# Introduction to Parallel and High Performance Computing

Oguz Kaya

Maˆıtre de Confˊerences
Universit'e Paris−Saclay and the LRI ParSys team, Orsay, France

## Outline

7 Around the development of parallel programs

# Outline

Around the development of parallel programs

# Objectives

## Objectives

- Getting to know parallel computing

université
PARIS-SACLAY

## Objectives

- Getting to know parallel computing
- Discover the applications that need computing power

## Objectives

- Getting to know parallel computing
- Discover the applications that need computing power Explore the
- modern architecture of a parallel computer

# Objectives

- Getting to know parallel computing
- Discover the applications that need computing power Explore the
- modern architecture of a parallel computer Introduce the main
- models of parallel programming

## Objectives

- Getting to know parallel computing
- Discover the applications that need computing power Explore the
- modern architecture of a parallel computer Introduce the main
- parallel programming models Introduce the basic concepts and
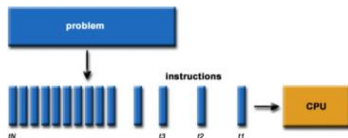- examples

## Objectives

- Getting to know parallel computing
- Discover the applications that need computing power Explore the
- modern architecture of a parallel computer Introduce the main
- parallel programming models Introduce the basic concepts and
- examples
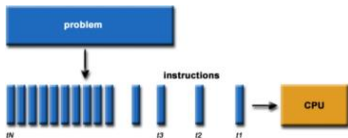- Provide some tips for developing effective parallel programs

# Equential programming

Traditionally, software is based on **sequential** calculation:

# Equential programming

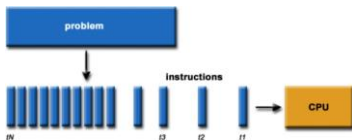Traditionally, software is based on **sequential** calculation:

- A problem is often `encountered` in instructions.

# Equential programming

Traditionally, software is based on **sequential** calculation:



- A problem is often *encountered* in instructions.
- These instructions *are* executed **sequentially** one after the other.

# Equential programming

Traditionally, software is based on **sequential** calculation:



- A problem is often *encountered* in instructions.
- These instructions *are* executed **sequentially** one after the other.
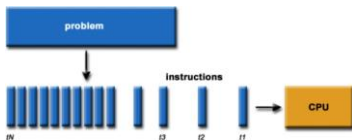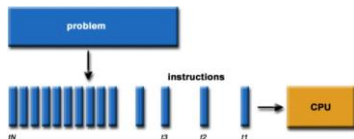- They are executed by **a single processor**.

# Equential programming

Traditionally, software is based on **sequential** calculation:



- A problem is often *encountered* in instructions.
- These instructions *are* executed **sequentially** one after the other.
- They are *executed* by **a single processor**. A g
- moment, only one instruction is ed
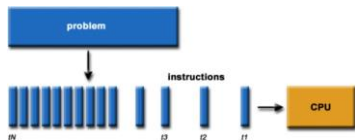
# Equential programming

Traditionally, software is based on **sequential** calculation:



- A problem is often *encountered* in instructions.
- These instructions *are* executed **sequentially** one after the other.
- They are *executed* by **a single processor**. At a given
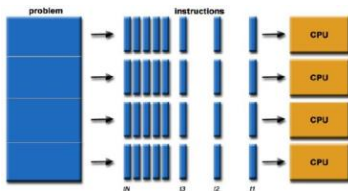- moment, only one instruction is ded
- The performance is mainly determined by **the frequency** (Hz) of the processor.
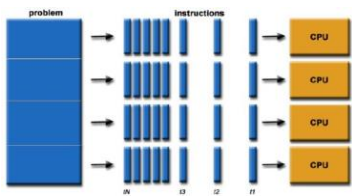
# Parallel programming

**Parallel programming** allows the use of several computing resources to solve a problem:

# Parallel programming

**Parallel programming** allows the use of several computing resources to solve a problem:

- A problem is divided into parts that can être launched in

université
PARIS-SACLAY

## Parallel programming

**Parallel programming** allows the use of several computing resources to solve a problem:

- A problem is divided into parts that can être launched
- Each part is still under instruction.

université
PARIS-SACLAY

## Parallel programming

**Parallel programming** allows the use of several computing resources to solve a problem:



- A problem is divided into parts that can ˆetre launched in
- Each part is still under instruction.
- The instructions of each part are *executed* **in parallel** using several processors.

## Parallel programming

**Parallel programming** allows the use of several computing resources to solve a problem:



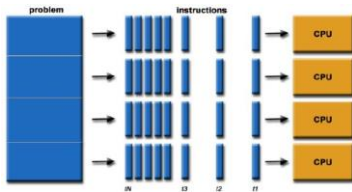- A problem is divided into parts that can ˆetre launched in
- Each part is still under instruction.
- The instructions of each part are *executed* **in parallel** using several processors.
- The performance is determined by:

université
PARIS-SACLAY

## Parallel programming

**Parallel programming** allows the use of several computing resources to solve a problem:

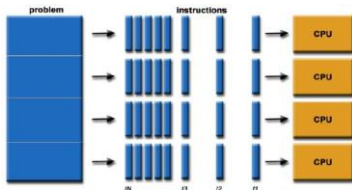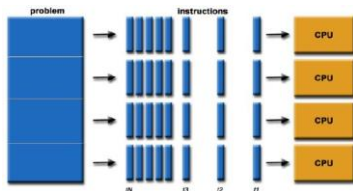- A problem is divided into parts that can être launched in
- Each part is still under instruction.
- The instructions of each part are *executed* **in parallel** using several processors.
- The performance is determined by:
  - The frequency of the processor

## Parallel programming

**Parallel programming** allows the use of several computing resources to solve a problem:



- A problem is divided into parts that can ˆetre launched in
- Each part is still under instruction.
- The instructions of each part are *executed* **in parallel** using several processors.
- The performance is determined by:
  - The frequency of the
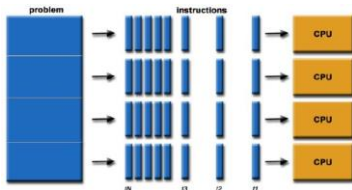  - processor The number of processors

université
PARIS-SACLAY

## Parallel programming

**Parallel programming** allows the use of several computing resources to solve a problem:
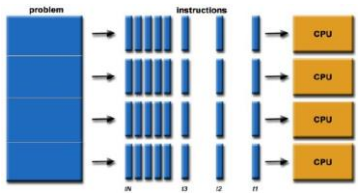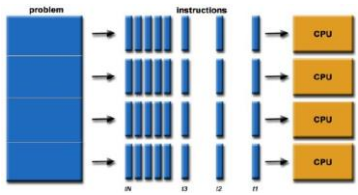


- A problem is divided into parts that can être launched in
- Each part is still under instruction.
- The instructions of each part are *executed* **in parallel** using several processors.
- The performance is determined by:
  - The frequency of the
  - processor The number of processors
  - The degree of parallelism of the problem

université
PARIS-SACLAY

# Outline
## Parallel programming

# Applications of parallel computing

Many time-consuming applications in various fields:

## Applications of parallel computing

Many time-consuming applications in various fields:



- Scientific computing: Simulations in physics, chemistry, biology, …

université
PARIS-SACLAY

**Applications of parallel computing**

Many time-consuming applications in various fields:



- Scientific computing: Simulations in physics, chemistry, biology, …
- Neural network processing

université
PARIS-SACLAY
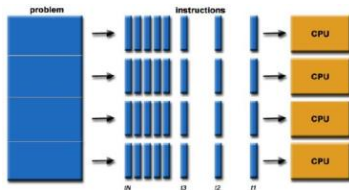
Introduction **Parallel computing, why?** Concepts of Parallel Computing Types of Parallelism Parallel Architecture Parallel Programming Parallel Program

○○○○　　　○●○○　　　　　○○○○○○○○○○○○○　　　○○○○○　　　○○○○○○○○○○○　　　○○○○○○○○○○○　　　○○○○○○○○○○

## Applications of parallel computing

Many time-consuming applications in various fields:



- Scientific computing: Simulations in physics, chemistry, biology, ...
- Neural network processing Graphics
- (rendering, video games, etc.)

