.

# Study complex phenomena

<https://hpc.llnl.gov/documentation/tutorials/introduction-parallel-computing-tutorial>



Rush Hour Traffic



Plate Tectonics

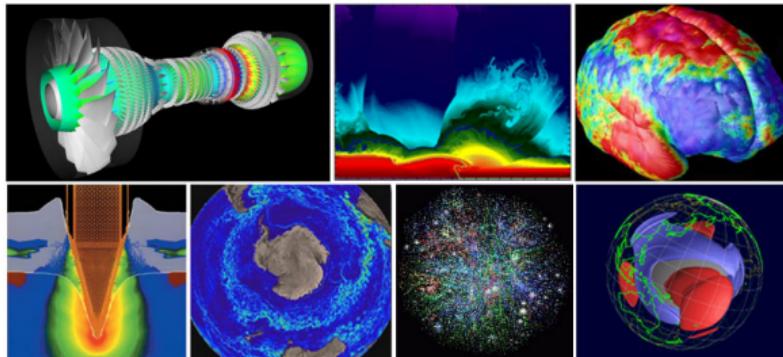


Weather

- ▶ In the natural world, many complex, interrelated events are happening at the same time, yet within a temporal sequence.
- ▶ Interrelated events happen concurrently.
- ▶ Compared to serial computing, parallel computing is much better suited for modelling, simulating and understanding complex, real world phenomena.

# Past: parallel computing = high-end science\*

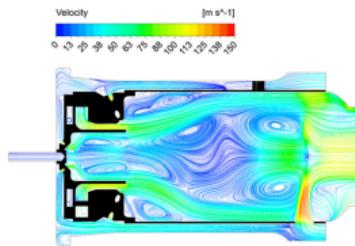
\*with special-purpose hardware



Parallel computing has been used to model difficult problems in many areas of science and engineering:

- ▶ Atmosphere, Earth, Environment
- ▶ Physics - applied, nuclear, particle, condensed matter, high pressure, fusion, photonics
- ▶ Bioscience, Biotechnology, Genetics, Chemistry, Molecular Sciences
- ▶ Geology, Seismology, Defense, Weapons
- ▶ ...

# Today: Parallel computing = every day (industry, academia) computing



Experiment



Macroscopic

Molecular Simulation



Microscopic

- ▶ Pharmaceutical design
- ▶ Data mining
- ▶ Financial modelling
- ▶ Advanced graphics, virtual reality
- ▶ Artificial intelligence/Machine learning
- ▶ ...



# Why parallel computing?

HPC helps with:

- ▶ Huge data
  - ▶ → Data management in memory
  - ▶ → Data management on disk
- ▶ Complex problems
  - ▶ → Time consuming algorithms
  - ▶ → Data mining
  - ▶ → Visualization
- ▶ Requires special purpose solutions in terms of:
  - ▶ → Processors
  - ▶ → Networks
  - ▶ → Storage
  - ▶ → Software
  - ▶ → Applications

# Why parallel computing\* (II)

\*List from P. Koumoutsakos

**Transmission speeds** Speed of a serial computer is directly dependent on how fast data can be moved through hardware. Absolute limits are the speed of light (30 cm/ns) and the transition limit of copper wire (9 cm/ns). Increased speeds necessitate increasing proximity of processing elements.

**Limits of miniaturization** Processor technology is allowing an increasing number of transistors to be placed on a chip. However, even with atomic-level components, a limit will be reached on how small components can be.

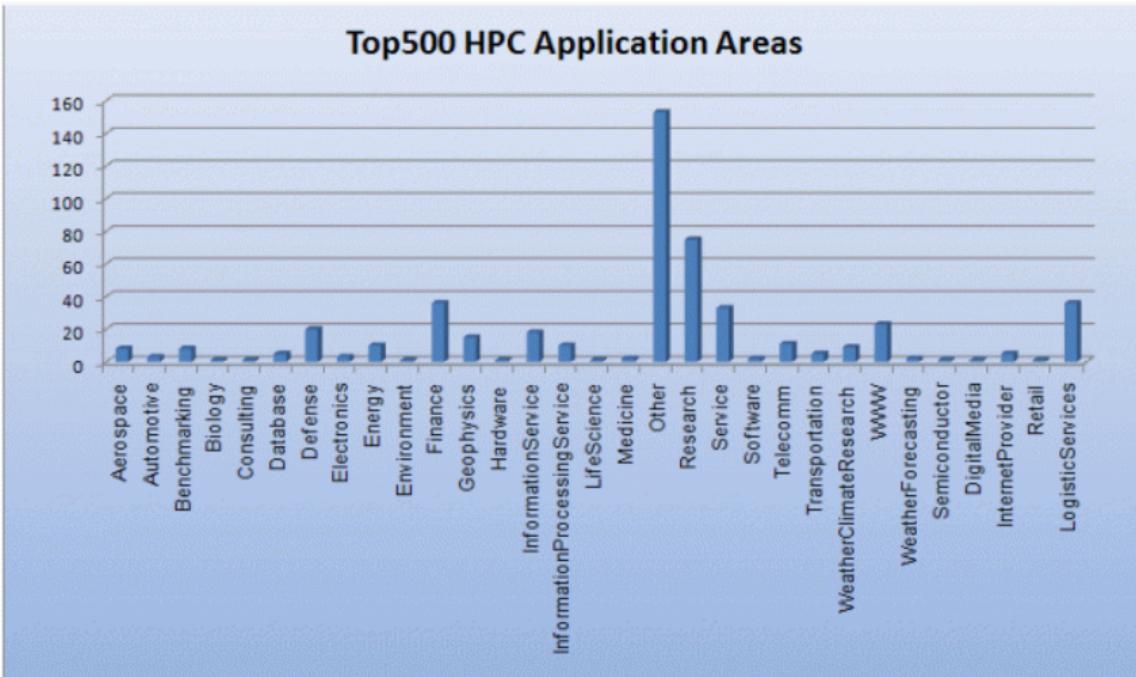
**Economic limits** It is increasingly expensive to make a single processor faster. Using a larger number of commodity processors to achieve same (or better) performance is less expensive.

**Energy limits** Limits imposed by cooling needs for chips and supercomputers.

## A way out

- ▶ Current computer architectures are increasingly relying upon hardware level parallelism.
- ▶ Multiple execution units.
- ▶ Multi-core.

# Who and what applies parallel computing?



# Why should YOU parallelize your code?

- ▶ A single core may be too slow to perform the required task(s) in a “tolerable” amount of time. The definition of “tolerable” certainly varies, but “overnight” is often a reasonable estimate.
- ▶ The memory requirements cannot be met by the amount of memory which is available on a single “Desktop”.



# Heterogeneity in Nature — Stylized modelling not always possible

- ▶ Supernova explosions - observable.
- ▶ Modelling in 1D:
  - ▶ → incl. fancy (GR/nuclear) physics.
  - ▶ → no explosions.
  - ▶ → only 1D models. due to lack of compute power (early 2000's).
- ▶ Only recently (last ~ 5y ): “realistic” 3D models possible to run on supercomputers
  - ▶ → resources now available.
  - ▶ → stylized models not sufficient.



February 24, 1987; SN 1987 A in the large Magellanic cloud. Closest SN next to earth since 1604 . Distance from earth:  $\sim 50\text{kpc}$ ;  $150'000 \text{ LY}$ .  
 $1\text{kpc} = 3.08 \times 10^{19} \text{ m}$ ,  $1\text{LY} = 9.46 \times 10^{15} \text{ m}$   
Progenitor:  $20M_{\text{sun}}$

# Supercomputers



- HPC is the “**third**” pillar of science — nowadays is the dawn of an era (in physics, chemistry, biology, ..., next to experiment and theory).
- “In Silico” experiments.
- Modern Supercomputers (e.g. world’s number 7 Summit at Oak Ridge National Lab)  $200 \times 10^{15}$  Flops/Sec. → **1 day vs. Laptop:  $\sim 30,000$  years.**

# Supercomputers



- → Move away from stylized models towards “realistically-sized” problems.
- **Creating “good” HPC software can be very difficult...**

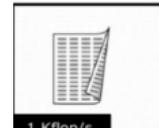
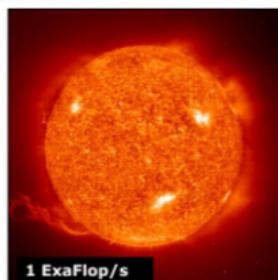
# Petascale Computers — How can we tap them?

