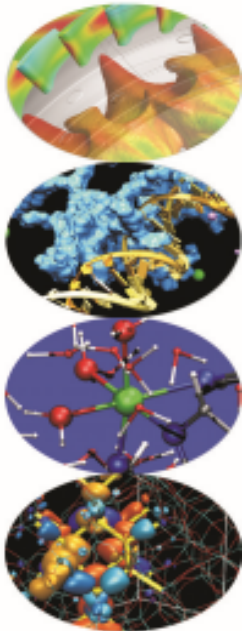


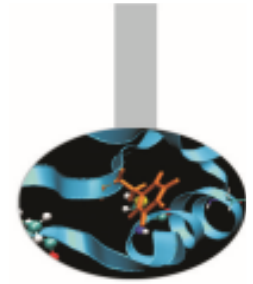
Parallel programming

Introduction to Parallel Computing with MPI and OpenMP

P.Dagna

May 2015





Designing Parallel Program

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.

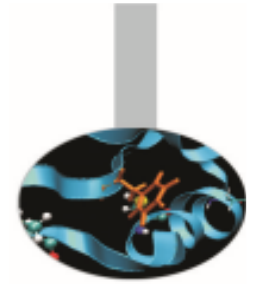
- **Identify the program's hotspots**

The majority of scientific and technical programs usually accomplish most of their work in a few places: profilers and performance analysis tools can help here. Focus on parallelizing the hotspots and ignore those sections of the program that account for little CPU usage.

- **Identify bottlenecks in the program**

For example, I/O is usually something that slows a program down.

May be possible to restructure the program or use a different algorithm to reduce or eliminate unnecessary slow areas

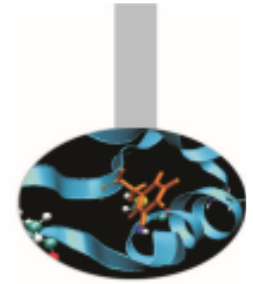


Designing Parallel Program

Overview

In order to design and develop a parallel program, we have to pay attention at several aspects:

- Partitioning
- Communications
- Synchronization
- Data Dependencies
- Load Balancing
- Granularity
- I/O



Designing Parallel Program

Partitioning

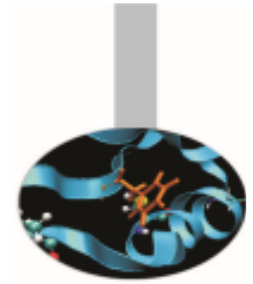
- One of the first steps in designing a parallel program is to break the problem into discrete "chunks" of work that can be distributed to multiple tasks. This is known as decomposition or partitioning.
- There are two basic ways to partition computational work among parallel tasks: ***domain decomposition*** and ***functional decomposition***.

Domain Decomposition:

- In this type of partitioning, the data associated with a problem is decomposed. Each parallel task then works on a portion of of the data.

Functional Decomposition:

- In this approach, the focus is on the computation that is to be performed rather than on the data manipulated by the computation. The problem is decomposed according to the work that must be done. Each task then performs a portion of the overall work.



Designing Parallel Program

Communications

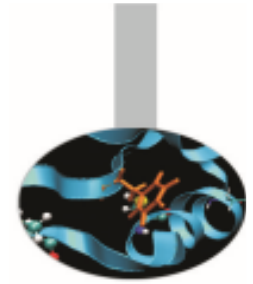
The need for communications between tasks depends upon your problem:

•You DON'T need communications

- Some types of problems can be decomposed and executed in parallel with virtually no need for tasks to share data: these types of problems are often called *embarrassingly parallel*. Very little inter-task communication is required.

•You DO need communications

- Most parallel applications are not quite so simple, and do require tasks to share data with each other. Changes to neighboring data has a direct effect on that task's data.



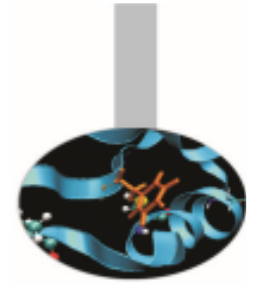
Designing Parallel Program

Communications

There are a number of important factors to consider when designing your program's inter-task communications:

•Cost of communications

- Inter-task communication virtually always implies overhead.
- Machine cycles and resources that could be used for computation are instead used to package and transmit data.
- Communications frequently require some type of synchronization between tasks, which can result in tasks spending time "waiting" instead of doing work.
- Competing communication traffic can saturate the available network bandwidth, further aggravating performance problems.



Designing Parallel Program

Communications

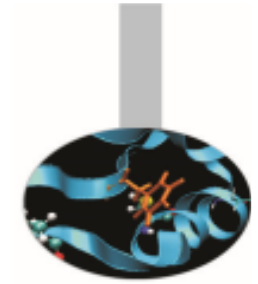
•Latency vs. Bandwidth

–**latency** is the time it takes to send a minimal (0 byte) message from point A to point B. Commonly expressed as microseconds.

