

#### The Panorama of Parallel and High Performance Computing

Concurrency and Parallelism — 2017-18 Master in Computer Science (Mestrado Integrado em Eng. Informática)

Joao Lourenço <joao.lourenco@fct.unl.pt>

Slides based in: https://computing.llnl.gov/tutorials/parallel\_comp/

## Bibliograpy

Chapter 1 of book

McCool M., Arch M., Reinders J.; Structured Parallel Programming: Patterns for Efficient Computation; Morgan Kaufmann (2012); ISBN: 978-0-12-415993-8





- 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 (sequentially)
    - Only one instruction may execute at any moment in time



- Is the simultaneous use of multiple compute resources to solve a computational problem:
  - To be executed using multiple processors



- Is the simultaneous use of multiple compute resources to solve a computational problem:
  - To be executed using multiple processors
  - A problem is broken into discrete parts that can be solved concurrently



- Is the simultaneous use of multiple compute resources to solve a computational problem:
  - To be executed using multiple processors
  - A problem is broken into discrete parts that can be solved concurrently
  - Each part is further broken own to a series of instructions



- Is the simultaneous use of multiple compute resources to solve a computational problem:
  - To be executed using multiple processors
  - A problem is broken into discrete parts that can be solved concurrently
  - Each part is further broken own to a series of instructions
  - Instructions from each part execute simultaneously on different processors



- Is the simultaneous use of multiple compute resources to solve a computational problem:
  - To be executed using multiple processors
  - A problem is broken into discrete parts that can be solved concurrently
  - Each part is further broken own to a series of instructions
  - Instructions from each part execute simultaneously on different processors
  - An overall control/coordination mechanism is employed



- Is the simultaneous use of multiple compute resources to solve a computational problem:
  - To be executed using multiple processors
  - A problem is broken into discrete parts that can be solved concurrently
  - Each part is further broken own to a series of instructions
  - Instructions from each part execute simultaneously on different processors
  - An overall control/coordination mechanism is employed



- The computational problem should be able to:
  - Be broken apart into discrete pieces of work that can be solved simultaneously
  - Execute multiple program instructions at any moment in time
  - Be solved in less time with multiple compute resources than with a single compute resource
- The computing resources might be:
  - A single computer with multiple processors
  - An arbitrary number of computers connected by a network (real or virtual systems)
  - A combination of both



![](_page_10_Picture_2.jpeg)

![](_page_10_Picture_3.jpeg)

login / remote partition server node

![](_page_10_Picture_5.jpeg)

gateway node

#### The Real World is Massively Parallel

- In the natural world, many complex, interrelated events are happening at the same time, yet within a temporal sequence.
- Compared to serial computing, parallel computing is much better suited for modeling, simulating and understanding complex, real world phenomena.
- For example, imagine modeling serially the following systems.

#### The Real World is Massively Parallel

![](_page_12_Picture_1.jpeg)

**Galaxy Formation** 

**Planetary Movments** 

![](_page_12_Picture_4.jpeg)

**Climate Change** 

![](_page_12_Picture_6.jpeg)

**Rush Hour Traffic** 

![](_page_12_Picture_8.jpeg)

**Plate Tectonics** 

![](_page_12_Picture_10.jpeg)

![](_page_12_Picture_11.jpeg)

Weather

![](_page_12_Picture_13.jpeg)

Auto Assembly

Jet Construction

## Uses for Parallel Computing

 Modeling difficult problems in many areas of science and engineering

![](_page_13_Picture_2.jpeg)

# Uses for Parallel Computing

#### Industrial and Commercial

![](_page_14_Picture_2.jpeg)

- Oil exploration
- Web search engines, web based business services
- Medical imaging and diagnosis
- Pharmaceutical design
  - Financial and economic modeling
- Management of national and multi-national corporations
- Advanced graphics and virtual reality, particularly in the entertainment industry
- Networked video and multi-media technologies
- Collaborative work environments

![](_page_14_Picture_12.jpeg)

# Why Use Parallel Computing?

- Save time and/or money
  - In theory, throwing more resources at a task will shorten its time to completion, with potential cost savings. Parallel computers can be built from cheap, commodity components.