\*Slide adjusted from O. Schenk

Modern Supercomputers (Summit)  $200 \times 10^{15}$  Flops/Sec.  $\rightarrow$  1 day  
vs. Laptop:  $\sim 30,000$ y.

**Let us say you can print:**

$10^{18}$  numbers (Exaflop) = 100,000,000 km  
(distance to the sun) stack printed per second  
**(Exaflop/s)**

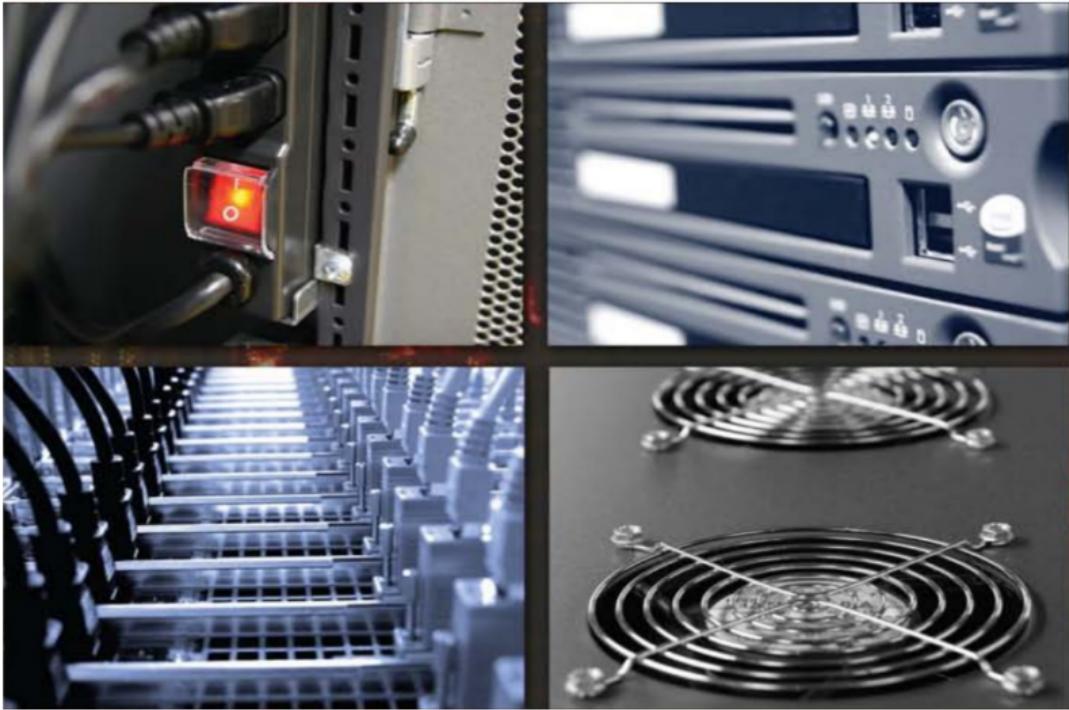


# June 2025 top500

<https://top500.org/lists/top500/2025/06/>

Rank	System	Cores	Rmax (PFlop/s)	Peak (PFlop/s)	Power (kW)
1	El Capitan - HPE Cray EX255a, AMD 4th Gen EPYC 24C 1.8GHz, AMD Instinct MI300A, Slingshot-11, T055, <b>HPE</b> DOE/NASA/LLNL United States	11,039,616	1,742.00	2,746.38	29,581
2	Frontier - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 4C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE Cray OS, <b>HPE</b> DOE/SC/Oak Ridge National Laboratory United States	9,066,176	1,353.00	2,055.72	24,607
3	Aurora - HPE Cray EX - Intel Exascale Compute Blade, Xeon CPU Max 9470 50C 2.4GHz, Intel Data Center GPU Max, Slingshot-11, <b>Intel</b> DOE/SC/Argonne National Laboratory United States	9,244,128	1,012.00	1,980.01	38,698
4	JUPITER Booster - BullSequana XH3000, GH Superchip 72C 30Hz, NVIDIA GH200 Superchip, Quad-Rail NVIDIA Infiniband NDR200, RedHat Enterprise Linux, <b>EVIDEN</b> EuroHPC/FZJ Germany	4,801,344	793.40	930.00	13,088
5	Eagle - Microsoft NDv5, Xeon Platinum 8480C 48C 20Hz, NVIDIA H100, NVIDIA Infiniband NDR, <b>Microsoft Azure</b> Microsoft Azure United States	2,073,600	561.20	846.84	
6	HPC6 - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 4C 2GHz, AMD Instinct MI250X, Slingshot-11, RHEL 8.9, <b>HPE</b> Eni S.p.A. Italy	3,143,520	477.90	606.97	8,461
7	Supercomputer Fugaku - Supercomputer Fugaku, A66FX 48C 2.2GHz, Tofu interconnect O, Fujitsu RIKEN Center for Computational Science Japan	7,630,848	442.01	537.21	29,899
8	Alps - HPE Cray EX254n, NVIDIA Grace 72C 3.1GHz, NVIDIA GH200 Superchip, Slingshot-11, HPE Cray OS, <b>HPE</b> Swiss National Supercomputing Centre (CSCS) Switzerland	2,121,600	434.90	574.84	7,124

## II. Contemporary hardware



# From a programming language to hardware

<http://cse-lab.ethz.ch/index.php/teaching/42-teaching/classes/577-hpcsei>

- ▶ A computer is a “stupid” device, only understands “on” and “off”.
- ▶ The symbols for these states are 0 and 1 (binary).
- ▶ First programmers communicated in 0 and 1.
- ▶ Later programs where developed to translate from symbolic notation to binary.  
The first was called “**assembly**”.

```
> add A, B (programmer writes in assembly language)  
>1000110010100000 (assembly translates to machine language)
```

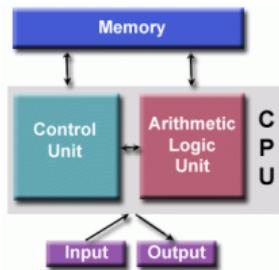
Advanced programming languages are better than “assembly”:

- ▶ programmer thinks in a more natural language.
- ▶ productivity of software development.
- ▶ portability.

# Basics: von Neumann Architecture

<https://hpc.llnl.gov/documentation/tutorials/introduction-parallel-computing-tutorial>

- ▶ Virtually all computers have followed this basic design. Comprised of **four main components**: **Memory, Control Unit, Arithmetic Logic Unit, Input/Output.**
- ▶ **Read/write, random access memory** is used to store both program instructions and data:
  - ▶ Program instructions are coded data which tell the computer to do something.
  - ▶ Data is simply information to be used by the program.
- ▶ **Control unit**
  - ▶ fetches instructions/data from memory, decodes the instructions and then sequentially coordinates operations to accomplish the programmed task.
- ▶ **Arithmetic unit**: performs basic arithmetic operations.
- ▶ **Input/Output**: interface to the human operator.



# Performance measures

- ▶ Execution time: time between start and completion of a task.
- ▶ Throughput/Bandwidth: total amount of work per elapsed time.
- ▶ Performance and execution time:

$$\text{Perf}_x = \frac{1}{\text{Exectime}_x}$$

$$\text{Perf}_x > \text{Perf}_y \implies \text{Exectime}_y > \text{Exectime}_x$$

$x$  is  $N$  times faster than  $y$  means:

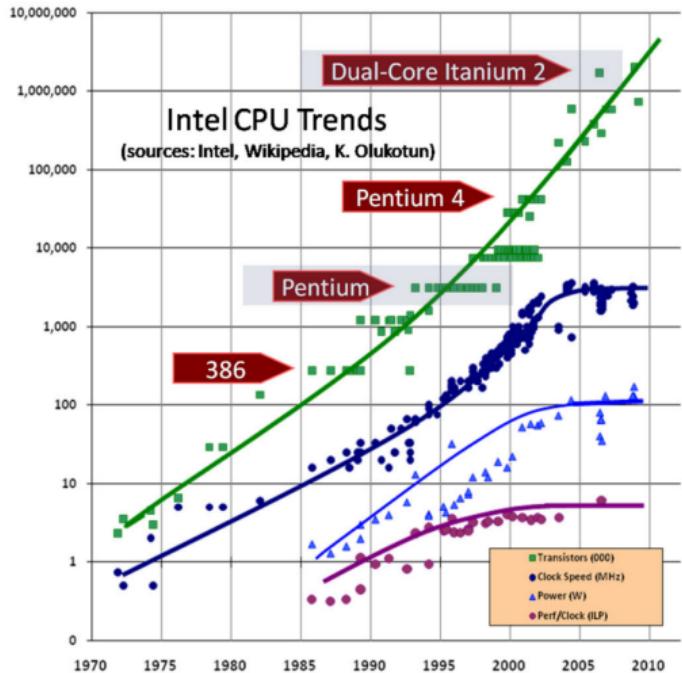
$$\frac{\text{Perf}_x}{\text{Perf}_y} = n = \frac{\text{Exectime}_y}{\text{Exectime}_x}$$

# The free lunch

- ▶ For a long time speeding up computations was a “free lunch”:
  - ▶ → the density of the transistors in chips increased, decreasing the size of integrated circuits.
  - ▶ → the clock speeds steadily rose, increasing the number of operations per second ( MHz to GHz ).
- ▶ But the free lunch has been over for a few years now:
  - ▶ → We are reaching the limitations of transistor density.
  - ▶ → Increasing clock frequency requires too much power.