–**bandwidth** is the amount of data that can be communicated per unit of time. Commonly expressed as megabytes/sec or gigabytes/sec.

–Sending many small messages can cause latency to dominate communication overheads. Often it is more efficient to package small messages into a larger message, thus increasing the effective communications bandwidth.

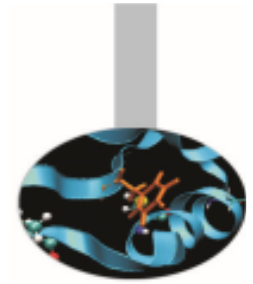
Designing Parallel Program



Communications

•Synchronous vs. asynchronous communications

- Synchronous communications require some type of "handshaking" between tasks that are sharing data. This can be explicitly structured in code by the programmer, or it may happen at a lower level unknown to the programmer.
- Synchronous communications are often referred to as blocking communications since other work must wait until the communications have completed.
- Asynchronous communications allow tasks to transfer data independently from one another. For example, task 1 can prepare and send a message to task 2, and then immediately begin doing other work. When task 2 actually receives the data doesn't matter.
- Asynchronous communications are often referred to as non-blocking communications since other work can be done while the communications are taking place.



Designing Parallel Program

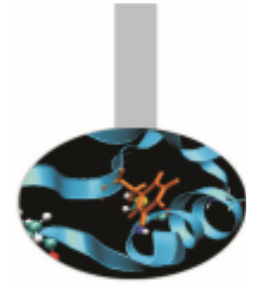
Communications

•Scope of communications

- Point-to-point: involves two tasks with one task acting as the sender/producer of data, and the other acting as the receiver/consumer.
- Collective: involves data sharing between more than two tasks, which are often specified as being members in a common group, or collective.

Both of the two scopings described can be implemented synchronously or asynchronously.

Designing Parallel Program



Synchronization

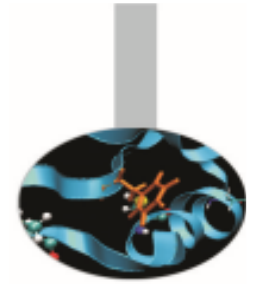
•Barrier

- Usually implies that all tasks are involved: each task performs its work until it reaches the barrier. It then stops, or "blocks". When the last task reaches the barrier, all tasks are synchronized.

•Lock / semaphore

- Can involve any number of tasks: typically used to serialize (protect) access to global data or a section of code. Only one task at a time may use (own) the lock / semaphore / flag.
- The first task to acquire the lock "sets" it. This task can then safely (serially) access the protected data or code. Other tasks can attempt to acquire the lock but must wait until the task that owns the lock releases it. Can be blocking or non-blocking.

•Synchronous communication operations



Designing Parallel Program

Data Dependencies

Definition:

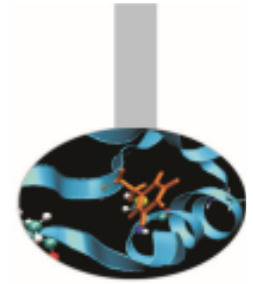
- A ***dependence*** exists between program statements when the order of statement execution affects the results of the program.
- A ***data dependence*** results from multiple use of the same location(s) in storage by different tasks.
- Dependencies are important to parallel programming because they are one of the primary inhibitors to parallelism.

How to Handle Data Dependencies:

- Distributed memory architectures: communicate required data at synchronization points.
- Shared memory architectures: synchronize read/write operations between tasks.

Designing Parallel Program

Load Balancing



- Load balancing refers to the practice of distributing work among tasks so that ***all*** tasks are kept busy ***all*** of the 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.

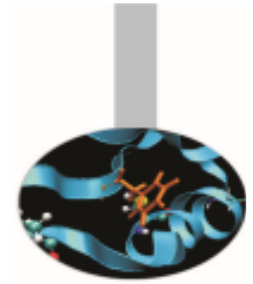
How to Achieve Load Balance:

• Equally partition the work each task receives

- For array/matrix operations where each task performs similar work, evenly distribute the data set among the tasks.
- For loop iterations where the work done in each iteration is similar, evenly distribute the iterations across the tasks.

• Use dynamic work assignment

- It may become necessary to design an algorithm which detects and handles load imbalances as they occur dynamically within the code.



Designing Parallel Program

Granularity

In parallel computing, granularity is a qualitative measure of the ratio of computation to communication.

Periods of computation are typically separated from periods of communication by synchronization events.

Fine-grain Parallelism:

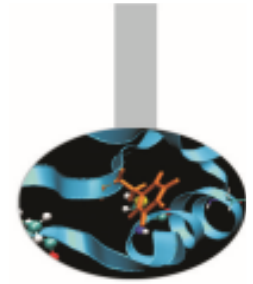
- Relatively small amounts of computational work are done between communication events: if granularity is too fine it is possible that the overhead required for communications and synchronization between tasks takes longer than the computation.