## Applications of parallel computing

Numerous time-consuming applications in various fields:



- Scientific computing: Simulations in physics, chemistry, biology, …
- Neural network processing Graphics
- (rendering, video games, etc.)
- Operating systems (Linux, Android, etc.)

## Applications of parallel computing

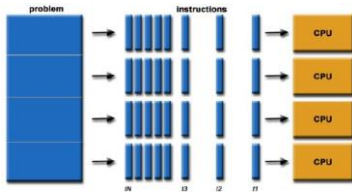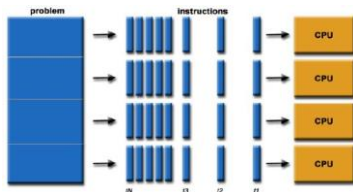Many time-consuming applications in various fields:



- Scientific computing: Simulations in physics, chemistry, biology, …
- Neural network processing Graphics
- (rendering, video games, etc.)
- Operating systems (Linux, Android, etc.) and
- many others...

université
PARIS-SACLAY

# Maximum frequency of a CPU in 2002

What is the maximum frequency of this CPU in 2002?



Intel Pentium 4, Northwood

# Maximum frequency of a CPU in 2002

What is the maximum frequency of this CPU in 2002?



- 3.06 GHz

Intel Pentium 4, Northwood

# Maximum frequency of a CPU

What is the maximum frequency of a CPU in 2020?



Intel Core i9-10900K, Comet Lake

# Maximum frequency of a CPU

What is the maximum frequency of a CPU in 2020?



Intel Core i9-10900K, Comet Lake

- 5.3 GHz

Introduction **Parallel computing, why?** Concepts of Parallel Computing Types of Parallelism Parallel Architecture Parallel Programming Parallel Program

○○○○ ○○○● ○○○○○○○○○○○○○ ○○○○○ ○○○○○○○○○○○ ○○○○○○○○○○○ ○○○○○○○○○○

# Maximum frequency of a CPU

What is the maximum frequency of a CPU in 2020?



Intel Core i9-10900K, Comet Lake

- 5.3 GHz
- The performance crˆete? 5.3 Gflops/s?

# Maximum frequency of a CPU

What is the maximum frequency of a CPU in 2020?



Intel Core i9-10900K, Comet Lake

- 5.3 GHz
- The performance crˆete? 5.3 Gflops/s?
- No! $\approx 1.7$Tflops/s (1700Gflops/s) How?

# Maximum frequency of a CPU

What is the maximum frequency of a CPU in 2020?



Intel Core i9-10900K, Comet Lake

- 5.3 GHz
- The performance crˆete? 5.3 Gflops/s?
- No! $\approx 1.7$Tflops/s (1700Gflops/s) How?
- Thanks to deadans parallelism (10 cores, each with 2 AVX256 virtual units)

universite
PARIS-SACLAY

# Maximum frequency of a CPU

# Moore's Law

The number of transistors in the `integrated` circuits doubles every 2 years.



Moore's Law

# Moore's Law

The number of transistors in the `integrated` circuits doubles every 2 years.



What is it ?

Moore's Law

Introduction Parallel computing, why? **Concepts of Parallel Computing** Types of Parallelism Parallel Architecture Parallel Programming Parallel Program

oooo          oooo                    o●oooooooooooo        ooooo           oooooooooooo          ooooooooooo          oooooooooo

# Moore's Law

The number of transistors in the `integrated` circuits doubles every 2 years.



Moore's Law

- What is it ?
- More transistors → more performance potential
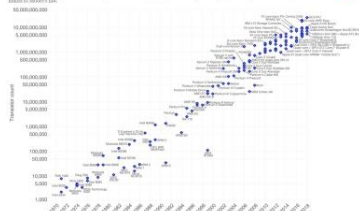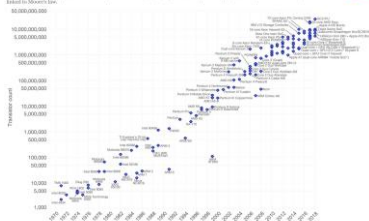
# Moore's Law

The number of transistors in the `integrated` circuits doubles every 2 years.



Moore's Law – The number of transistors on integrated circuit chips (1971-2018)

Moore's Law

- What is it ?
- More transistors → more performance potential
- Moore's law is still valid today (although a bit slowed down).

## Dennard Scaling

If the engraving size is reduced by -30% (0.7x) in each generation, this gives



40 Years of Microprocessor Trend Data

Transistors (thousands)

Single-Thread Performance (SpecINT x 10$^3$)

Frequency (MHz)

Typical Power (Watts)

Number of Logical Cores

Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

Dennard Scaling

## Dennard Scaling

If the engraving size is reduced by -30% (0.7x) in each generation, this gives

A decrease in the circuit area of 50%.



Dennard Scaling

## Dennard Scaling

If the engraving size is reduced by -30% (0.7x) in each generation, this gives



40 Years of Microprocessor Trend Data

Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

Dennard Scaling

50% decrease in circuit area 30%
- decrease in latency
- 

université
PARIS-SACLAY

## Dennard Scaling

If the engraving size is reduced by -30% (0.7x) in each generation, this gives



40 Years of Microprocessor Trend Data

Transistors (thousands)

Single-Thread Performance (SpecINT x 10³)

Frequency (MHz)

Typical Power (Watts)

Number of Logical Cores

Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

Dennard Scaling

- 50% decrease in circuit area 30%
- decrease in latency
- Consequently, a 40% increase in frèquence ($10/7 \approx 1.4$).

## Dennard Scaling

If the engraving size is reduced by -30% (0.7x) in each generation, this gives



40 Years of Microprocessor Trend Data

Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

- 50% decrease in circuit area 30%
- decrease in latency
- Consequently, a 40% increase in frequence
  ($10/7 \approx 1.4$).
  - Not possible anymore because of quantum effects!

université
PARIS-SACLAY

## Dennard Scaling

If the engraving size is reduced by -30% (0.7x) in each generation, this gives



40 Years of Microprocessor Trend Data

Transistors (thousands)

Single-Thread Performance (SpecINT x 10³)

Frequency (MHz)

Typical Power (Watts)

Number of Logical Cores

Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

Dennard Scaling
Dennard Scaling

- 50% decrease in circuit area 30%
- decrease in latency
- Consequently, a 40% increase in frequence $(10/7 \approx 1.4)$.
  - Not possible anymore because of quantum effects!
  - Use either multi-core or large-core.

université
PARIS-SACLAY

## Dennard Scaling

If the engraving size is reduced by -30% (0.7x) in each generation, this gives



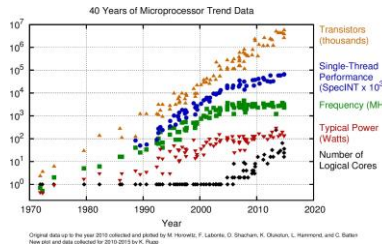40 Years of Microprocessor Trend Data

Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

- 50% decrease in circuit area 30%
- decrease in latency
- Consequently, a 40% increase in frequence
  $(10/7 \approx 1.4)$.
  - Not possible anymore because of quantum effects!
  - Use either multi-core or large-core.

- *energy ~ frequency2*

université
PARIS-SACLAY

## Acceleration and efficiency

If the engraving size is reduced by -30% (0.7x) in each generation, this gives



40 Years of Microprocessor Trend Data

Transistors (thousands)

Single-Thread Performance (SpecINT x 10³)

Frequency (MHz)

Typical Power (Watts)

Number of Logical Cores

Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

Dennard Scaling

Cray-1 Supercomputer (1975) with peak performance of ≈ 160Mflops/s

- 50% decrease in circuit area 30%
- decrease in latency
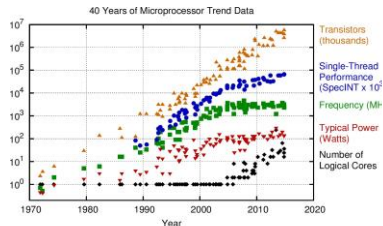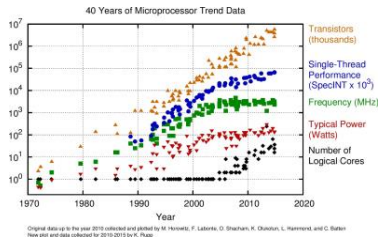- Consequently, a 40% increase in frèquence ($10/7 \approx 1.4$).
    - Not possible anymore because of quantum effects!
    - Use either multi-core or large-core.