→ We used to focus on floating point operations per second. Now we also think about floating point operations per Watt.

# The free lunch is over



## Clock speed cannot increase More:

1. heat (too much of it and too hard to dissipate).
2. current leakage.
3. power consumption (too high - also memory must be considered).

Processor performance doubles every 18 month:

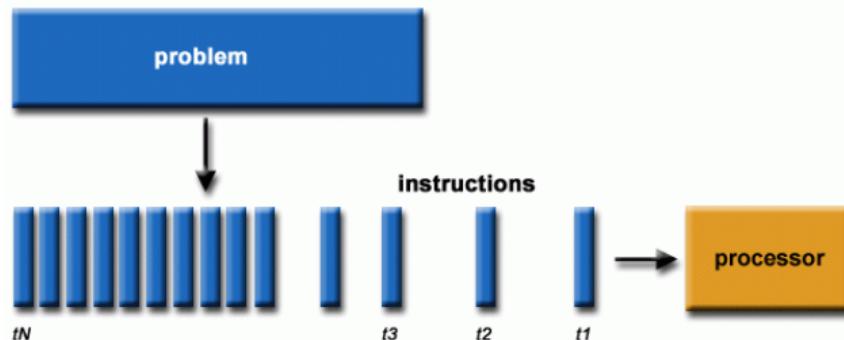
- ▶ Still true for the number of transistors.
- ▶ Not true for clock speed ( $W \sim \beta$  + hardware failures).
- ▶ Not true for memory/storage.

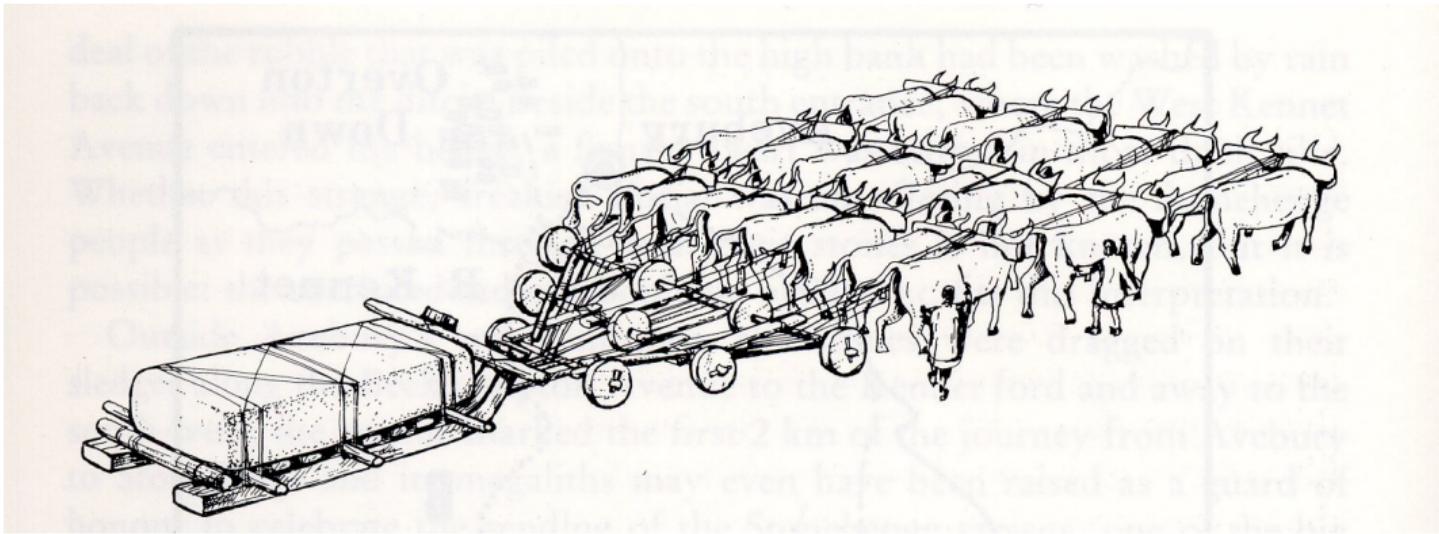
# What is parallel computing?

## Serial Computing:

**Traditionally, software has been written for serial computation.**

- ▶ A problem is broken into a discrete series of instructions.
- ▶ Instructions are executed serially one after another.
- ▶ Executed on a single processor.
- ▶ Only one instruction may execute at any moment in time.



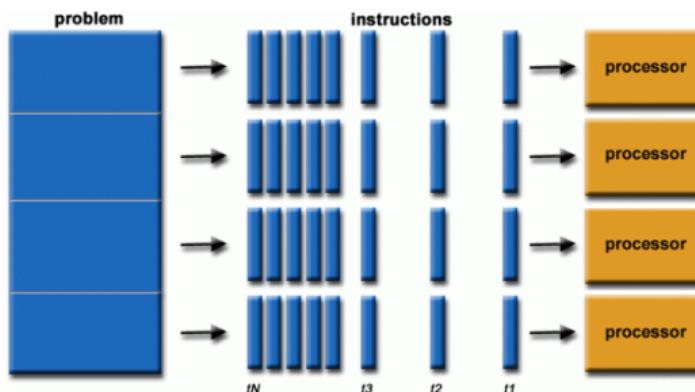


**"To pull a bigger wagon, it is easier to add more oxen than to grow a gigantic ox"** — (Skjellum et al. 1999)

# What is parallel computing II?

**Parallel Computing:** In the simplest sense, parallel computing is the simultaneous use of multiple compute resources to solve a computational problem.

- ▶ A problem is broken into a discrete series of instructions that can be solved concurrently.
- ▶ Each part is further broken down to a series of instructions.
- ▶ Instructions from each part execute simultaneously on different processors.
- ▶ An overall control/coordination mechanism is employed.



# Flynn's Taxonomy

<https://hpc.llnl.gov/documentation/tutorials/introduction-parallel-computing-tutorial>

- ▶ There are different ways to classify parallel computers.
- ▶ One of the most commonly used (since 1966) is **Flynn's Taxonomy**.
- ▶ It distinguishes multiprocessor computer architectures according to how they can be classified along the two independent dimensions of **Instruction Stream** and **Data Stream**.
- ▶ Each of these dimensions can have only one of two possible states: **Single** or **Multiple**.

<b>SISD</b> Single Instruction stream Single Data stream	<b>SIMD</b> Single Instruction stream Multiple Data stream
<b>MISD</b> Multiple Instruction stream Single Data stream	<b>MIMD</b> Multiple Instruction stream Multiple Data stream

# Single Instruction, Single Data (SISD)

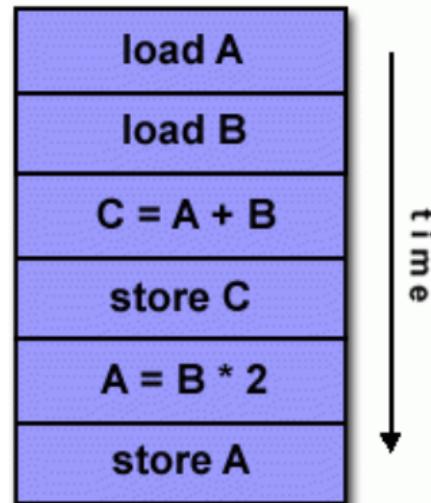
A serial (non-parallel) computer.

- ▶ **Single Instruction:**

- ▶ → Only one instruction stream is being acted on by the CPU during any one clock cycle.