Coarse-grain Parallelism:

- Relatively large amounts of computational work are done between communication/synchronization events but it's harder to load balance efficiently

Designing Parallel Program

I/O

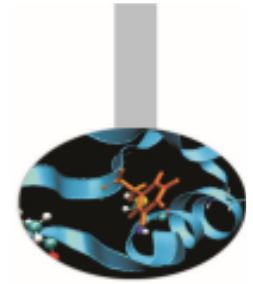


- I/O operations are generally regarded as inhibitors to parallelism
- Parallel I/O systems may be immature or not available for all platforms
- In an environment where all tasks see the same file space, write operations can result in file overwriting
- Read operations can be affected by the file server's ability to handle multiple read requests at the same time

A few advices:

- Reduce overall I/O as much as possible
- Writing large chunks of data rather than small packets is usually significantly more efficient.
- Confine I/O to specific serial portions of the job, and then use parallel communications to distribute data to parallel tasks.
- Use local, on-node file space for I/O if possible. For example, each node may have /tmp filespace which can be used. This is usually much more efficient than performing I/O over the network to one's home directory.

Parallel programs

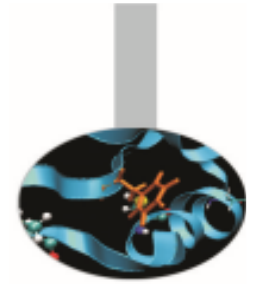


Generally speaking a program parallelisation implies a subdivision of the problem model.

After subdivision the computing tasks can be distributed among more processes.

Two main approaches may be distinguished:

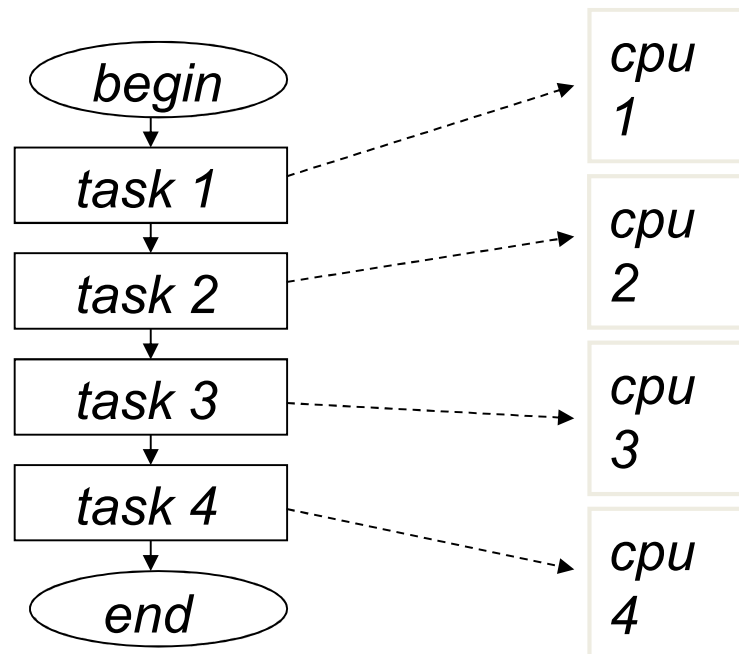
- **Thread** level parallelism
- **Data** level parallelism

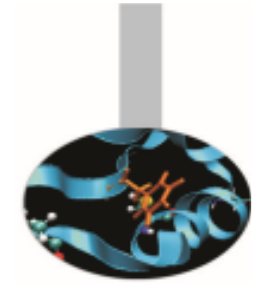


Task parallelism

Thread (or task) parallelism is based on parting the operations of the algorithm.

If an algorithm is implemented with series of independent operations these can be spread throughout the processors thus realizing program parallelisation.