- *energy ~ frequency2*
    - Use more cores at a lower frequency
    
      → more performance in the mˆeme fenˆetre energetics

universite PARIS-SACLAY
universite PARIS-SACLAY

# Dennard Scaling

How much faster does a parallel program run? How well are computing resources

# Acceleration and efficiency

How much faster does a parallel program run? How well are computing resources



Cray-1 Supercomputer (1975) with
peak performance of ≈ 160Mflops/s

- Sequential execution time: $T(1)$

# Acceleration and efficiency

How much faster does a parallel program run? How well are computing resources



Cray-1 Supercomputer (1975) with
peak performance of $\approx$ 160Mflops/s

- Sequential execution time: $T(1)$
- The execution time on $N$ processors: $T(N)$

# Acceleration and efficiency

How much faster does a parallel program run? How well are computing resources



Cray-1 Supercomputer (1975) with peak performance of ≈ 160Mflops/s

- Sequential execution time: $T(1)$
- The execution time on $N$ processors: $T(N)$
- **Accèlèration**: $S(N) = T(1)/T(N)$
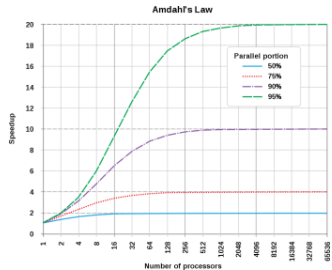
# Acceleration and efficiency

How much faster does a parallel program run? How well are computing resources ？



- Sequential execution time: $T(1)$
- The execution time on $N$ processors: $T(N)$
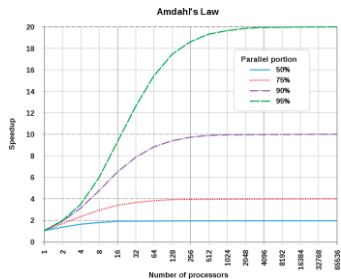- **Acceleration**: $S(N) = T(1)/T(N)$
- **Efficiency**: $E(N) = S(N)/N$

Cray-1 Supercomputer (1975) with peak performance of $\approx$ 160Mflops/s

universite
PARIS-SACLAY

# Amdahl's Law



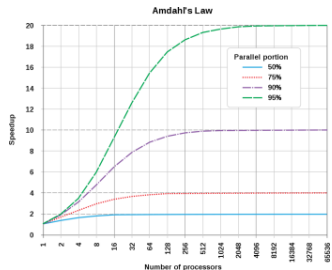Amdahl's Law

# Amdahl's Law

How much faster could a program run? What is the limit?



Amdahl's Law

- Introduced in 1967 by Gene Amdahl

# Amdahl's Law



Amdahl's Law

Introduced in 1967 by Gene Amdahl

- Let $s$ and $p$ be the **equential** fraction **s**
- and
  **parallelizable** of a program ($s + p = 1.0$).

université
PARIS-SACLAY

# Amdahl's Law
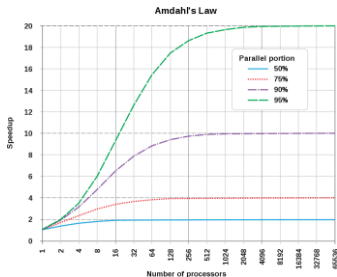
How much faster could a program run? What is the limit?



- Introduced in 1967 by Gene Amdahl
- Let $s$ and $p$ be the **equential** fraction **s** and **parallelizable** of a program ($s + p = 1.0$).
- The achievable acceleration is defined by
$$A(N) = \frac{T(1)}{T(N)} \leq \frac{T(1)}{pT(1)/N + sT(1)} = \frac{1}{p/N + s}$$

# Amdahl's Law

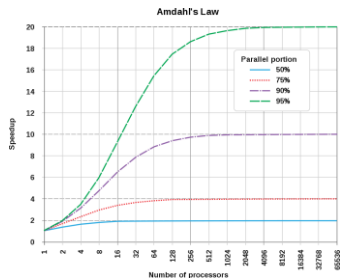How much faster could a program run? What is the limit?



Amdahl's Law

- Introduced in 1967 by Gene Amdahl
- Let $s$ and $p$ be the **equential** fraction **s** and **parallelizable** of a program ($s + p = 1.0$).
- The achievable acceleration is defined by

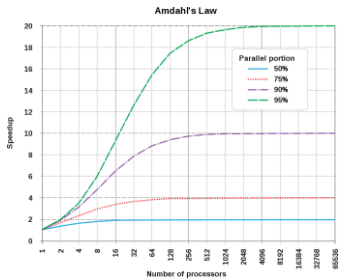$$A(N) = \frac{T(1)}{T(N)} \leq \frac{T(1)}{pT(1)/N + sT(1)} = \frac{1}{p/N + s}$$

Amdahl's Law

# Amdahl's Law

**Strong scaling:**

Parallel scalability for a fixed size problem

# Amdahl's Law

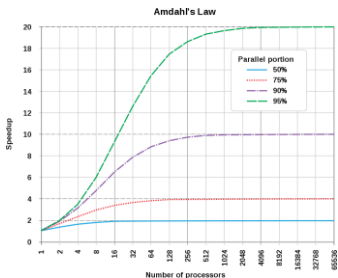How much faster could a program run? What is the limit?



Amdahl's Law

- Introduced in 1967 by Gene Amdahl
- Let $s$ and $p$ be the **equential** fraction **s** and **parallelizable** of a program ($s + p = 1.0$).
- The achievable acceleration is defined by

$$A(N) = \frac{T(1)}{T(N)} \leq \frac{T(1)}{pT(1)/N + sT(1)} = \frac{1}{p/N + s}$$

- **Strong scaling:** Parallel scalability for a fixed size problem

# Amdahl's Law

- Hope for parallel computing?

# Amdahl's Law

How much faster could a program run? What is the limit?



Amdahl's Law

- Introduced in 1967 by Gene Amdahl
- Let $s$ and $p$ be the **equential** fraction **s** and **parallelizable** of a program ($s + p = 1.0$).
- The achievable acceleration is defined by

$$A(N) = \frac{T(1)}{T(N)} \leq \frac{T(1)}{pT(1)/N + sT(1)} = \frac{1}{p/N + s}$$

Amdahl's Law

# Amdahl's Law

**Strong scaling:**
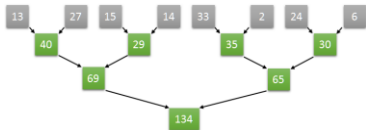
Parallel scalability for a fixed size problem

Hope for parallel computing?

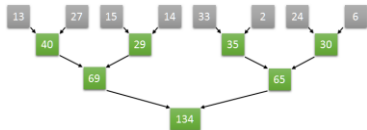Typically, *S* decreases with the size of the problem.

universite
PARIS-SACLAY

# Example: Sum of an array

$$\sum_{i=0}^{L-N-1} A[i]$$



Sum in parallel

## Example: Sum of an array

$$\sum_{i=0}^{N-1} A[i]$$

- $N - 1$ additions



Sum in parallel

# **Example: Sum of an array**



Sum in parallel

$$\sum_{i=0}^{N-1} A[i]$$

- $N$ - 1 additions
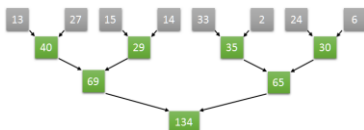- For $N = 8$, 7 additions of which 3 are equal,
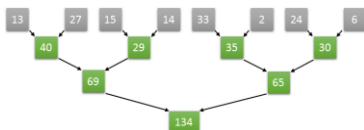
$S \leq 7/3 = 2.34$

## Example: Sum of an array



Sum in parallel

$$\sum_{i=0}^{N-1} A[i]$$

- $N$ - 1 additions
- For $N = 8$, 7 additions of which 3 are equal,

  $S \leq 7/3 = 2.34$
- For $N = 16$, 15 additions of which 4 are equential, $S \leq 15/4 = 3.75$

# **Example: Sum of an array**



Sum in parallel

$$\overset{N-1}{\underset{i=0}{L-}} A[i ]$$

- $N$ - 1 additions
- For $N = 8$, 7 additions of which 3 are equal, $S \leq 7/3 = 2.34$
- For $N = 16$, 15 additions of which 4 are equential, $S \leq 15/4 = 3.75$
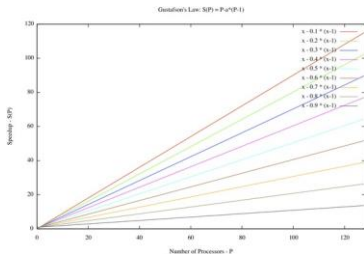- In the general case, $N$ - 1 additions of which $\log N$ equals, $S \leq (N - 1)/\log N$

université
PARIS-SACLAY

# Gustafson's law

How much acceleration could be obtained if the problem size increases with the number of processors?