- ▶ **Single Data:**

- ▶ → Only one data stream is being used as input during any one clock cycle.
- ▶ This is the oldest type of computer.
- ▶ Examples: single processor/core PCs.



# Single Instruction, Multiple Data (SIMD)

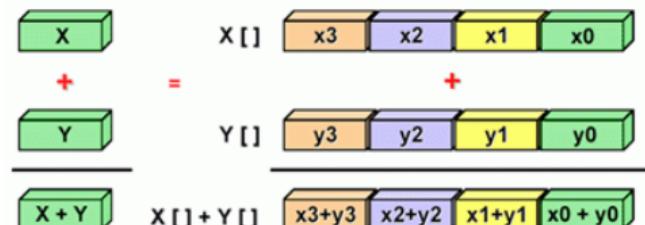
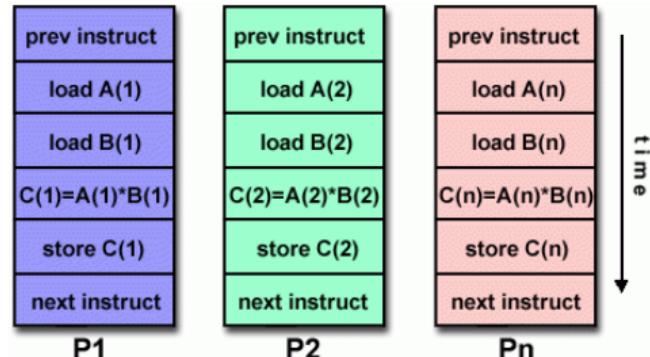
A type of parallel computer.

- ▶ **Single Instruction:**

- ▶ → All processing units execute the same instruction at any given clock cycle.

- ▶ **Multiple Data:**

- ▶ Each processing unit can operate on a different data element.
- ▶ Best suited for specialized problems characterized by a high degree of regularity.
- ▶ Most modern computers, particularly those with graphics processor units (GPUs) employ SIMD instructions and execution units.



# Multiple Instruction, Multiple Data (MIMD)

A type of parallel computer.

- ▶ **Multiple Instruction:**

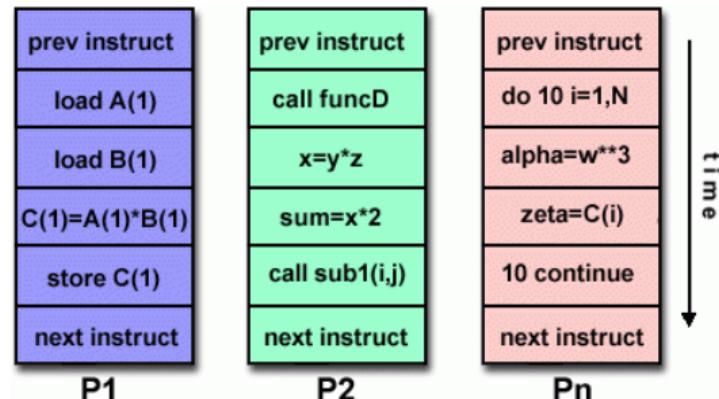
- ▶ → Every processor may be executing a different instruction stream.

- ▶ **Multiple Data:**

- ▶ → Every processor may be working with a different data stream.

- ▶ **Currently, the most common type of parallel computer most modern supercomputers fall into this category.**

- ▶ Note: many MIMD architectures also include SIMD execution subcomponents.



# Speedup, Efficiency & Amdahl's Law

$T(p, N)$  := time to solve problem of total size  $N$  on  $p$  processors.

**Parallel speedup:**  $S(p, N) = \frac{T(1, N)}{T(p, N)}$

→ Compute same problem with more processors in shorter time.

**Parallel Efficiency:**  $E(p, N) = \frac{S(p, N)}{p}$

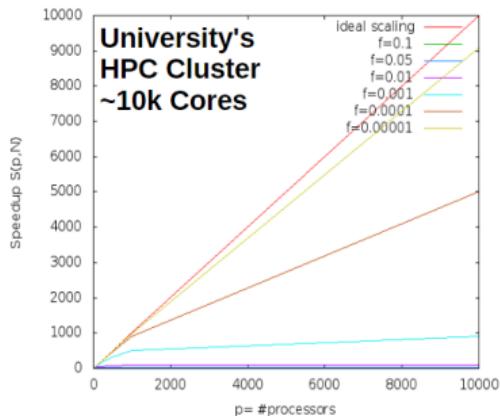
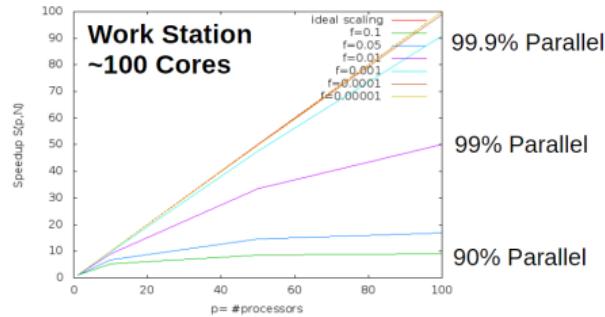
**Amdahl's Law:**  $T(p, N) = f * T(1, N) + \frac{(1-f)T(1, N)}{p}$

f...sequential part of the code that can not be done in parallel.

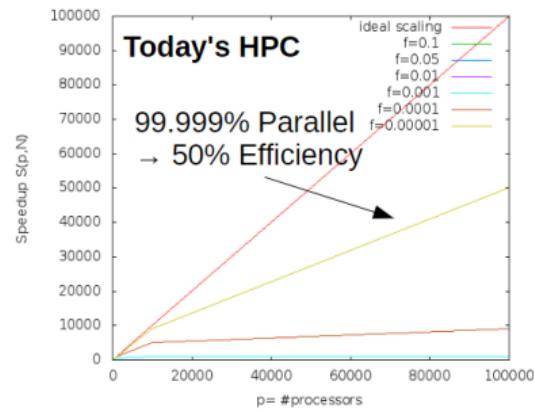
$$S(p, N) = \frac{T(1, N)}{T(p, N)} = \frac{1}{\left(f + \frac{1-f}{p}\right)}$$

For  $p \rightarrow \infty$ , speedup is limited by  $S(p, N) < \frac{1}{f}$

# Amdahl's Law: Scaling is tough



For  $p \rightarrow \infty$ , speedup is limited by  
 $S(p, N) < \frac{1}{f}$



# Weak versus strong scaling

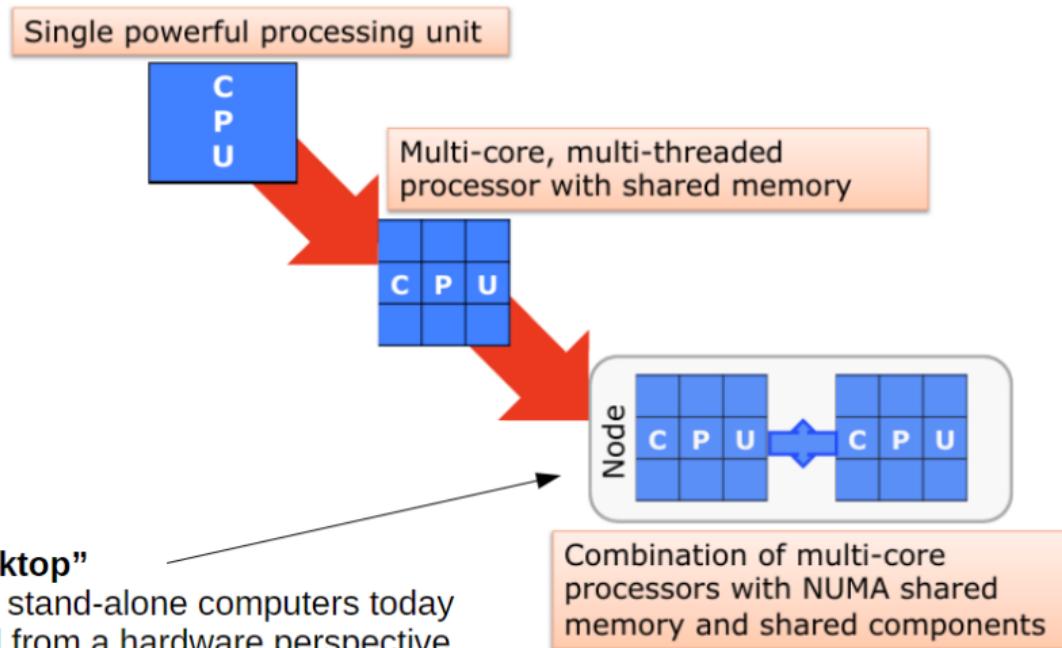
## **Strong scaling:**

Is defined as how the solution time varies with the number of processors for a fixed total problem size.