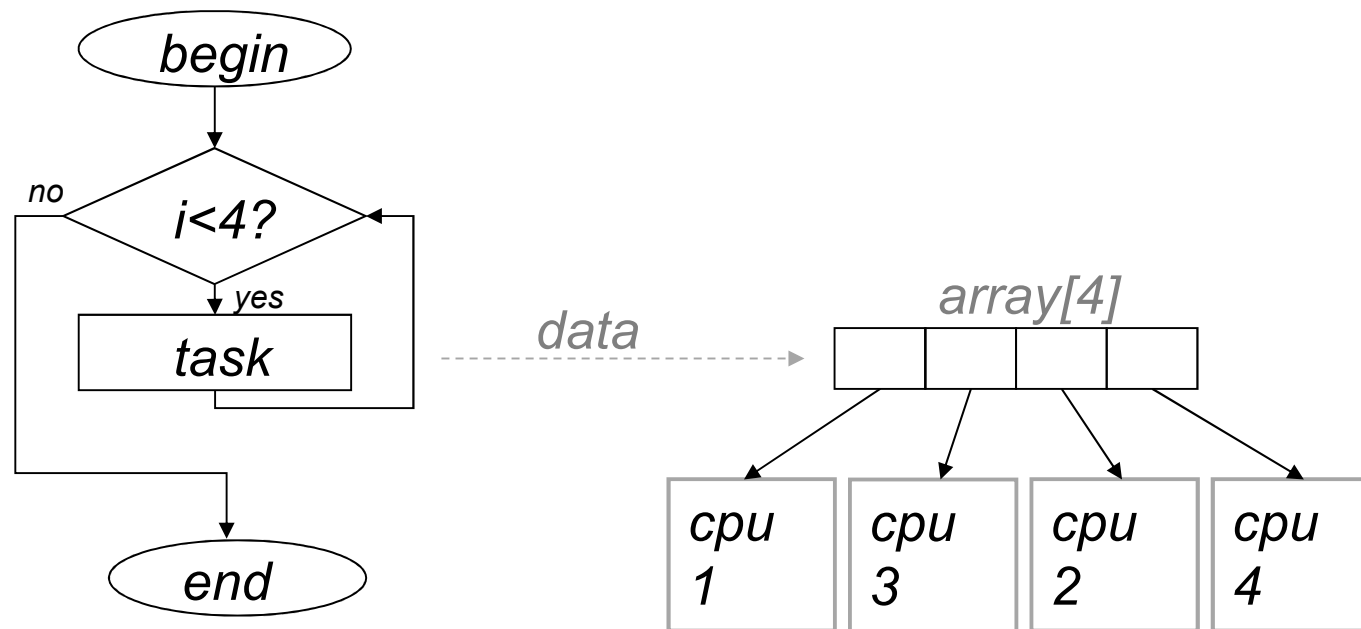


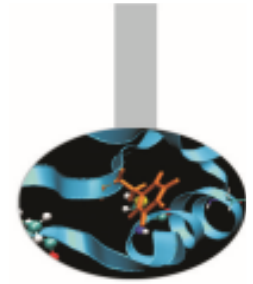


Data parallelism

Data parallelism means spreading data to be computed through the processors.

The processors execute merely the same operations, but on diverse data sets. This often means distribution of array elements across the computing units.





Parallel, concurrent, distributed

What is the difference between parallel, concurrent and distributed programming?

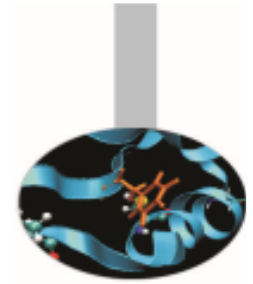
A program is said to be **concurrent** if multiple threads are generated during execution.

A **parallel** program execution is carried on by multiple, tightly cooperating threads.

A program is **distributed** when independent processes do cooperate to complete execution.

Anyhow there are not unique definitions and authors may give different versions. The definitions herein cited are those held by P. Pacheco, “An introduction to parallel programming”.

Parallel, concurrent, distributed



Based on the preceding definitions, parallel and distributed programs are *concurrent* programs, because multiple independent threads are working together to complete computation.

Often a program is said to be *parallel* if it is executed on computing units that share the same memory or are elsewhere connected by a high speed network and usually are very closed together.

Distributed programs instead are executed on processors physically distributed in a (wide) geographical area and connected by a (not so fast) network. Program processes are therefore considered rather independent each other.

Processes, threads and multitasking

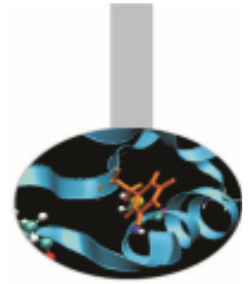


Operating systems are sets of programs that manage software and hardware resources in a computer. Operating systems control the usage of processor time, mass storage, I/O devices and other resources.

When a program execution is started, the operating system generates one or more processes. These are instances of the computer program and contain:

- Executable machine code
- A memory area, often divided in stack, heap and other parts
- A list of computer resources allocated to enable program execution
- Security data to access hardware and software resources
- Informations on the state of the process, i.e. executing, waiting for a resource availability, memory allocation and so on

Processes, threads and multitasking

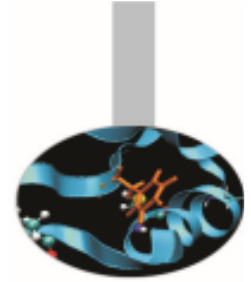


If the operating system is able to manage the execution of multiple processes at one time, it is said to be **multitasking**. On high performance parallel computers multi-tasking is usually of the pre-emptive type, i.e. slices of CPU time are dedicated in turn to each process, unless enough multiple computing units are available.

This means that parallel programs can be executed by concurrent **processes** and the operating system is able to manage their requests. If a computing resource is temporarily unavailable, the requiring process is halted. Anyhow program execution may still be carried on because time slices are granted to the processes that have the availability of the resource.

Parallel programs launched on systems where processors share a global memory are often executed as one process containing multiple **threads**, that share the computing resources of the process including process memory and devices.