Gustafson's law

# Gustafson's law

How much acceleration could be obtained if the problem size increases with the number of processors?



Gustafson's law

- $S(N) = N - (1 - N)s$

## Gustafson's law

How much acceleration could be obtained if the problem size increases with the number of processors?



Gustafson's law

- $S(N) = N - (1 - N)s$
- Amdalh: It could ˆetre difficult accélér un calculation with more processors

Introduction Parallel computing, why? **Concepts of Parallel Computing** Types of Parallelism Parallel Architecture Parallel Programming Parallel Program

0000    0000    000000●000000    00000    000000000000    0000000000    0000000000

# Gustafson's law

How much acceleration could be obtained if the problem size increases with the number of processors?



Gustafson's law

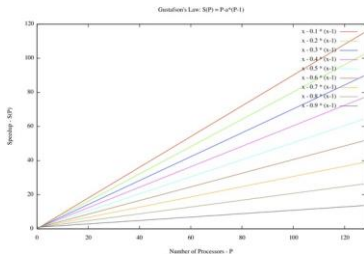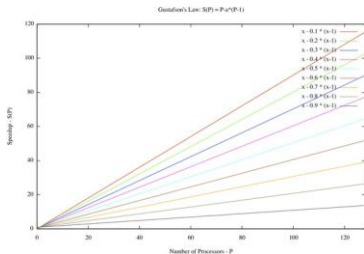- $S(N) = N - (1 - N)s$
- Amdalh: It could être difficult accèler un calculation with more processors
- Gustafson: It is possible to do more computation as fast with more processors

## Gustafson's law

How much acceleration could be obtained if the problem size increases with the number of processors?



Gustafson's law

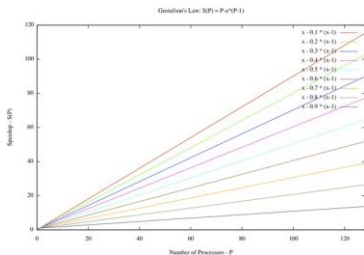- $S(N) = N - (1 - N)s$
- Amdalh: It could être difficult accéler un calculation with more processors
- Gustafson: It is possible to do more computation as fast with more processors
- Which one is more relevant in practice?

## Gustafson's law

How much acceleration could be obtained if the problem size increases with the number of processors?



Gustafson's law

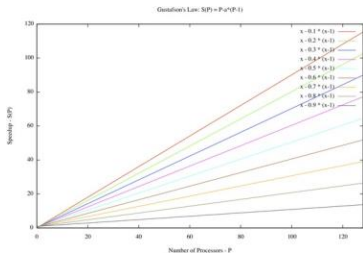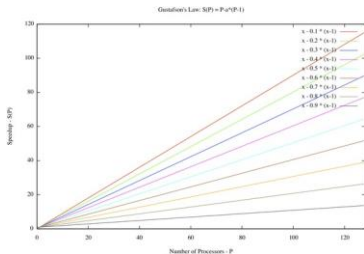- $S(N) = N - (1 - N)s$
- Amdalh: It could ˆetre difficult accéler un calculation with more processors
- Gustafson: It is possible to do more computation as fast with more processors
- Which one is more relevant in practice?
- **Weak scaling**Scalability parallel to the The problem size is increasing.

université
PARIS-SACLAY

# Flynn's classification

A $2 \times 2$ matrix to classify parallel machines.



| SISD | SIMD |
|------|------|
| Single Instruction stream<br>Single Data stream | Single Instruction stream<br>Multiple Data stream |
| MISD | MIMD |
| Multiple Instruction stream<br>Single Data stream | Multiple Instruction stream<br>Multiple Data stream |

Flynn's classification

# Flynn's classification

A $2 \times 2$ matrix to classify parallel machines.



| SISD | SIMD |
|------|------|
| Single Instruction stream Single Data stream | Single Instruction stream Multiple Data stream |
| MISD | MIMD |
| Multiple Instruction stream Single Data stream | Multiple Instruction stream Multiple Data stream |

Flynn's classification

- Extended classification, used since 1966

# Flynn's classification

A $2 \times 2$ matrix to classify parallel machines.



- Extended classification, used since 1966 The X
- axis: single, multiple instruction

Flynn's classification

# Flynn's classification

A $2 \times 2$ matrix to classify parallel machines.



| SISD | SIMD |
|------|------|
| Single Instruction stream Single Data stream | Single Instruction stream Multiple Data stream |
| MISD | MIMD |
| Multiple Instruction stream Single Data stream | Multiple Instruction stream Multiple Data stream |

Flynn's classification

- Extended classification, used since 1966 X
- axis: single (simple), multiple instruction Y
- axis: single (simple), multiple data

université
PARIS-SACLAY

Introduction Parallel computing, why? **Concepts of Parallel Computing** Types of Parallelism Parallel Architecture Parallel Programming Parallel Program

0000    0000    000000000●0000    00000    000000000000    00000000000    0000000000

# Single Instruction, Single Data (SISD)



SISD Processor

Introduction Parallel computing, why? **Concepts of Parallel Computing** Types of Parallelism Parallel Architecture Parallel Programming Parallel Program

0000 0000 00000000●0000 00000 000000000000 00000000000 0000000000

## Single Instruction, Single Data (SISD)

| load A |
|--------|
| load B |
| C = A + B |
| store C |
| A = B * 2 |
| store A |

time

SISD Processor

- A completely sequential **m**

## Single Instruction, Single Data (SISD)



| load A |
| load B |
| C = A + B |
| store C |
| A = B * 2 |
| store A |

time

SISD Processor

- A completely sequential m
- Single" instructionOnly one instruction is executed each clock cycle
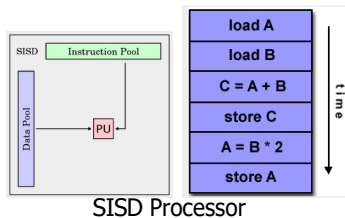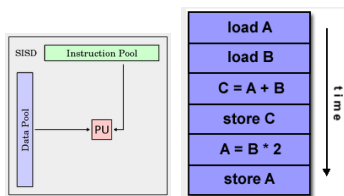
## Single Instruction, Single Data (SISD)



SISD Processor

- A completely sequential m
- Single" instructionOnly one instruction is executed each clock cycle
- Single data: Only one data stream is used

université
PARIS-SACLAY

## Single Instruction, Single Data (SISD)



SISD Processor

- A completely sequential m
- Single" instructionOnly one instruction is executed each clock cycle
- Single data: Only one data stream is used
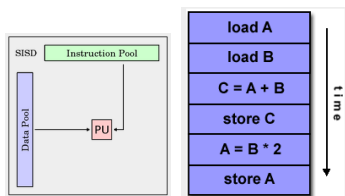- Execution of eterminist

## Single Instruction, Single Data (SISD)



SISD Processor

- A completely sequential m**a**
- Single" instruction**O**nly one instruction is `executed` **ea**ch clock cycle
- Single data: Only one `data` stream is used
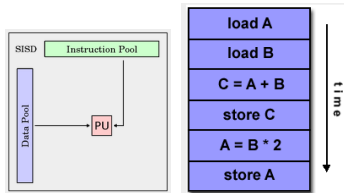- Execution of eterminist
- The oldest computers (i.e., architecture courses)