## **Weak scaling:**

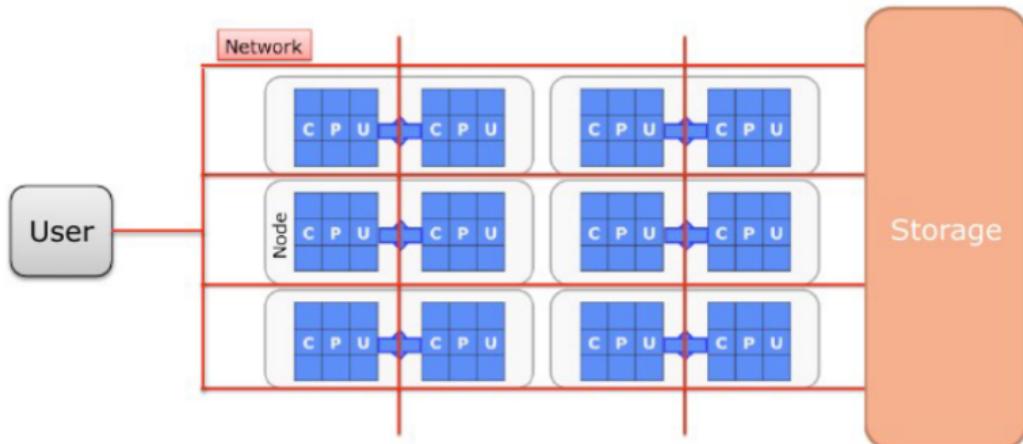
Is defined as how the solution time varies with the number of processors for a fixed problem size per processor.

# HPC systems: off-shelf components



# HPC systems

- ▶ Each compute node is a multi-processor parallel computer itself.
- ▶ Multiple compute nodes are networked together.
- ▶ Building blocks need to communicate (via the network) and give feedback (User/Storage).



# Strengths of Commodity Clusters

- ▶ Excellent performance to cost.
- ▶ Exploits economy of scale.
  - ▶ Through mass-production of components
  - ▶ Many competing cluster vendors
- ▶ Flexible just-in-place configuration.
  - ▶ Scalable up and down
- ▶ Rapid tracking of technology.
  - ▶ First to exploit newest components
- ▶ Programmable.
  - ▶ Uses industry standard programming languages and tools
- ▶ User empowerment.
  - ▶ Low cost, ubiquitous systems
  - ▶ Programming systems make it relatively easy to program for expert users
- ▶ nearly all of TOP-500 deployed systems commodity clusters.

# Two problems with this model

1. pure CPU-systems are becoming too expensive:
  - ▶ Power consumption.
  - ▶ Cooling and infrastructural costs (can be  $O$  (mio \$/ year)) → e.g. one Olympic swimming pool per hour cooling water needed.
2. It is hard to increase the performance:
  - ▶ Speed-up the processor.
  - ▶ Increase network bandwidth.
  - ▶ Make I/O faster.
  - ▶ ...

# HPC system's power consumption

...how to get 1000x more performance without building a nuclear power plant next to your HPC? (J. Shalf, September 09, Lausanne/Switzerland)

**Coming HPC systems are anticipated to draw enormous amounts of electrical power...**

**Example: N.1 Top 500**

Year (Nov.)	System	Performance (Tflop/s)	Electrical power (KW)
2002	Hearth Simulator	41	3200
2005	Blue Gene/L	367	1433
2008	Roadrunner	1105	2483
2009	Jaguar	2331	6950
2011	K computer	11280	12659
2012	Titan	27112	8209
2014	Tianhe-2	54902	17808
→2020	?????????	1000000	>100000 (100 MW!!!)

⇒ Not sustainable

# Improving performance: multi & many-cores

- ▶ Adding more and more “classical” cores can increase the computing power without having to increase the clock speed. However:
- ▶ → Hardware strongly limits the scalability.
- ▶ Progressively, a new generation of processors with less general, more specialized architecture, capable superior performance on a narrower range of problems, is growing:
- ▶ Many-core accelerators:
  - ▶ Nvidia Graphics Processor Units ([GPUs](#)).
  - ▶ Intel Xeon Phi Many Integrated Cores ([MIC](#)) architecture.
  - ▶ TPU (Tensor Processing Units)
  - ▶ *Idots*
- ▶ They all support a massively parallel computing paradigm.

# Why are accelerators more efficient than CPUs?

A few reasons:

- ▶ Lower frequencies of GPU clocks.
- ▶ More of the transistors in a GPU are working on the computation. **CPUs are more general purpose. They require energy to support such flexibility.**
- ▶ Single Instruction Multiple Data (SIMD) approach, which leads to a simpler architecture and instruction set.

# More hardware accelerators for ML\*



Microsoft

FPGA – field programmable gate arrays



TPU — Tensor Processing Unit

**ALTERA**  
Part of INTEL



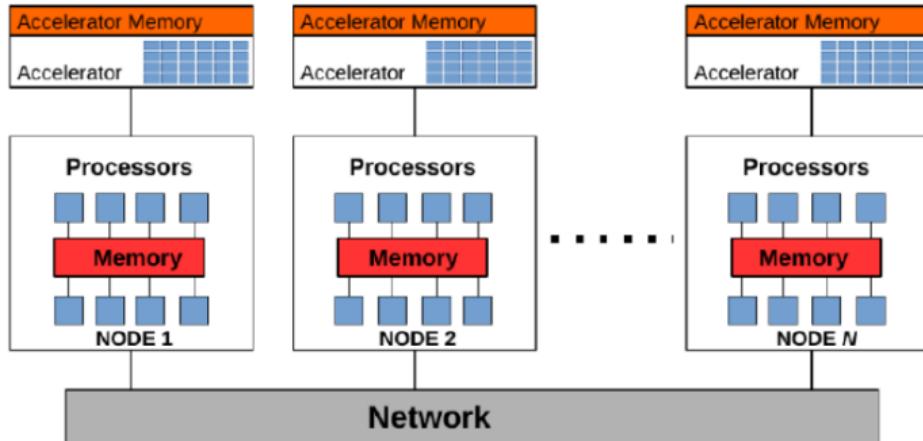
facebook

intel®

NVIDIA.