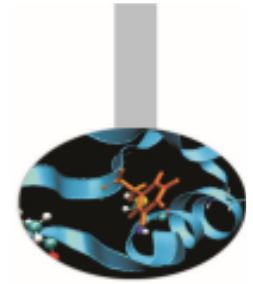
Process interactions



Process interactions may be classified as:

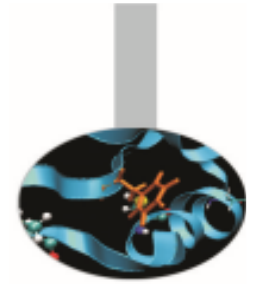
- Cooperation
- Competition
- Interference
- Mutual exclusion
- Deadlock

Cooperation



This kind of interaction is predictable and desirable. Cooperating processes exchange short signals or heavier data transfers.

Process interaction leads to synchronisation and hence to a communication if data are transferred.



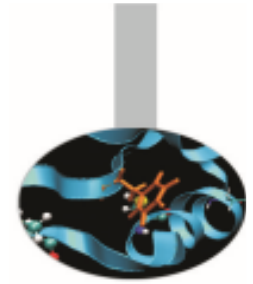
Competition

This kind of interaction is undesirable but nonetheless predictable and unavoidable. It may happen when more processes need to access a common resource that can not be shared (as an example updating a unique counter). Competition may be managed with so called critical sections.

Also contending processes exchange signals and synchronize but in a way different from cooperation.

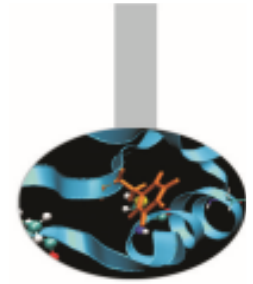
We can distinguish direct or explicit synchronisation (coming from cooperation) from indirect or implicit synchronisation (caused by competition).

Interference



Interference is an unpredictable and undesirable kind of interaction usually arising from errors in developing a parallel program. Errors could come from interactions not required by the implemented algorithm or from interactions not properly handled.

This kind of interaction may show up or not depending by process execution flowing.

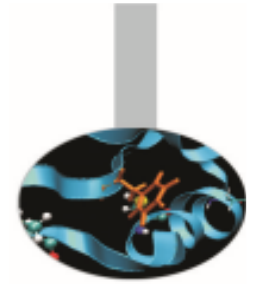


Mutual exclusion

Whenever more processes should not access concurrently a computing resource the problem of realising mutual exclusion has to be managed. This may come up from accessing devices such as writing a disk file or from updating a common memory space.

This kind of problem is often solved using critical sections.

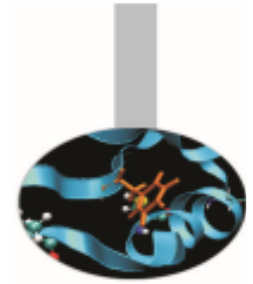
Critical sections do ensure that processes can execute the instructions contained therein but only one at a time.



Deadlock

This undesired situation is always due to programming errors and arises when one or more processes are compelled to wait for something that will never happen.

Processes often enter a deadlock state if they encounter a synchronising point while some other process follow a different executing stream. As an example a program could contain two distinct barriers but processes can reach both of them concurrently.



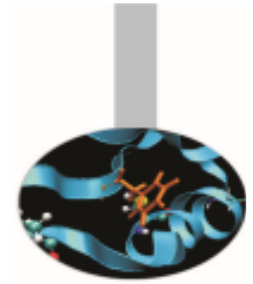
Parallel program performance

The goal of program parallelisation is to reduce execution elapsed time. This is accomplished by distributing execution tasks across the independent computing units. To measure the goodness of the parallelisation effort the time spent in execution by the sequential version of the program (i.e. the program before parallelisation optimisation) must be compared to the time spent by the parallelised version of the program.

Let us call T_{serial} the execution elapsed time of the sequential version of a program and $T_{parallel}$ the execution elapsed time of the parallel version. In an ideal case if we run the program with p computing units (or cores):

$$T_{parallel} = \frac{T_{serial}}{p}$$

If that is true it is said the (parallel) program has a **linear speed-up**.



Speed-up and efficiency

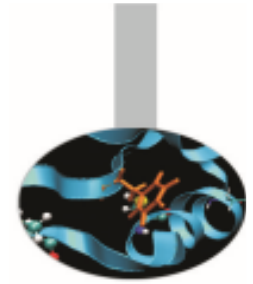
In a real program a linear speed-up is difficult to gain. It has to be considered that the execution flow of the sequential version of the program does not encounter troubles that the parallel version does.

Overheads in a parallel program are introduced by simply dividing the program execution stream. Moreover there is often need of synchronisation and data exchange; furthermore critical sections have to be implemented.

Speed-up is defined as:

$$S = \frac{T_{serial}}{T_{parallel}}$$

The program has a linear speed-up if $S=p$, where p is the number of cores used in executing the program.