# Single Instruction, Multiple Data (SIMD)



SISD Processor

# Single Instruction, Multiple Data (SIMD)



SISD Processor

- A completely sequential m**a**

# Single Instruction, Multiple Data (SIMD)



SISD Processor

- A completely sequential m**a**
- Single" instruction Only one instruction is executed **e**ach clock cycle

# Single Instruction, Multiple Data (SIMD)
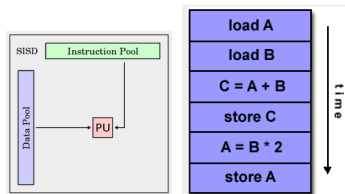


SISD Processor

- A completely sequential m**a**
- Single" instructionOnly one instruction is executed **e**ach clock cycle
- Single data: Only one data stream is used

université
PARIS-SACLAY

# Single Instruction, Multiple Data (SIMD)



SISD Processor

- A completely sequential m
- Single" instructionOnly one instruction is executed each clock cycle
- Single data: Only one data stream is used
- Execution of eterminist

# Single Instruction, Multiple Data (SIMD)



SISD Processor

- A completely sequential m**a**
- Single" instructionOnly one instruction is `executed` **a**ach clock cycle
- Single data: Only one `data` stream is used
- Execution of eterminist
- The oldest computers (i.e., architecture courses)

# Single Instruction, Multiple Data (SIMD)

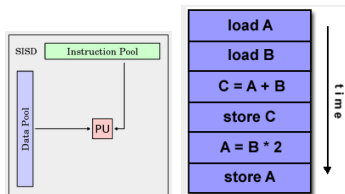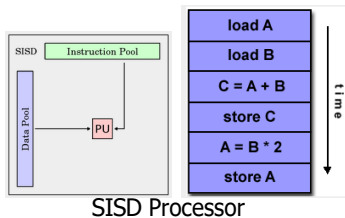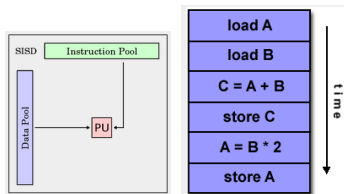- An example of a parallel machine
- Instruction "Single" : All the computing *units* execute the *same* instruction at each clock cycle
- Multiple data: Each computing unit can operate on a different data Compatible for quite
- regular problems like image processing

# Multiple Instruction, Single Data (MISD)

- An example of a parallel machine
- Multiple Instruction: Each computing unit operates on independent data via different instruction flows
- Single" data: A single data stream feeds several computing units Some
- examples of use are :
  - several frequency filters operating on a single signal
  - several cryptographic algorithms trying to decrypt a single coded message

# Multiple Instruction, Multiple Data (MIMD)

- An example of a parallel machine
- Multiple Instructioneach computing *unit* can execute a different instruction flow
- Multiple *data*: each processor can *operate* on a *different data* stream Execution can
- be synchronous or asynchronous, deterministic or non-deterministic The majority
- of modern parallel machines belong to this *category* Several MIMD architectures
- also include SIMD components



universite
PARIS-SACLAY

## Outline

# Parallelism at bit level

This is the word-size of the processor.



Intel C4004, 4-bit processor (1971)

# Parallelism at bit level

This is the word-size of the processor.



Intel C4004, 4-bit processor (1971)

- Relevant for the `years` 1970..1986.

# Parallelism at bit level

This is the word-size of the processor.



Intel C4004, 4-bit processor (1971)

- Relevant for the `years` 1970..1986.
- The word size: 4-bit $\rightarrow$ 8-bit $\rightarrow$ 16-bit $\rightarrow$ 32-bit

# Parallelism at bit level

This is the word-size of the processor.



Intel C4004, 4-bit processor (1971)

- Relevant for the `years` 1970..1986.
- The word size: 4-bit $\rightarrow$ 8-bit $\rightarrow$ 16-bit $\rightarrow$ 32-bit
- Convergée in 64-bit today.

# Parallelism at the instruction level (ILP)



Ex'ecution without pipeline



Superscalar execution with pipeline

# Parallelism at the instruction level (ILP)

Allows to execute independent instructions simultaneously.



- a = b + c; c = d + e;
  f = a + c; g = c - a;

Ex'ecution without pipeline

Superscalar execution with pipeline

# Parallelism at the instruction level (ILP)



Ex'ecution without pipeline



Superscalar execution with pipeline

- **a = b + c; c = d + e;**
  **f = a + c; g = c - a;**
- Possibilité 1: Pipeline. Différents étapes de deux instructions indépendantes peuvent être exécutées dans le pipeline en même temps.

# Parallelism at the instruction level (ILP)

Allows to execute independent instructions simultaneously.



Ex'ecution without pipeline



Superscalar execution with pipeline

- **a = b + c; c = d + e;**
- **f = a + c; g = c - a;**

- Possibility 1: Pipeline. Different stages of two ind'ependant instructions can ^etre ex'ecut'ees in the pipeline en m^eme temps.

- Possibility 2: Superscalar execution: Independent instructions can use several execution units (i.e. ALU) in a **superscalar** processor.

université
PARIS-SACLAY

# Data Parallelism

Allows to execute independent instructions simultaneously.

Ex'ecution without pipeline

Superscalar execution with pipeline

- **a = b + c; c = d + e;**
  **f = a + c; g = c - a;**

- Possibility 1: Pipeline. Different stages of two ind'ependant instructions can ^etre exècutèes in the pipeline en m^eme temps.

- Possibility 2: Superscalar execution: Independent instructions can use several execution units (i.e. ALU) in a **superscalar** processor.

- Most modern processors are superscalar.

université
PARIS-SACLAY

# Parallelism at the instruction level (ILP)

Carry out the mêmes independent operations of the data tables.



Addition of the two tables

## Data Parallelism

Carry out the mˆemes ind'ependent operations of the data tables.



- Suitable for the SIMD paradigm.

Addition of the two tables

# Data Parallelism

Carry out the mêmes ind'ependent operations of the data tables.



Addition of the two tables

- Suitable for the SIMD paradigm.
- Well adapted hardware for efficient execution (CPU and GPU vector unit)

université
PARIS-SACLAY

Introduction Parallel computing, why? Concepts of Parallel Computing **Types of Parallelism** Parallel Architecture Parallel Programming Parallel Program

oooo        oooo        oooooooooooo        ooooo        ooooooooooo        ooooooooooo        oooooooooo

## Data Parallelism

Carry out the mˆemes ind'ependent operations of the data tables.



index   0   1   2   3   4   5   6   7

Array A | 1 | 2 | 3 | 1 | 4 | 1 | 6 | 7 |

**+**

B | 1 | 2 | 3 | 1 | 5 | 2 | 6 | 1 |

**=**

C | 2 | 4 | 6 | 2 | 9 | 3 | 12 | 8 |

Addition of the two tables

- Suitable for the SIMD paradigm.
- Well adapted hardware for efficient execution (CPU and GPU vector unit)
- Extensive use in scientific computing (matrices, vectors), graphics (rendering), image and signal processing (filters).

# Task Parallelism



Task parallelism

## Task Parallelism

It is the execution of ind'ependent `tasks` in a program.

- Examples of ind'ependent `tasks`?



Task parallelism

## Task Parallelism



Task parallelism

- Examples of ind'ependent `tasks`?
  - Calls to functions: **a = f(x); b = g(y);**

université
PARIS-SACLAY

# Task Parallelism

It is the execution of ind'ependent `tasks` in a program.



Task parallelism

- Examples of ind'ependent `tasks`?
  - Appels aux fonctions: **a = f(x); b =**
  - **g(y);** Blocs de code ind'ependents dans une mˆeme fonction.

université
PARIS-SACLAY

## Task Parallelism

It is the execution of ind'ependent `tasks` in a program.