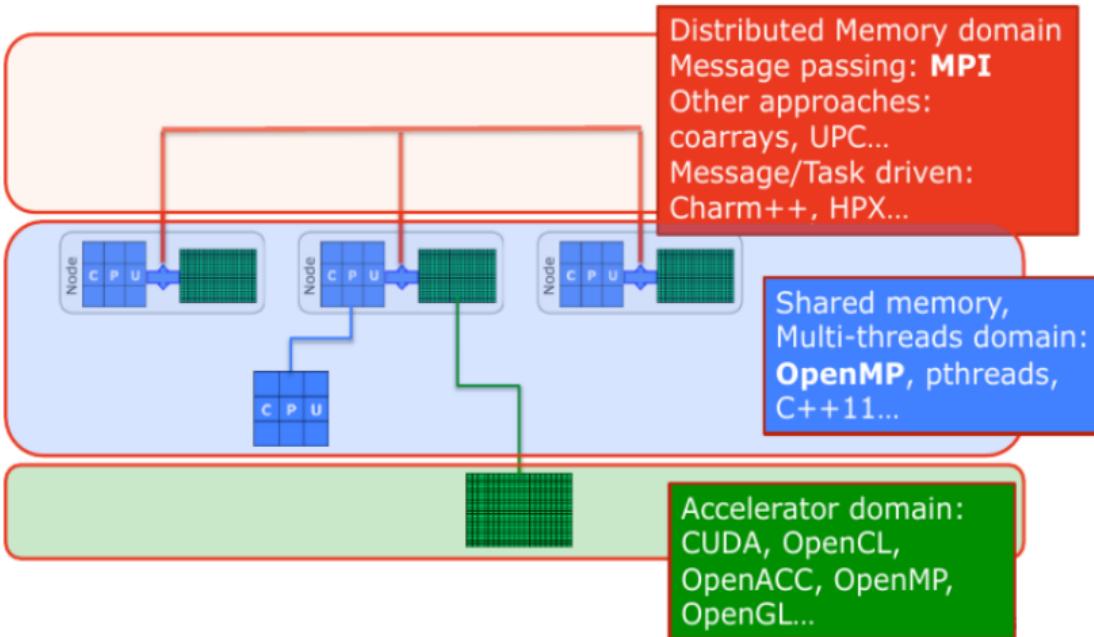
+ nervana

# Today's HPC systems



# Overall picture of programming models

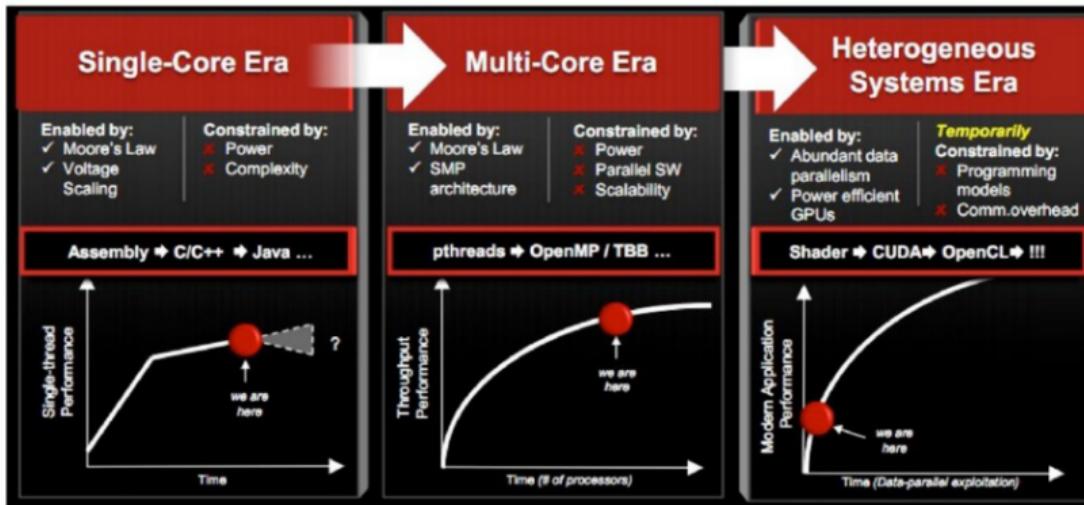
(Slide from C. Gheller)



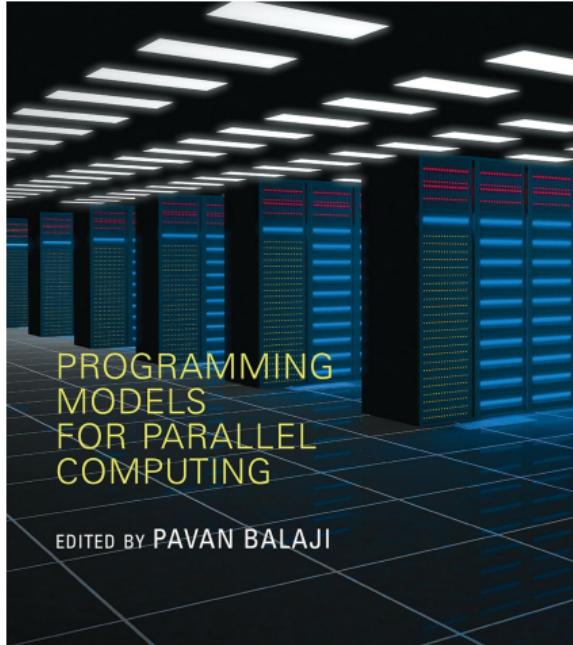
# Where are we going?

(Slide from C. Gheller)

**Heterogeneous systems are the next future. Right now, they represent the only viable solution for efficient high performance computing**

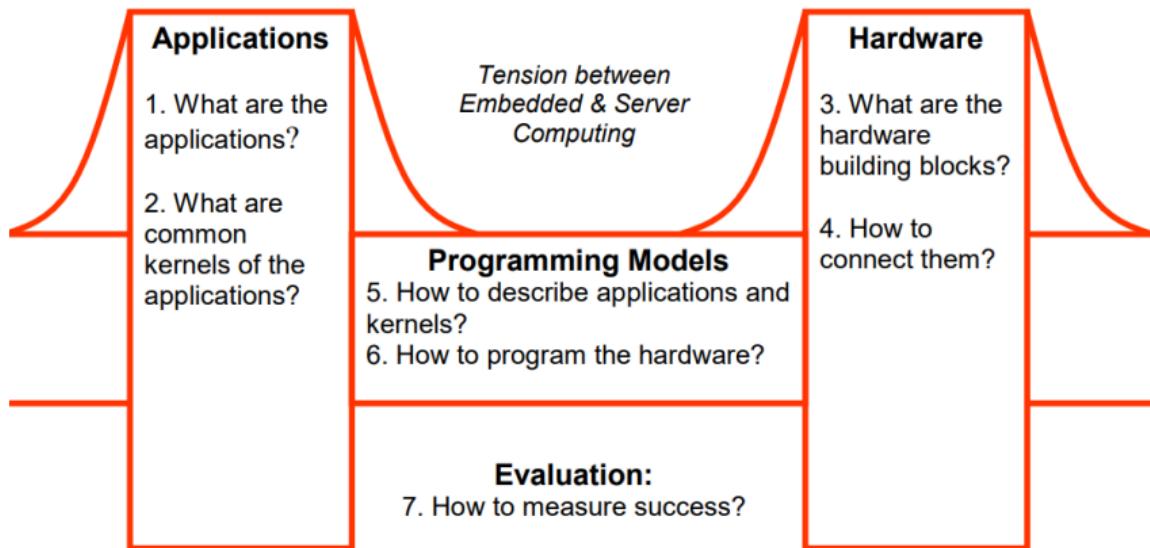


### 3. Programming Models



# Questions for parallel programming

Asanovic et. al. (2006) — The landscape of Parallel Computing Research: A View from Berkeley



**Figure 1.** A view from Berkeley: seven critical questions for 21<sup>st</sup> Century parallel computing.  
(This figure is inspired by a view of the Golden Gate Bridge from Berkeley.)

# Programming models: Overview

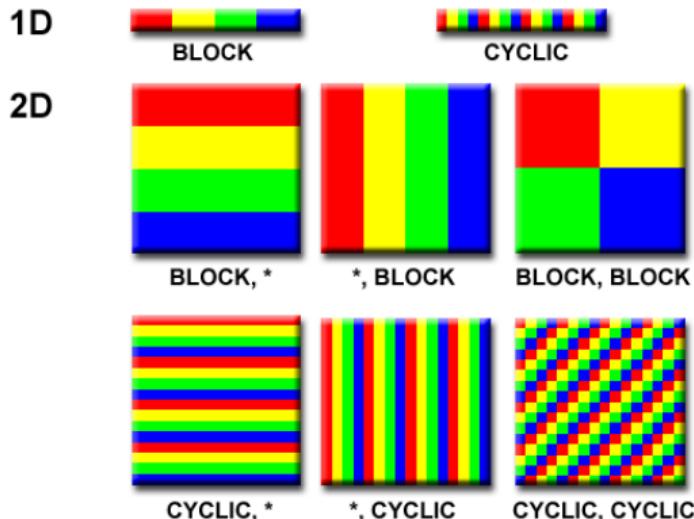
- ▶ There are several parallel programming models in common use:
  1. Shared Memory (without threads)
  2. Threads
  3. Distributed Memory / Message Passing
  4. Data Parallel
  5. Hybrid
  6. Single Program Multiple Data (SPMD)
  7. Multiple Program Multiple Data (MPMD)
- ▶ Parallel programming models exist as an abstraction above hardware and memory architectures.
- ▶ Although it might not seem apparent, these models are NOT specific to a particular type of machine or memory architecture. In fact, any of these models can (theoretically) be implemented on any underlying hardware.
- ▶ Which model to use?
- ▶ This is often a combination of what is available and personal choice. There is no "best" model, although there certainly are better implementations of some models over others.