Speed-up and efficiency

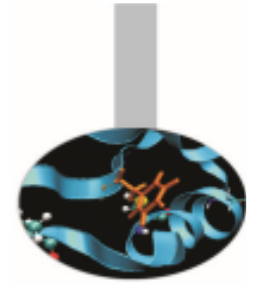
It could be difficult to get a linear speed-up because of the overheads due to synchronisations, communications and often because of an unbalanced distribution of the computing tasks.

This leads to decreasing speed-up while growing the number of cores, because each core brings added overhead.

Efficiency is said to be the ratio between speedup and number of cores:

$$E = \frac{S}{p} = \left(\frac{\frac{T_{serial}}{T_{parallel}}}{p} \right) = \frac{T_{serial}}{p \cdot T_{parallel}}$$

Usually more cores are added, less efficiency is measured.



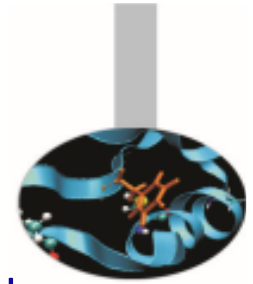
Overhead

Overheads are a significant issue in parallel programs and strongly affect program efficiency.

If overhead delays have to be considered elapsed execution times could be calculated according to:

$$T_{parallel} = \frac{T_{serial}}{p} + T_{overhead}$$

Amdahl's law



If we can analyze a program and measure the portion of code that must be executed sequentially and the part of code that can be distributed across the cores we are able to foresee the program speed-up.

As an example, if it would be possible to parallelize 90% of a program, the remaining 10% of code runs sequentially; then according to Amdahl law:

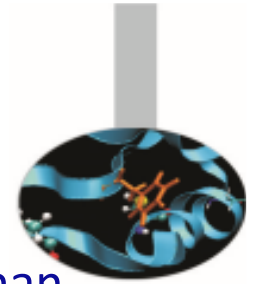
$$T_{parallel} = (0.9 \times T_{serial}) / p + 0.1 \times T_{serial} = 18/p + 2$$

where $p = \text{number of available cores}$

If $T_{serial} = 20 \text{ sec}$ and $p = 6$, then speed-up will be: $S = 20 / (18/p + 2) = 4$.

The time spent in the parallel portion of code decreases as the number of cores increases. Eventually this time tends to zero, but the time spent in the sequential part of the code still remains and strongly limits the program speed-up.

Amdahl's law



As a consequence Amdahl's law tells that speed-up will always be less than $1/r$, where r is the sequential portion of the program.

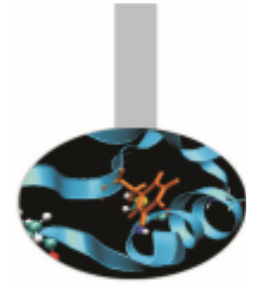
But let us not worry too much!

In real parallel computing world we have to take account of many facets and one of the most important is *problem dimension*. If we consider this we can be interested in Gustafson's (or Gustafson-Barsis') law:

$$S_p^G = p - a(p-1)$$

This formula can be applied to problems for which execution time can be kept constant increasing parallel cores as the problem dimensions increase. This actually applies to many real cases.

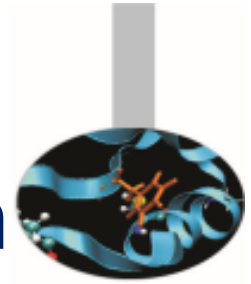
Problem dimensions



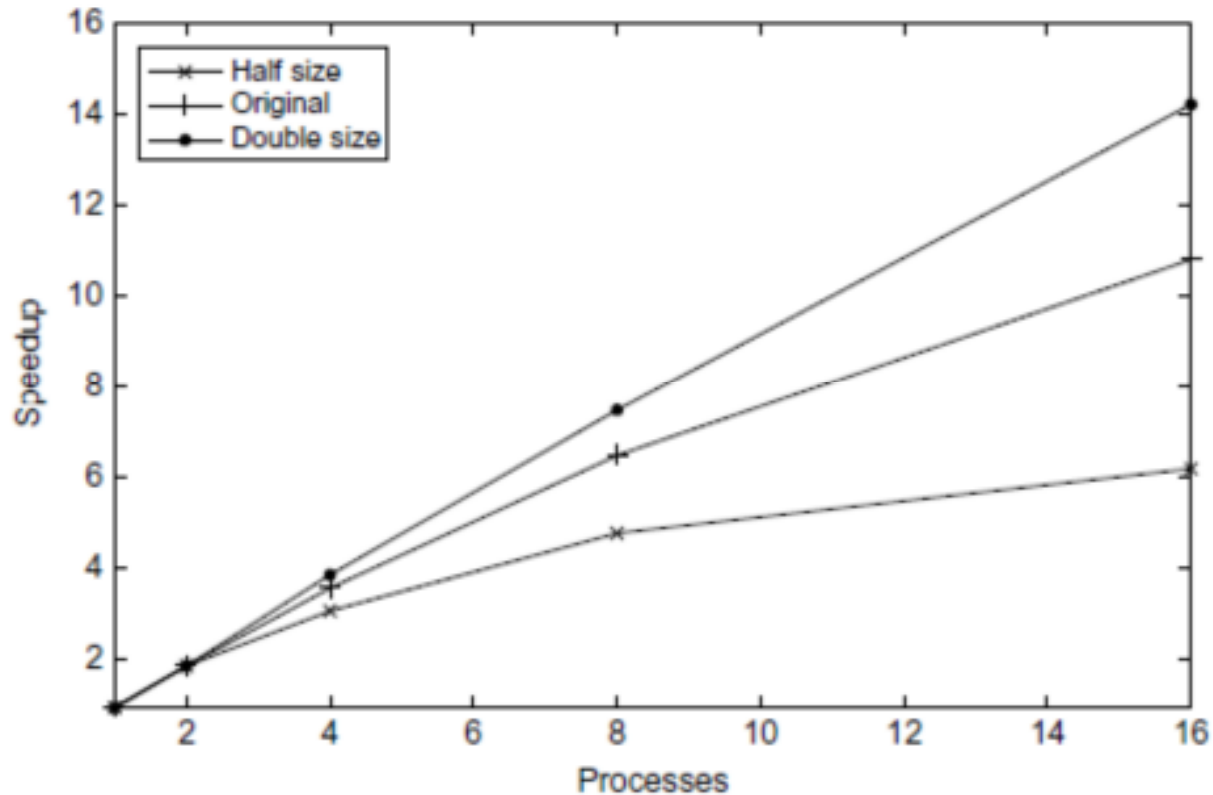
Problem dimension is important because size of data to be computed increases the processors computing time. It is possible to lower global elapsed time by distributing the work across more processors.