![](_page_15_Picture_3.jpeg)

# Why Use Parallel Computing?

- Solve larger problems
  - Many problems are so large and/or complex that it is impractical or impossible to solve them on a single computer. For example:
    - "Grand Challenge" (en.wikipedia.org/wiki/Grand\_Challenge) problems requiring PetaFLOPS and PetaBytes of computing resources.
    - Web search engines/databases processing millions of transactions per second.

![](_page_16_Picture_5.jpeg)

# Why Use Parallel Computing?

- Use of non-local resources
  - Using compute resources on a wide area network, or even the Internet when local compute resources are scarce. For example:
    - SETI@home (setiathome.berkeley.edu) over 1.6 million users, 4 million computers, in nearly every country in the world. Source: https://setiathome.berkeley.edu/stats.php (Sep, 2016).
    - Folding@home (folding.stanford.edu) uses over 320,000 computers globally (Sep, 2016)

![](_page_17_Figure_5.jpeg)

### Limits to serial computing

- Both physical and practical reasons pose significant constraints to simply building ever faster serial computers:
  - Transmission speeds
    - the speed of a serial computer is directly dependent upon how fast data can move through hardware. Absolute limits are the speed of light (30 cm/nanosecond) and the transmission limit of copper wire (9 cm/nanosecond). Increasing speeds requires increasing proximity of processing elements.
  - Limits to miniaturization
    - processor technology is allowing an increasing number of transistors to be placed on a chip. However, even with molecular or atomic-level components, a limit will be reached on how small components can be.

#### Limits to serial computing

![](_page_19_Figure_1.jpeg)

## Limits to serial computing

- Both physical and practical reasons pose significant constraints to simply building ever faster serial computers:
  - Economic limitations
    - it is increasingly expensive to make a single processor faster.
      Using a larger number of moderately fast commodity processors to achieve the same (or better) performance is less expensive.
  - Current computer architectures are increasingly relying upon hardware level parallelism to improve performance:
    - Multiple execution units
    - Pipelined instructions
    - Multi-core

![](_page_20_Figure_8.jpeg)

#### The Future

- During the past 20+ 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.
- In this same time period, there has been a greater than 1000x increase in supercomputer performance, with no end currently in sight.
- The race is already on for Exascale Computing! (10<sup>18</sup> FLOPS)

#### Performance development

![](_page_22_Figure_1.jpeg)

## Top500.org

Rank	System	Cores	Rmax (TFlop/s)	Rpeak (TFlop/s)	Power (kW)
1	<b>Sunway TaihuLight</b> - Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway , NRCPC National Supercomputing Center in Wuxi China	10,649,600	93,014.6	125,435.9	15,371
2	Tianhe-2 (MilkyWay-2) - TH-IVB-FEP Cluster, Intel Xeon E5-2692 12C 2.200GHz, TH Express-2, Intel Xeon Phi 31S1P, NUDT National Super Computer Center in Guangzhou China	3,120,000	33,862.7	54,902.4	17,808
3	<b>Piz Daint</b> - Cray XC50, Xeon E5-2690v3 12C 2.6GHz, Aries interconnect , NVIDIA Tesla P100 , <b>Cray Inc</b> . Swiss National Supercomputing Centre (CSCS) <b>Switzerland</b>	361,760	19,590.0	25,326.3	2,272
4	<b>Titan</b> - Cray XK7, Opteron 6274 16C 2.200GHz, Cray Gemini interconnect, NVIDIA K20x , <b>Cray Inc.</b> DOE/SC/Oak Ridge National Laboratory <b>United States</b>	560,640	17,590.0	27,112.5	8,209
5	<b>Sequoia</b> - BlueGene/Q, Power BQC 16C 1.60 GHz, Custom , IBM DOE/NNSA/LLNL United States	1,572,864	17,173.2	20,132.7	7,890

#### Sunway TaihuLight System

![](_page_24_Picture_1.jpeg)

#### Figure 4: Overview of the Sunway TaihuLight System

#### Cores per socket

![](_page_25_Figure_1.jpeg)

#### Accelerators

![](_page_26_Figure_1.jpeg)

#### The Future

![](_page_27_Figure_1.jpeg)

#### The END