Task parallelism

- Examples of ind'ependent `tasks`?
  - Appels aux fonctions: **a = f(x); b = g(y);** Blocs de code ind'ependents dans une mˆeme fonction.
  - Multiple execution of a mˆeme program with different parameters (i.e. simulation).

## Task Parallelism

It is the execution of ind'ependent `tasks` in a program.



Task parallelism

- Examples of ind'ependent `tasks`?
  - Appels aux fonctions: **a = f(x); b =**
  - **g(y);** Blocs de code ind'ependents dans une mˆeme fonction.
  - Multiple execution of a mˆeme program with different parameters (i.e. simulation).
- Each `task` can be executed on a `separate` calculation unit.

université
PARIS-SACLAY

## Task Parallelism

It is the execution of ind'ependent `tasks` in a program.



Task parallelism

- Examples of ind'ependent `tasks`?
    - Appels aux fonctions: **a = f(x); b =**
    - **g(y);** Blocs de code ind'ependents dans une
      mˆeme fonction.
    - Multiple execution of a mˆeme program with
      different parameters (i.e. simulation).
- Each `task` can be executed on a `separate`
  calculation unit.
- It is necessary to respect the d'ependences if there are
  any.

université
PARIS-SACLAY

## Task Parallelism

It is the execution of ind'ependent `tasks` in a program.



- Examples of ind'ependent `tasks`?
  - Appels aux fonctions: **a = f(x); b =**
  - **g(y);** Blocs de code ind'ependents dans une mˆeme fonction.
  - Multiple execution of a mˆeme program with different parameters (i.e. simulation).
- Each `task` can be executed on a `separate` calculation unit.
- It is necessary to respect the d'ependences if there are any.
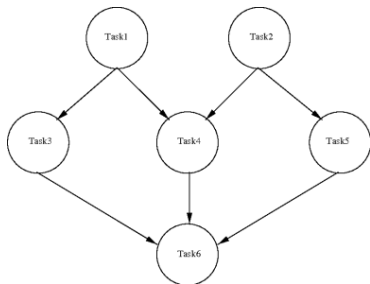
## Task Parallelism

- Load balancing?

Task parallelism

## Task Parallelism

It is the execution of ind'ependent `tasks` in a program.



- Exemples des tˆaches ind'ependentes?
  - Appels aux fonctions: **a = f(x); b =**
  - **g(y);** Blocs de code ind'ependents dans une mˆeme fonction.
  - Multiple execution of a mˆeme program with different parameters (i.e. simulation).
- Each `task` can be executed on a `separate` calculation unit.
- It is necessary to respect the d'ependences if there are any.
  - Load balancing?
  - 

université
PARIS-SACLAY

# Task Parallelism

Task parallelism

A whole field of research (scheduling).

Task parallelism

# Outline
## Task Parallelism

# CPU

"Central Processing Unit, a `general` computing unit consisting 6



Zen 2 architecture

Introduction Parallel computing, why? Concepts of Parallel Computing Types of Parallelism **Parallel Architecture** Parallel Programming Parallel Program

0000  0000  000000000000  00000  0●0000000000  00000000000  0000000000

# CPU

"Central Processing Unit, a `general` computing unit consisting o

- Several execution units (hearts)



Zen 2 architecture

## CPU

"Central Processing Unit, a `general` computing unit consisting 6



Zen 2 architecture

- Several execution units (hearts)
- Multiple memory levels (registers, L1, L2, L3, RAM)

université
PARIS-SACLAY

## CPU

"Central Processing Unit, a `general` computing unit consisting fo



Zen 2 architecture

- Several execution units (hearts)
- Multiple memory levels (registers, L1, L2, L3, RAM)
- Several execution ports in each core (ALUs, vector units)

université
PARIS-SACLAY

# **CPU**

"Central Processing Unit, a `general` computing unit consisting 6



Zen 2 architecture

- Several execution units (hearts)
- Multiple memory levels (registers, L1, L2, L3, RAM)
- Several execution ports in each core (ALUs, vector units)
- Vector units (AVX2, AVX512, Arm Neon, ...)

## CPU

"Central Processing Unit, a general computing unit consisting of



Zen 2 architecture

- Several execution units (hearts)
- Multiple memory levels (registers, L1, L2, L3, RAM)
- Several execution ports in each core (ALUs, vector units)
- Vector units (AVX2, AVX512, Arm Neon, ...)
- Execution of some threads (1-4) simultaneously

université
PARIS-SACLAY

# CPU

"Central Processing Unit, a `general` computing unit consisting of



Zen 2 architecture

- Several execution units (hearts)
- Multiple memory levels (registers, L1, L2, L3, RAM)
- Several execution ports in each core (ALUs, vector units)
- Vector units (AVX2, AVX512, Arm Neon, ...)
- Execution of some threads (1-4) simultaneously
- Able to exploit parallelism at the instruction level (micro-op buffer, instruction renaming, register renaming, ...)

## CPU

"Central Processing Unit, a `general` computing unit consisting of



Zen 2 architecture

- Several execution units (hearts)
- Multiple memory levels (registers, L1, L2, L3, RAM)
- Several execution ports in each core (ALUs, vector units)
- Vector units (AVX2, AVX512, Arm Neon, ...)
- Execution of some threads (1-4) simultaneously
- Able to exploit parallelism at the instruction level (micro-op buffer, instruction renaming, register renaming, ...)
- A considerable part of the circuit is dedicated to ILP +

## GPU cache

"Graphical Processing Unit", a specific virtual computing unit consisting fo

# CPU



The Nvidia Ampere architecture

# GPU

"Graphical Processing Unit", a specific virtual computing unit consisting 6

- Multiple execution units (symmetric multiprocessors (SM))



The Nvidia Ampere architecture

# GPU



The Nvidia Ampere
architecture

- Multiple execution units (symmetric multiprocessors (SM))
- Several memory levels (registers, shared memory, L1, L2, RAM)

# GPU

"Graphical Processing Unit", a specific virtual computing unit consisting of



The Nvidia Ampere architecture

- Multiple execution units (symmetric multiprocessors (SM))
- Several memory levels (registers, shared memory, L1, L2, RAM)
- Several large (2-4) vector units (16-32 floats) in each DM

## GPU

"Graphical Processing Unit", a specific virtual computing unit consisting of



- Multiple execution units (symmetric multiprocessors (SM))
- Several memory levels (registers, shared memory, L1, L2, RAM)
- Several large (2-4) vector units (16-32 floats) in each DM
- Execution of hundreds of simultaneous threads

The Nvidia Ampere architecture

## GPU

"Graphical Processing Unit", a specific virtual computing unit consisting f



The Nvidia Ampere architecture

- Multiple execution units (symmetric multiprocessors (SM))
- Several memory levels (registers, shared memory, L1, L2, RAM)
- Several large (2-4) vector units (16-32 floats) in each DM
- Execution of hundreds of simultaneous threads Large
- register array (65K)

université
PARIS-SACLAY

# GPU

"Graphical Processing Unit", a specific virtual computing unit consisting of



The Nvidia Ampere architecture

- Multiple execution units (symmetric multiprocessors (SM))
- Several memory levels (registers, shared memory, L1, L2, RAM)
- Several large (2-4) vector units (16-32 floats) in each DM
- Execution of hundreds of simultaneous threads Large
- register array (65K)
- Very fast thread exchange

université
PARIS-SACLAY

## GPU

"Graphical Processing Unit", a specific virtual computing unit consisting of



The Nvidia Ampere
architecture

- Multiple execution units (symmetric multiprocessors (SM))
- Several memory levels (registers, shared memory, L1, L2, RAM)
- Several large (2-4) vector units (16-32 floats) in each DM
- Execution of hundreds of simultaneous threads Large
- register array (65K)
- Very fast thread exchange
- Most of the circuit is devoted to vector units

université
PARIS-SACLAY

# Supercomputer / Cluster

"Graphical Processing Unit", a specific virtual computing unit consisting of



The Nvidia Ampere
architecture

- Multiple execution units (symmetric multiprocessors (SM))
- Several memory levels (registers, shared memory, L1, L2, RAM)
- Several large (2-4) vector units (16-32 floats) in each DM
- Execution of hundreds of simultaneous threads Large
- register array (65K)
- Very fast thread exchange
- Most of the circuit is devoted to vector units
- Parall'elisation takes the effort

université
PARIS-SACLAY
université
PARIS-SACLAY

# GPU

A set of machines (CPU+GPU) connected



Jolio Curie supercomputer, 300K
CPU cores, 1024 GPUs

## Supercomputer / Cluster

A set of machines (CPU+GPU) connected



- Connection by a network with a particular topology (ring, grid, torus, clique, etc.)