# Parallel program design: Understand the problem

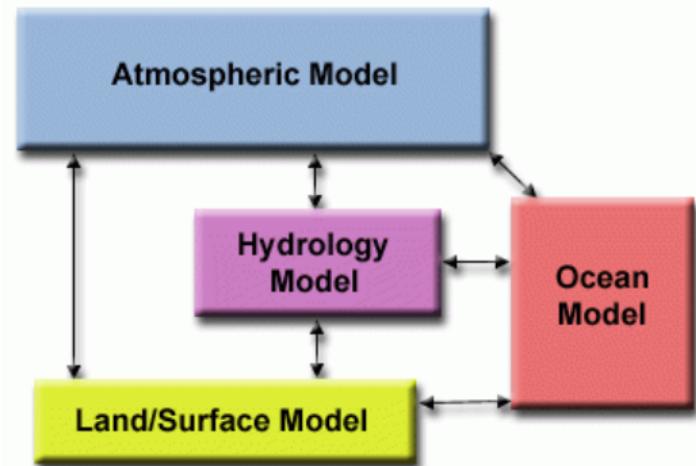
- ▶ The first step in developing parallel software is to first understand the problem that you wish to solve in parallel.
- ▶ Before spending time in an attempt to develop a parallel solution for a problem, determine whether or not the problem is one that can actually be parallelized.
  - ▶ Example of Parallelizable Problem: Monte Carlo integration
  - ▶ Example of a Non-parallelizable Problem: Calculation of the Fibonacci series  $(0, 1, 1, 2, 3, 5, 8, 13, 21, \dots)$  by use of the formula:  $F(n) = F(n - 1) + F(n - 2)$ 
    - ▶ This is a non-parallelizable problem because the calculation of the Fibonacci sequence as shown would entail dependent calculations rather than independent ones. The calculation of the  $F(n)$  value uses those of both  $F(n - 1)$  and  $F(n - 2)$ . These three terms cannot be calculated independently and therefore, not in parallel.
- ▶ **Identify hotspots** (where is the real work being done?)
- ▶ **Identify bottlenecks** in the program
  - ▶ are there disproportionately slow regions in code, e.g. I/O; can restructuring of program help?

# Decomposition of Problem

<https://hpc.llnl.gov/documentation/tutorials/introduction-parallel-computing-tutorial>



e.g. domain decomposition (e.g. for PDEs)



Functional decomposition (e.g. for climate modelling)

# Communication

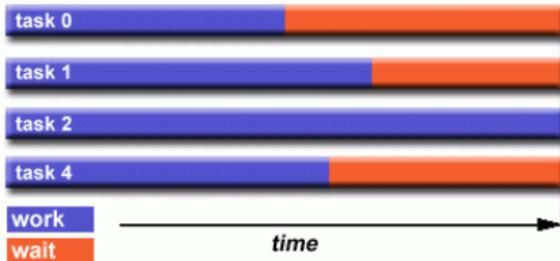
## Is communication needed?

The need for communication between tasks depend upon your problem

- ▶ you don't need communications: embarrassingly parallel (e.g. Monte Carlo)
- ▶ you need communications (e.g. stencil computation in a finite difference scheme (data dependency))
- ▶ Cost of communication (e.g. transmit data implies overhead)

# Load balancing

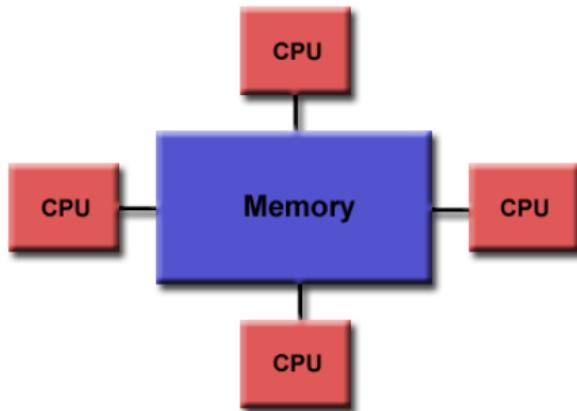
<https://hpc.llnl.gov/documentation/tutorials/introduction-parallel-computing-tutorial>



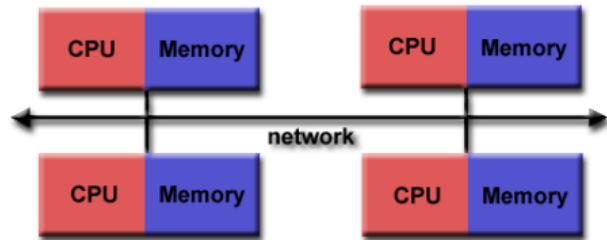
- ▶ Load balancing refers to the practice of distributing work among tasks so that all tasks are kept busy all of the time.
- ▶ **It can be considered a minimization of task idle time.**
- ▶ Load balancing is important to parallel programs for performance reasons. For example, if all tasks are subject to a barrier synchronization point, the slowest task will determine the overall performance.

# We will consider mainly OpenMP\* & MPI\*\*

\*<https://hpc-tutorials.llnl.gov/openmp/> (Open Multi-Processing) and \*\*<https://hpc-tutorials.llnl.gov/mpi/> (Message Passing Interface)



Shared Memory (OpenMP)



Distributed Memory (MPI)

# Using SIMD instruction sets → Vectorization

- ▶ Vectorization performs multiple operations in parallel on a core with a single instruction (SIMD - single instruction multiple data).
- ▶ Data is loaded into vector registers that are described by their width in bits:
  - ▶ 256 bit registers: 8x float, or 4x double
  - ▶ 512 bit registers: 16x float, or 8x double
- ▶ Vector units perform arithmetic operations on vector registers simultaneously.
- ▶ Vectorization is key to maximising computational performance.

# Vectorization illustrated

sequential

```
a [i] = b[i] + c[i] * d[i];
```

**a[i]**  
=  
**b[i]**  
+  
**c[i]**  
\*  
**d[i]**

8×vectorization

```
a [i:8] = b[i:8] + c[i:8] * d[i:8];
```

**a[i]**    **a[i+1]**    **a[i+2]**    **a[i+3]**    **a[i+4]**    **a[i+5]**    **a[i+6]**    **a[i+7]**  
=            =            =            =            =            =            =            =  
**b[i]**    **b[i+1]**    **b[i+2]**    **b[i+3]**    **b[i+4]**    **b[i+5]**    **b[i+6]**    **b[i+7]**  
+            +            +            +            +            +            +            +  
**c[i]**    **c[i+1]**    **c[i+2]**    **c[i+3]**    **c[i+4]**    **c[i+5]**    **c[i+6]**    **c[i+7]**  
\*            \*            \*            \*            \*            \*            \*            \*  
**d[i]**    **d[i+1]**    **d[i+2]**    **d[i+3]**    **d[i+4]**    **d[i+5]**    **d[i+6]**    **d[i+7]**

In an optimal situation all this is carried out by the compiler automatically.

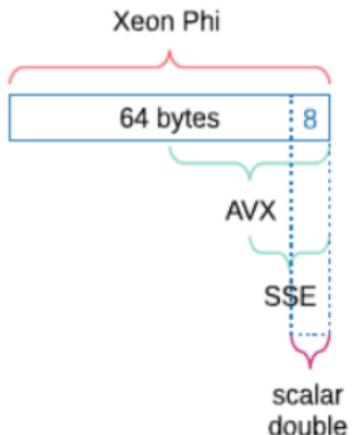
Compiler directives can be used to give hints as to where vectorization is safe and/or beneficial.

---

```
1 ! vectorized part
2 rest = mod(N,4)
3 do i=1,N-rest,
4   load R1 = [x(i),x(i+1),x(i+2),x(i+3)]
5   load R2 = [y(i),y(i+1),y(i+2),y(i+3)]
6   ! "packed" addition (4 SP flops)
7   R3 = ADD(R1,R2)
8   store [r(i),r(i+1),r(i+2),r(i+3)] = R3
9 enddo
10 ! remainder loop
11 do i=N-rest+1,N
12   r(i) = x(i) + y(i)
13 enddo
```

---

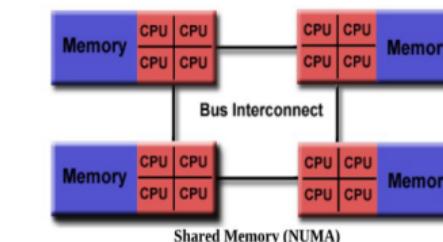
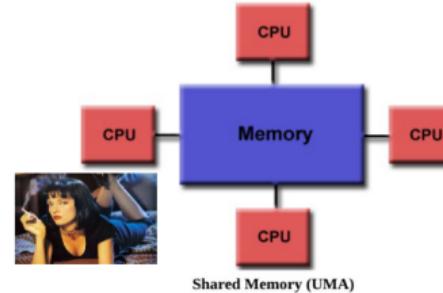
# Advanced Vector Extensions (AVX)



**Fig. 7.** Vector registers on modern CPUs: a scalar program can utilize only 1/4 of computational parallelism on AVX-enabled CPUs, e.g. the SandyBridge.

