



22nd Summer School on **PARALLEL** COMPUTING

Performance Evaluation

Giovanni Erbacci - [g.erbacci@cineca.it](mailto:g.erbacci@ Cineca.it)

Supercomputing, Applications & Innovation Department - CINECA





Outline

Performance Evaluation

Definition of parallel system

- size of the problem,
- serial execution time
- parallel execution time

Performance indices

- speed-up
- efficiency



Art of Performance evaluation

→ Performance evaluation of HPC Systems

- Benchmarks
- Performance models

→ Performance evaluation of Application Codes



Achieving Efficiency

- Parallel computers allow
 - faster solutions to problems
 - larger problems to be addressed
- An efficient parallel program
 - maximises the amount of work each processor does, and
 - minimises the amount of communication between processes
- How the problem is *decomposed is critical*
 - different ways exist depending on problem



Performance Measures

- **Scientist:**
 - size of the problem,
 - accuracy of the solution, etc. ...
 - Number of operations per unit time (flop / s),
- **Computational scientist:**
 - execution time,
 - speed-up,
 - efficiency.



Quantifying Performance

- Serial computing concerned with complexity
 - how execution time varies with problem size N
 - adding two arrays (or *vectors*) is $O(N)$
 - matrix times vector is $O(N^2)$, *matrix-matrix* is $O(N^3)$
- Look for clever algorithms
 - naïve sort is $O(N^2)$
 - divide-and-conquer approaches are $O(N \log(N))$
- Parallel computing *also concerned with scaling*
 - how time varies with number of processors P
 - different algorithms can have different scaling behaviour
 - but always remember that we are interested in minimum time!



Parallel System

A **parallel system** is the implementation of a parallel algorithm on a specific parallel architecture.

"How to scale the parallel system?"

The **size of the problem** - **W** - is the number of basic operations required by the fastest known sequential algorithm to solve the problem itself, which is equivalent to the concept of computational complexity.



Observations

Sequential algorithm

- computational complexity
- execution time

$f(W)$ (function of the amount of data supplied as input)

Parallel algorithm

- computational complexity
- execution time

$f(W, p, \text{arch})$

p = number of processors used

arch = architecture on which the algorithm is implemented,
(topology of the interconnection network)



Execution Time

Serial execution time - T_s - is the time that elapses between the beginning and the end of the program on a single processor.

Parallel execution time - T_p - is the time that elapses between the beginning of the parallel execution and the time when the last processor finishes his execution.

I/O Problem



Speed-up

$$\text{Speedup} = \frac{\text{Sequential execution time of the best sequential algorithm known}}{\text{Execution time on } P \text{ processors}}$$

- A more honest measure of performance
- Avoids picking an easily parallelizable algorithm with poor sequential execution time

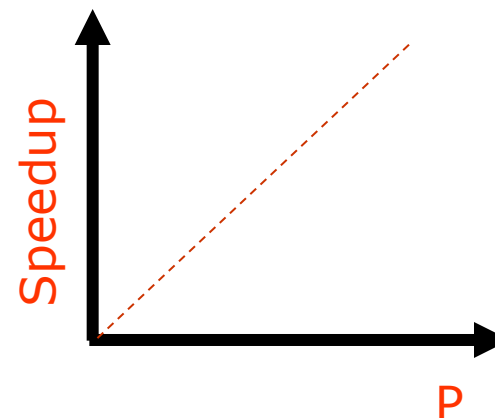


Speed-up

speed-up is defined as:

$$S(W,p) = T_s(W) / T_p(W,p)$$

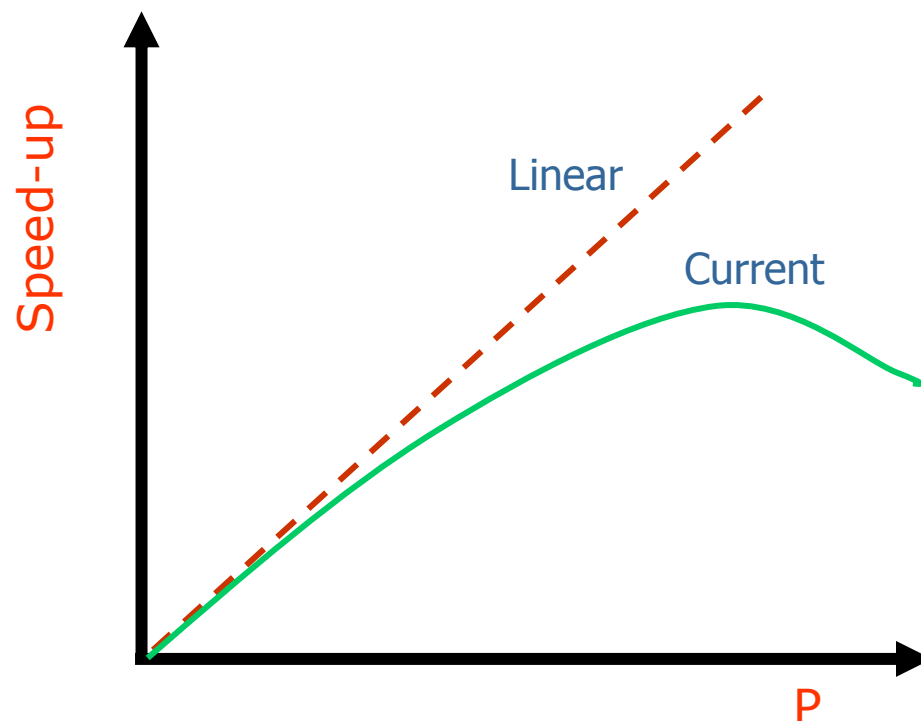
- What does it represent?
- In which interval varies?