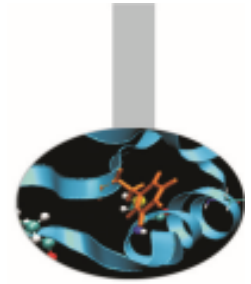
But overheads due to parallelisation stuff will not grow as much, hence speed-up is likely to increase.

Usually, as the dimension of the problem grows, speed-up will grow as well, if enough parallel processors are added.

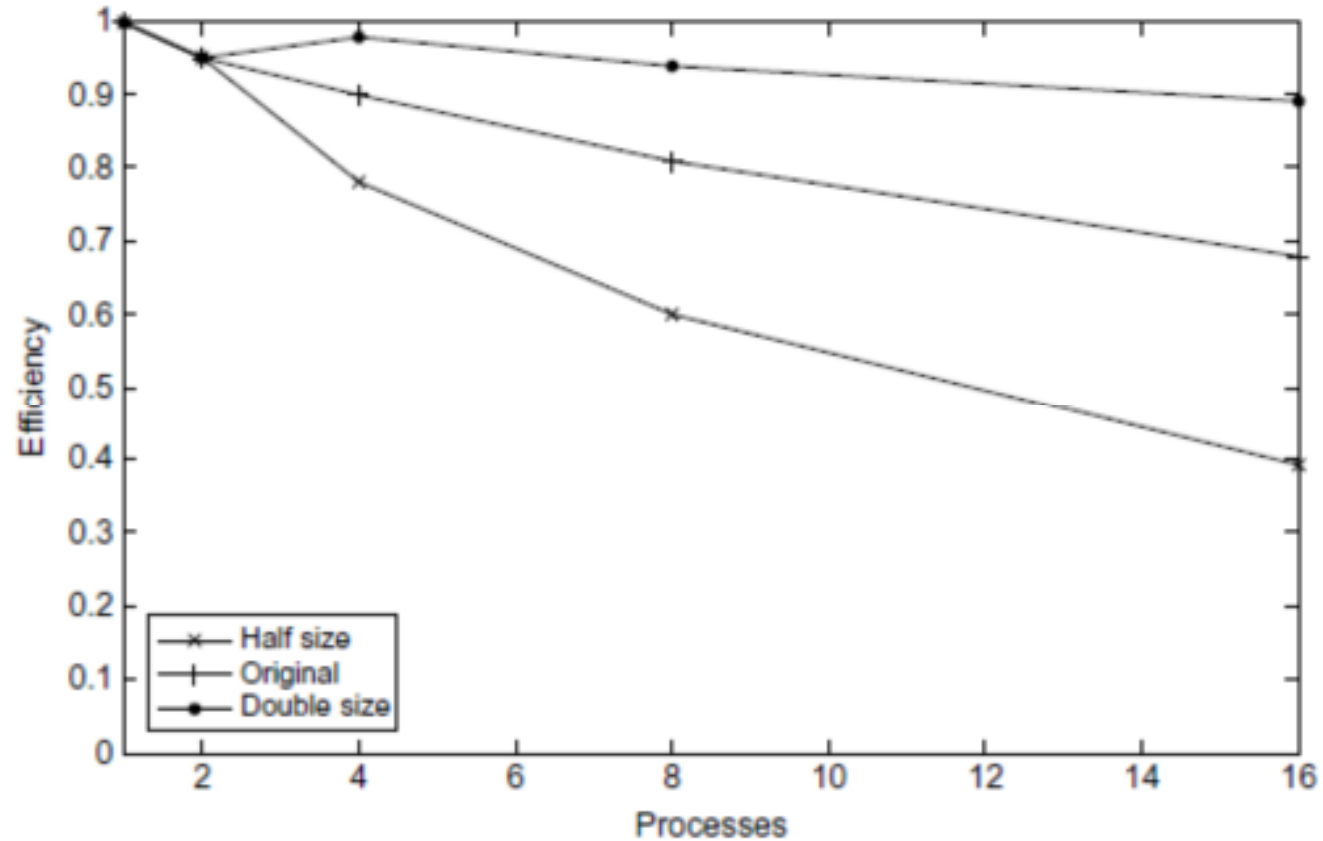


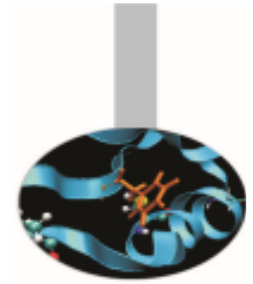
Speed-up and problem dimension





Efficiency and problem dimension



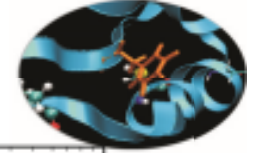


Scalability

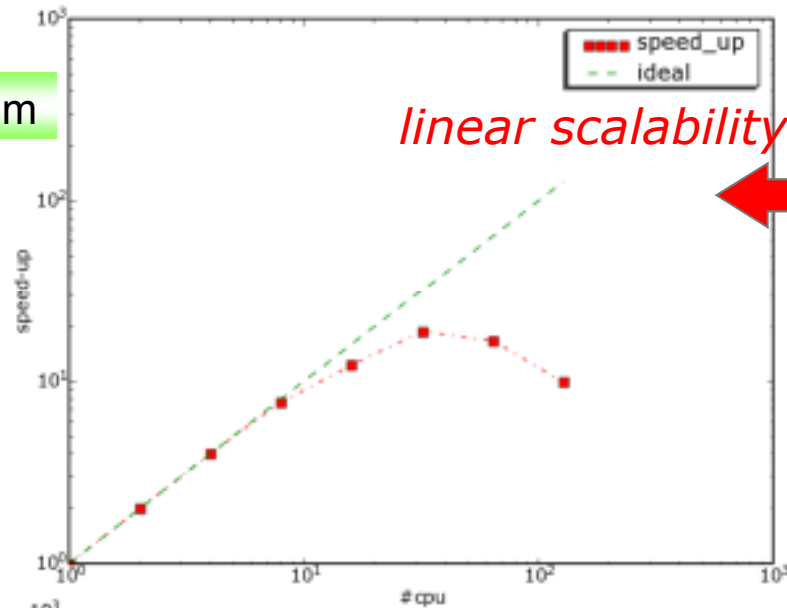