# Shared memory systems

- ▶ Process can access same **GLOBAL** memory
- ▶ Uniform Memory Access (**UMA**) model
  - ▶ Access time to memory is uniform.
  - ▶ Local cache, all other peripherals are shared.
- ▶ Non-Uniform Memory Access (**NUMA**) model
  - ▶ Memory is physically distributed among processors.
  - ▶ Global virtual address spaces accessible from all processors.
  - ▶ **Access time to local and remote data is different.**
- ▶ **OpenMP**, but other solutions available (e.g. Intel's TBB).



N-



# Shared Memory — Pros & Cons

## Pros:

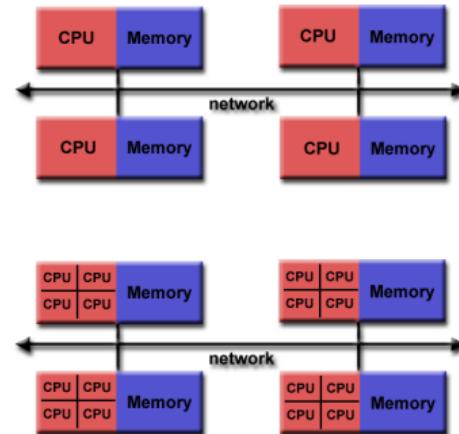
- ▶ Global address space provides a **user-friendly programming perspective** to memory.
- ▶ Data sharing between tasks is both fast and uniform due to the proximity of memory to CPUs.

## Cons:

- ▶ Lack of scalability between memory and CPUs. Adding more CPUs geometrically increases traffic on the shared memory-CPU path...
- ▶ Programmer responsibility for synchronization constructs that ensure “correct” access of global memory.

# Distributed-memory parallel programming (with MPI)

- ▶ We need to use explicit message passing, i.e., communication between processes:
  - ▶ Most tedious and complicated but also the most flexible parallelization method.
- ▶ Message passing is required if a parallel computer is of **distributed-memory** type, i.e., if **there is no way for one processor to directly access the address space of another**.
- ▶ However, it can also be regarded as a programming model and used on shared-memory or **hybrid systems** as well.
- ▶ → **Message Passing Interface (MPI)**.



# Parallel programming with MPI

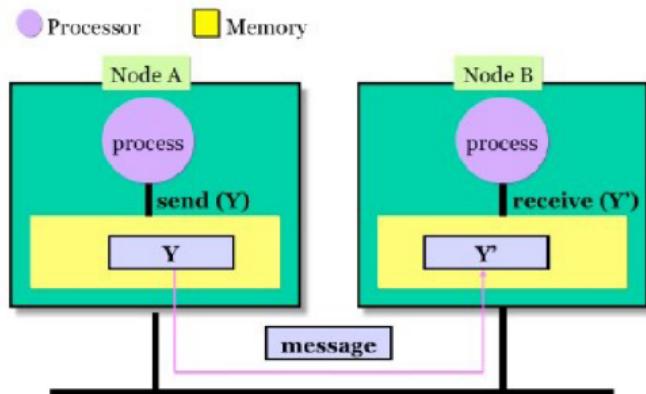
- ▶ The Message Passing Interface Standard (MPI) is a **message passing library** standard based on the consensus of the **MPI Forum**, which has over 40 participating organizations, including vendors, researchers, software library developers, and users.
- ▶ The goal of the MPI is to establish a **portable**, efficient, and flexible standard for message passing that will be widely used for writing message passing programs. As such, MPI is the first standardized, vendor independent, message passing library.
- ▶ **Several MPI implementations exist** (Mpich, OpenMPI, IntelMPI, CrayMPI, ...).
- ▶ The current MPI standard currently defines over 500 functions, and it is beyond the scope of this introduction even try to cover them all. Here, I **will concentrate on the important concepts of message passing and MPI in particular**, and **provide** some knowledge that will enable you to consult more advanced textbooks and tutorials.

# Message Passing Interface (MPI)

- Resources are LOCAL (different from shared memory).
- Each process runs in an “isolated” environment. Interactions require **Messages** to be exchanged.
- Messages can be: **instructions, data, synchronization**.
- MPI works also on Shared Memory systems.
- Time to exchange messages is much larger than accessing local memory.

Message Passing is a COOPERATIVE Approach, based on 3 operations:

- SEND (a message)
- RECEIVE (a message)
- SYNCHRONIZE



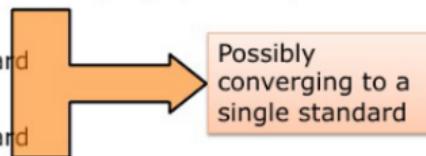
# MPI availability

- ▶ MPI is standard defined in a set of documents compiled by a consortium of organizations : <http://www.mpi-forum.org/>
- ▶ In particular the MPI documents define the APIs for C, C++, FORTRAN77 and FORTRAN 90.
- ▶ Bindings available for Perl, Python, Java...
- ▶ In all systems MPI is implemented as **a library of subroutines/functions** over the network drivers and primitives.

# GPU Programming

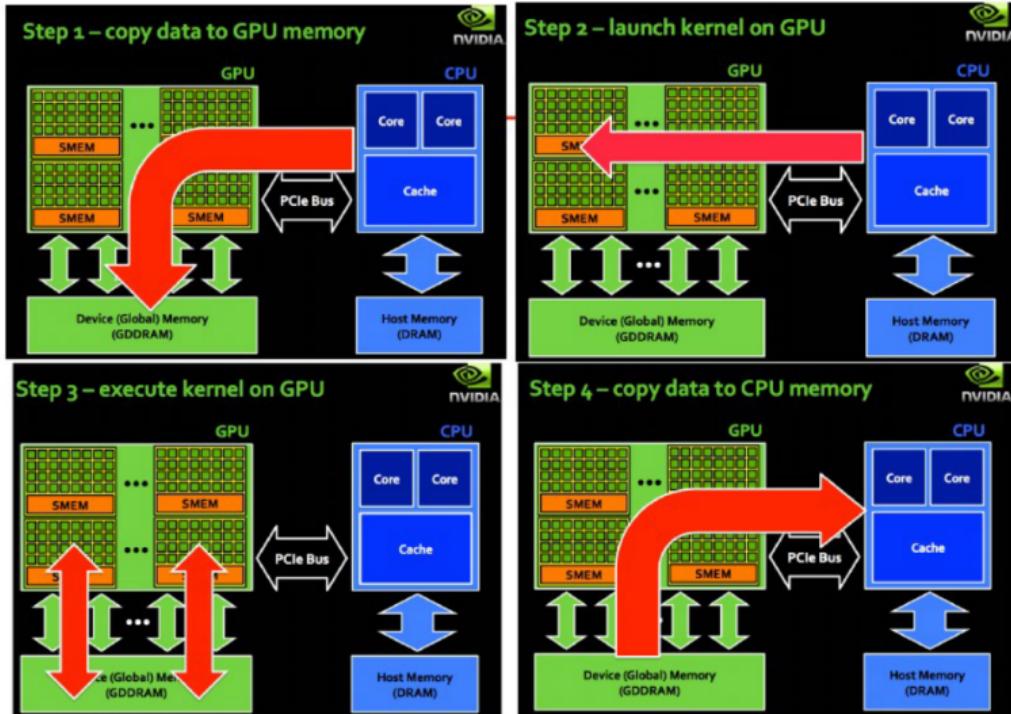
## API currently available

- **CUDA**
  - NVIDIA product, best performance, low portability
- **OpenCL**
  - Open standard, API similar to CUDA, high portability
- **OpenACC**
  - Directive based open standard
- **OpenMP**
  - Directive based open standard
- **Libraries**
  - MAGMA, cuFFT, Thrust, CULA, cuBLAS, cuSOLVER...
- **C++11**
  - High level abstraction



**Supported languages:** **C, C++, Fortran but also Java, Python (not from all the APIs)**

# GPU programming II



# Degrees of Parallelism

- ▶ Current supercomputers are hybrid.
- ▶ There are 4 main degrees of parallelism on multi-core (HPC) systems.
- ▶ Parallelism among nodes/sockets (e.g. MPI).
- ▶ Parallelism among cores in a socket (e.g. OpenMP).
- ▶ Vector units in each core (e.g. AVX).
- ▶ Accelerators...

All should be exploited at once (→lecture next week)

# Questions?