What Speedup you get?

- **Linear Speed-up :**
 - Speed-up = N , with N processors
- **Sub-linear Speed-up :**
 - More normal, due to the overhead of initialization, synchronization, communication, etc..
- **Speed-up in flexion:**
 - Decreases as the number of processors grows.
- **super-linear Speed-up**





Super Linear Speed-up

- The speed-up is **superlinear**. when $s > p$,
- This behavior is due to the fact that the program uses the memory hierarchies (due HW),
 - Subparts fit into cache / memory of each node
 - Whole problem does not fit in cache / memory of a single node
- or there is a better code optimization with regard to the scheduling of instructions (due SW) or not determinism eg. in search problems. One thread finds near-optimal solution very quickly => leads to drastic pruning of search space



Efficiency

Efficiency is defined as

$$E(W, p) = S(W, p) / p$$

is the ratio between the sequential execution time and the execution time on p processors, multiplied by p

What amount is?

What does it represent?

In which interval varies?



Scalability of Parallel computers

- There is no precise definition of **scalability**

An Architecture is scalable if it continues to have the same performance per processor when the number of processors and the size of the computational problem to be solved increase.

Scalable MPP systems are designed in such a way that larger versions of the same machine (systems with a greater number of nodes) can be constructed or extended from the same design.

A program scales for a number of processors P , if moving from $p-1$ to p processors an improvement in terms of speed-up is observed .

- Improving load balance / algorithm increases the turn-over to a higher numbers of PEs
- better scaling = ability to utilise larger computers



Sources of overhead

The **efficiency** of real parallel systems often (unless trivial parallel algorithms, embarrassingly parallel) is not maximum because in a parallel system appears **sources of overhead** such as:

- **extra computations** for the parallel algorithm compared to the best sequential algorithm,
- need for **interprocessor communications**,
- **workload imbalance** and more.



Granularity

- How long does it take to communicate? Relevant network metrics:
 - **Bandwidth:** number of bits per second that can be transmitted through the network
 - **Latency:** time to make a message transfer through the network
- Message-passing parallel programs can minimize communication delays by partitioning the program into processes and considering the **granularity** of the process on the machine.

$$\textit{granularity} = \frac{t_{\textit{computation}}}{t_{\textit{communication}}}$$



Serial and parallel fraction

The **serial fraction** of a program, f_s , is the ratio between the time spent in the code inherently sequential and $T_s(W)$.

We define **parallel fraction** of a program, f_p , the ratio between the time spent in the code parallelizable and $T_s(W)$.

Obviously $f_s = (1 - f_p)$.