Jolio Curie supercomputer, 300K CPU cores, 1024 GPUs

université
PARIS-SACLAY

## Supercomputer / Cluster

A set of machines (CPU+GPU) connected



- Connection by a network with a particular topology (ring, grid, torus, clique, etc.)
- Topology-adapted communications libraries

Jolio Curie supercomputer, 300K CPU cores, 1024 GPUs

université
PARIS-SACLAY

## Supercomputer / Cluster

A set of machines (CPU+GPU) connected



Jolio Curie supercomputer, 300K
CPU cores, 1024 GPUs

- Connection by a network with a particular topology (ring, grid, torus, clique, etc.)
- Topology-adapted communications libraries
- Able to address very large problems

université
PARIS-SACLAY

## Supercomputer / Cluster

A set of machines (CPU+GPU) connected



Jolio Curie supercomputer,
300K CPU cores, 1024 GPUs

- Connection by a network with a particular topology (ring, grid, torus, clique, etc.)
- Topology-adapted communications libraries
- Able to address very large problems
- Today, at the exaflops scale ($^{1018}$ floating operations per second)

université
PARIS-SACLAY

# Shared memory: general characteristics

- All processors can access the memory as a global address space

- Plusieurs processeurs peuvent opérer de façon indépendante mais partagent les mêmes ressources moire

- Modifications by a processor in a memory area are visible to all other processors

universite
PARIS-SACLAY

# Mèmoire partagèe : Uniform Memory Access (UMA)

- Presented by SMP (Symmetric Multiprocessor) machines Identical
- processors
- Time of access the memory ègal for all processors



**Shared Memory (UMA)**

université
**PARIS-SACLAY**

# Non-Uniform Memory Access (NUMA)

- In most casesphysically built by linking two or more SMPs An SMP can directly
- access the memory of another SMP
- All processors do not have an `equal` access time for all memories Memory
- accesses the memory of another SMP are slower
- If the cache coherence is assured, we speak of cc-NUMA



Shared Memory (NUMA)

# Mémoire partagée:

1. Benefits:
   - The global address space allows a simpler programming from the point of view of memory management
   - The sharing of data between threads is fast and uniform grˆace `a proximit'e de la m'emoire du CPU

2. Inconvénients :
   - Lack of scalability between memory and CPUs: adding more CPUs will increase the use of the shared bus and for systems with a coherent cache, it will increase the effort of managing the coherence between the cache and the memory
   - It is the programmer's responsibility to make the necessary synchronizations to ensure "correct" accesses to the global memory

université
PARIS-SACLAY

# Distributed memory: general characteristics

- Distributed memory do not require a network to ensure the connection between the processors

- Each processor has its own memory and address space

- It is up to the programmer to explicitly indicate how and when the data should ˆetre communicated and when the synchronizations between the precessors should ˆetre effectu'ees

# Distributed memory: general characteristics

- The network used to transfer data between processors is very varied, it can be as simple as Ethernet

## Mémoire distribuée

1. Benefits:
   - Memory is scalable with the number of processors
   - Chaque processeur accède à mémoire rapidement sans ni interference avec les autres processeurs ni de additionnel pour maintenir une cohérence globale du cache

2. Inconvénients :
   - The programmer is responsible for several details associated with data communication between processors
   - Il peut être difficile de distribuer des structures de données conçues sur une base de mémoire globale sur cette nouvelle organisation mémoire

# Hybrid memory: general characteristics

- The most efficient machines in the world are machines that use shared and distributed memories
- Les composantes à mémoire partagée peuvent être des machines à mémoire partagée ou des GPUs (graphics processing units)
- The processors of a computing node share the same memory space nécessitent
- communications to échanger les données entre les noeuds

## Outline

# Parallel programming models

Parallel programming models exist as an abstraction on top of parallel architectures

1. Mémoire partagée:
   - **Intrinsics** SIMD instructions (Intel SSE2, ARM NEON), low level (**More details in upcoming courses**)
   - Posix Threads library
   - **OpenMP** is based on compiler directives to be played in a sequential code (**More details in the next courses**)
   - CUDA (**More details in upcoming courses**) OpenCL

2. Mémoire distribuée:
   - Library sockets, low level
   - **MPI** Message Passing Interface the standard for distributed memory architectures, the parallel code is generally very different from the serial code (**More details in the next courses**)

université
PARIS-SACLAY

# Parallel programming models

Parallel programming models exist as an abstraction on top of parallel architectures

1. Mémoire partagée:
   - **Intrinsics** SIMD instructions (Intel SSE2, ARM NEON), low level (**More details in upcoming courses**)
   - Posix Threads library
   - **OpenMP** is based on compiler directives to be played in a sequential code (**More details in the next courses**)
   - CUDA (**More details in upcoming courses**) OpenCL

2. Mémoire distribuée:
   - Library sockets, low level
   - **MPI** Message Passing Interface the standard for distributed memory architectures, the parallel code is generally very different from the serial code (**more than**

**What kind of pro        ramming used?**
**More details in the next courses**)

université
PARIS-SACLAY

**Thread model**

- IEEE POSIX Threads (PThreads)
  - A Standard UNIX API, also exists under Windows
  - More than 60 functions: *pthread create*, *pthread join*, *pthread exit*, ...

- OpenMP
  - A higher level interface, based on
    - compiler directives library
    - functions a runtime
  - An orientation towards high performance computing (HPC) applications

# OpenMP

## Fork-join execution model



## Memory model

## Message Passing Interface

- Specification and management by the MPI forum
  - The library provides a set of communication primitives: point-to-point or collective
  - C/C++ and Fortran

- A low-level programming model
  - data distribution and communications must be done manually Primitives are easy to
  - use but the development of parallel programs can be quite difficult

- Communications
  - Point to point (messages between two processors)
  - Collective (messages in groups of processors)

université
PARIS-SACLAY

# Scalar product : Sequentiel

```c
# include < s t d i o . h>
# define SIZE 256

int main () {
    double sum , a [ SIZE ] , b [ SIZE ] ;

// Initialization
    sum = 0 ;
    for ( s i z e t i = 0 ; i < SIZE ; i ++) {
        a [ i ] = i  *   0 . 5 ;
        b [ i ] = i  *   2 . 0 ;
    }

    // Computation
    for ( s i z e t i = 0 ; i < SIZE ; i ++) sum =
        sum + a [ i ]*b [ i ] ;

    p r i n t f ( " sum u=u% g\ n" , sum ) ;
    return 0 ;
}
```

# Scalar product: SSE instructions

```cpp
#include <immintrin.h>
#include <iostream>
#include <algorithm>
> #include <numeric>

int main()
{
  std::size_t const size = 4 * 5; s
  td::srand(time(nullptr));
  float *array0 = static_cast<float *>(_mm_malloc(size * sizeof(float), 16)); float
  *array1 = static_cast<float *>(_mm_malloc(size * sizeof(float), 16)); std::ge
  nerate_n(array0, size, []() { return std::rand()%10; });
  std::generate_n(array1, size, []() { return std::rand()%10; });

  auto r0 = _mm_mul_ps(_mm_load_ps(&array0[0]),          _mm_load_ps(&array1[0]));

  for (std::size_t i = 0; i < size; i+=4)
  {
    r0 = _mm_add_ps(r0, _mm_mul_ps(_mm_load_ps(&array0[i]), _mm_load_ps(&array1[i])));

  }

  float tmp[4] attribute((aligned(16)));
  _mm_store_ps(tmp, r0);
  auto res = std::accumulate(tmp, tmp + 4, 0.0f);

  _mm_free(array0);
```

# Scalar product : Pthreads

```c
# include < stdio.h >
# include < pthread.h >
# define SIZE 256
# define NUM THREADS 4
# define CHUNK SIZE/NUM THREADS

int id [ NUM_THREADS ] ;
double sum , a [ SIZE ] , b [ SIZE ] ; p t
h r e a d t t i d [ NUM_THREADS ] ;
p t h r e a d m u t e x t mutex sum ;

void * dot ( void * id ) {
    size t i ;
    int my first = * ( int *) id  *  CHUNK;
    int my last = ( *( int *) id + 1 ) * CHUNK;
    double sum local = 0 ;

    // Computation
    for ( i = my first ; i < my last ; i ++) su
      m local = sum local + a [ i ]*b [ i ] ;

    p t h r e a d m u t e x l o c k (& mutex sum ) ;
    sum = sum + sum local ;
    p t h r e a d m u t e x u n l o c k (& mutex sum ) ;
    return NULL ;
}
```

```c
int main () {
    size t i ;

    // Initialization
    sum = 0 ;
    for ( i = 0; i < SIZE; i ++) {
      a [ i ] = i  *  0 . 5 ;
      b [ i ] = i  *  2 . 0 ;
    }

    p t h r e a d m u t e x i n i t (&mutex sum , NULL ) ;

    for ( i = 0 ; i < NUM THREADS; i ++) {
      id [ i ] = i ;
      p t h r e a d c r e a t e (& tid [ i ] , NULL , dot ,
                   ( void *)& id [ i ]) ;
    }

    for ( i = 0 ; i < NUM THREADS; i ++)
      p t h r e a d j o i n ( tid [ i ] , NULL ) ;

    p t h r e a d m u t e x d e s t r o y (& mutex

    sum ) ; p r i n t f ( " sum u=u% g \ n " , sum ) ;
    return 0 ;
}
```

# Scalar product : OpenMP

```c
#include <stdio.h>
#define SIZE 256

int main () {
    double sum , a [ SIZE ] , b [ SIZE ] ;
// Initialization
    sum = 0 ;
    for ( size_t i = 0 ; i < SIZE ; i ++ ) {
        a [ i ] = i  *   0 . 5 ;
        b [ i ] = i  *   2 . 0 ;
    }

    // Computation
    #pragma omp parallel for reduction (+: sum )
    for ( size_t i = 0 ; i < SIZE ; i ++ ) {
        sum = sum + a [ i ] * b [ i ] ;
    }

    printf ( "sum u=u% g\ n" , sum ) ;
    return 0 ;
}
```

# Scalar product : MPI

```c
#include <stdio.h>
#include "mpi.h" #
define SIZE 256

int main ( int argc , char* argv[] ) {
    int numprocs , my rank , my first , my last ;
    double sum , sum local , a[SIZE] , b[SIZE] ; M PI
    Init (& argc , &argv ) ;
    MPI Comm size (MPI COMM WORLD, &numprocs );
    MPI Comm rank (MPI COMM WORLD, &my rank ); m y
    first = my rank_ *  SIZE/ numprocs ;
    my last = ( my rank + 1 )  *   SIZE/ numprocs ;

    // Initialization
    sum local = 0. ;
    for ( size t i = 0 ; i < SIZE ; i ++) {
        a [ i ] = i  *  0 . 5 ;
        b [ i ] = i  *  2 . 0 ;
    }

    // Computation
    for ( size t i = my first ; i < my last ; i ++) s u
        m local = sum local + a [ i ]*b [ i ] ;
    M PI A l l r e d u c e (& sum local , &sum , 1 , MPI DOUBLE , MPI SUM , MPI COMM WORLD ) ;

    if ( my rank == 0 )
        printf ("sum u=u% g\ n", sum ) ;
```
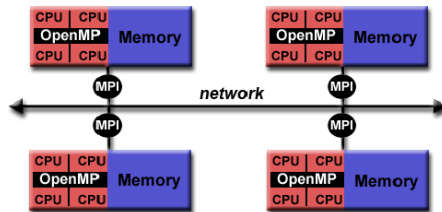
# Hybrid model

- Several MPI processes, each managing a number of threads
    - Inter-process communication via message sending (MPI)
    - Intra-process (thread) communication via shared memory

- Well suited to hybrid architectures
    - one process per node
    - one thread per core

# Outline

# Plan

université
PARIS-SACLAY

7 Around the development of parallel programs

# A performance model; the roofline model

La performance qu'on peut atteindre (Gflop/s) est bornée par

The performance crˆete of the machine,

The maximum bandwidth × operational intensity

où the operational intensity is the number of float operations performed per byte of DRAM transferred

# A performance model; the roofline model

La performance qu'on peut atteindre (Gflop/s) est bornée par

The performance crˆete of the machine,

The maximum bandwidth × operational intensity

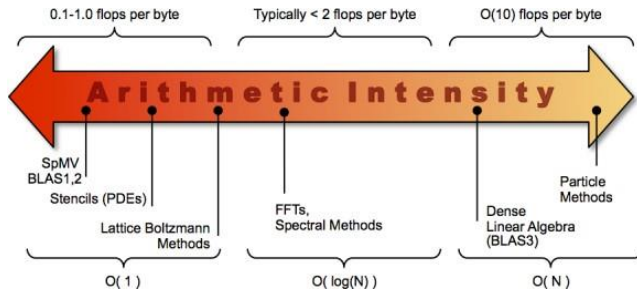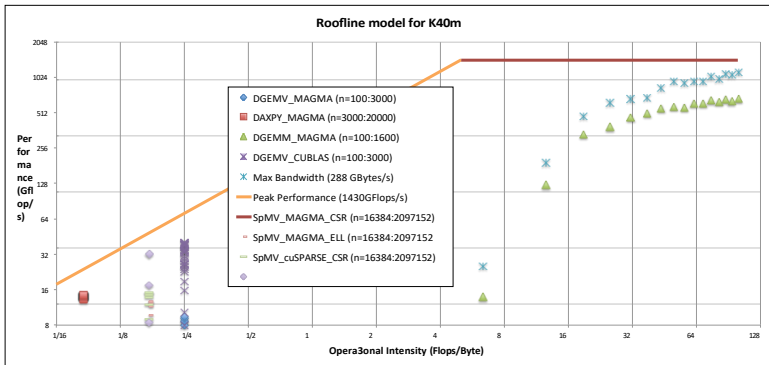the operational intensity is the number of float operations performed per byte of DRAM transferred

- depends on the algorithm and the target architecture
- dense matrix-vector product: $I_{dgemv} \leq \frac{1}{4}$
- Hollow matrix-vector product: $I_{SpMV} \leq \frac{1}{6}$, depending on the storage format (CSR here)
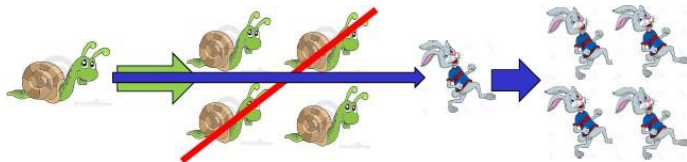
université
PARIS-SACLAY

# the roofline model: operational intensity

# Performance model for NVIDIA Tesla K40

# What? : Core Optimization



## Performance at the heart

- Reducing cache errors: blocking, tilling, loop ordering, ... Vectorization
- (SSE/AVX units)

université
PARIS-SACLAY

## What to do next?

- Identify program <span style="color:red">bottlenecks</span>:
  - to know which parts consume the most running time
  - performance analysis tools can help here ( profilers . . . ) Focus on
  - parallelizing bottlenecks

- Re-structure the program or use/develop another algorithm to reduce the parts that are very slow

- Use existing: optimized parallel software and libraries (IBM's ESSL, Intel's MKL, AMD's AMCL, LAPACK, C++ Parallel STD, . . . )

**université**
**PARIS-SACLAY**

# Qu'est ce qui doit être considéré? : quelques éléments

- Data distribution: 1D, 2D, block, block cyclic, tiles ... Granularity
- Communications
- Synchronization
- Overlapping of calculations and communications Load
- balancing between threads and/or processors
- . . .

université
PARIS-SACLAY

# References

**Contact**

Oguz Kaya
Université Paris-Saclay and LRI, Paris, France
oguz.kaya@lri.com
www.oguzkaya.com