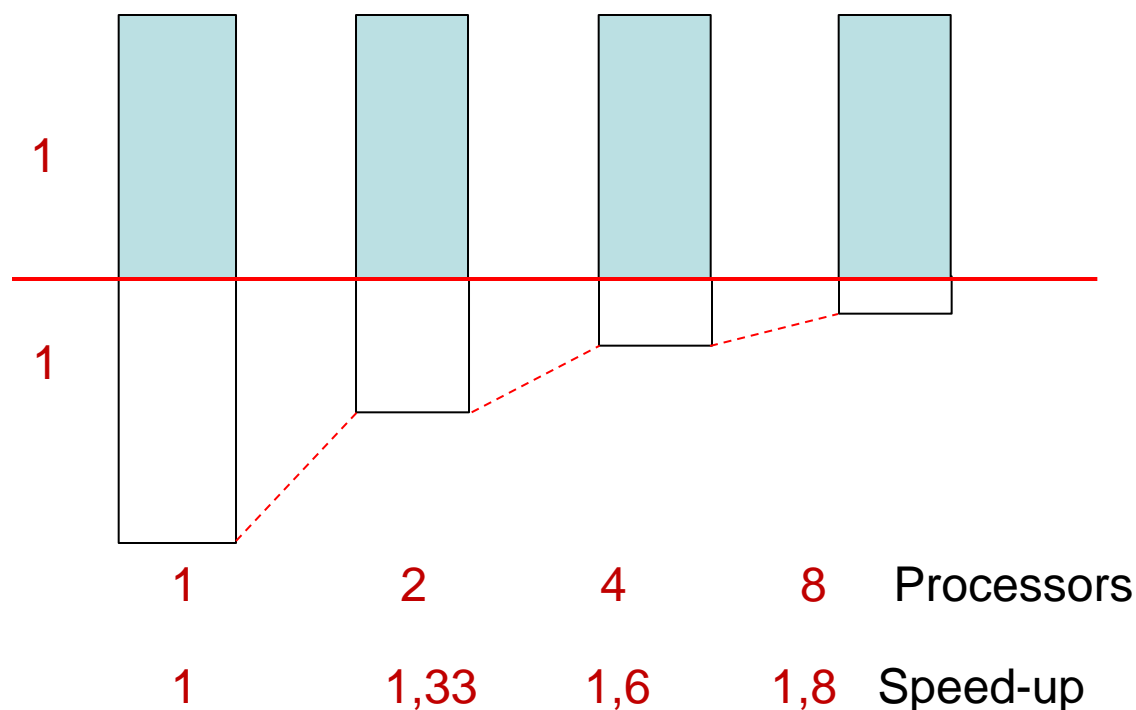
Sequential vs. Parallel

- Sequential execution time: t seconds
- Start-up overhead of parallel execution: t_{st} seconds (depends on architecture)
- Parallel execution time (ideal): $t/p + t_{st}$
- If $t/p + t_{st} > t$, no gain!



The Serial Component

- Amdahl's law
- *“the performance improvement to be gained by parallelisation is limited by the proportion of the code which is serial”*
- Gene Amdahl, 1967





Amdahl's Law

- Assume a fraction f_s is completely serial
 - time is sum of serial and potentially parallel

- Parallel time

- parallel part 100% efficient

$$T_p(W,p) = T_s(W) f_s + (T_s(W) f_p / p)$$

- Parallel speedup

$$\begin{aligned} S(W,p) &= T_s(W) / T_p(W,p) = \\ &= T_s(W) / (f_s \times T_s(W) + (1-f_s) \times T_s(W) / p) \\ &= p / (1 + (p-1) \times f_s) \\ S(W,p) &\rightarrow 1 / f_s \text{ per } p \rightarrow \infty \end{aligned}$$

- for $f_s = 0$, $S = P$ as expected (ie $E = 100\%$)
- otherwise, speedup limited by $1 / f_s$ for any P
- Eg. if 5% of the code is sequential, the speed-up will never exceed 20 even with an infinite number of processors.
- *the performance improvement to be gained by parallelisation is limited by the proportion of the code which is serial*
- impossible to effectively utilise large parallel machines?



Amdahl Law confutation

- Sometimes even a limited speed-up can be a very important milestone for certain applications.

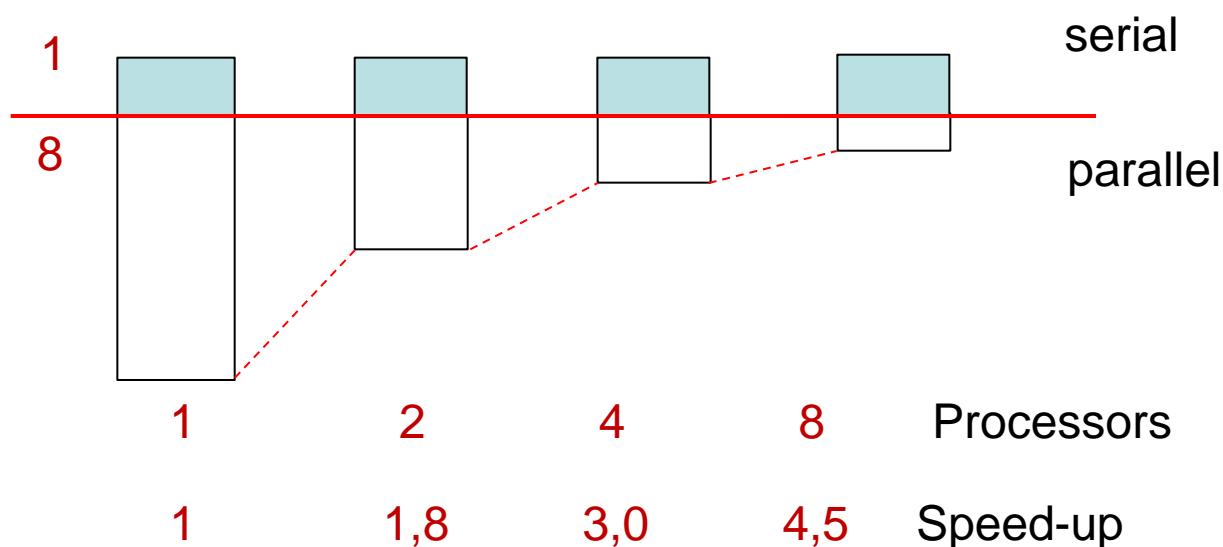
In addition, applications can scale as the number of processors increases:

- a system with a larger number of processors in general allows to solve the biggest problems in a reasonable time
- instead of assuming fixed the **size of the problem** we assume that **the parallel execution** time is fixed
- **Gustafson's Law.**



Gustafson's Law

- Need larger problems for larger numbers of CPUs
- to maintain constant efficiency we need to scale the problem size with the number of CPUs





Performance Models

Fixed-size model: to find the best parallel system by fixing W and varying p .

Fixed-time model: identify on the curve of the execution time, the pairs (W, p) keeping fixed $T_p(W, p)$.

Fixed-memory model: we always work with all available memory.