In conclusion, there are basically two ways of evaluating scalability of a program.

If global problem dimension is fixed and efficiency does not decrease while increasing the number of cores, then it is said that the program is **strongly scalable**.

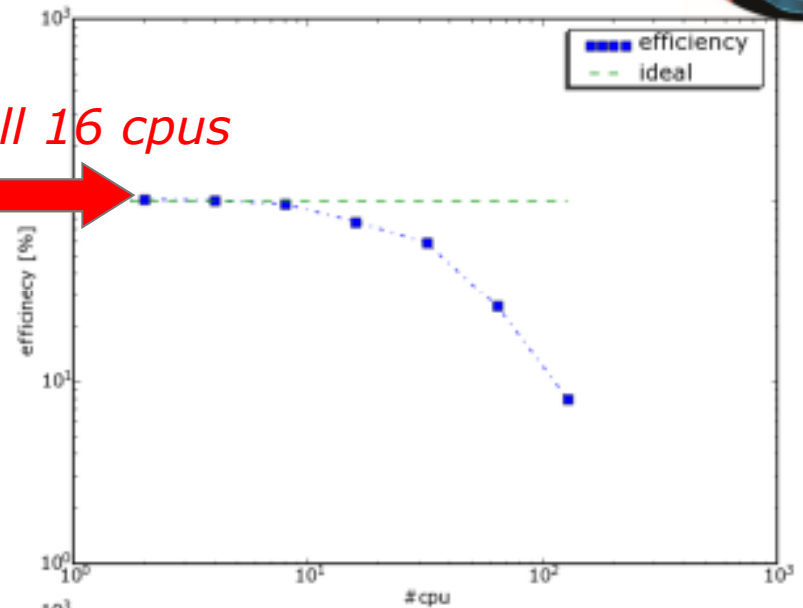
If the efficiency does not decrease when problem dimension per processor (i.e. global dimension has to be augmented as the number of processors increases) is kept almost unchanged, then the program is said to be **weakly scalable**.



Medium

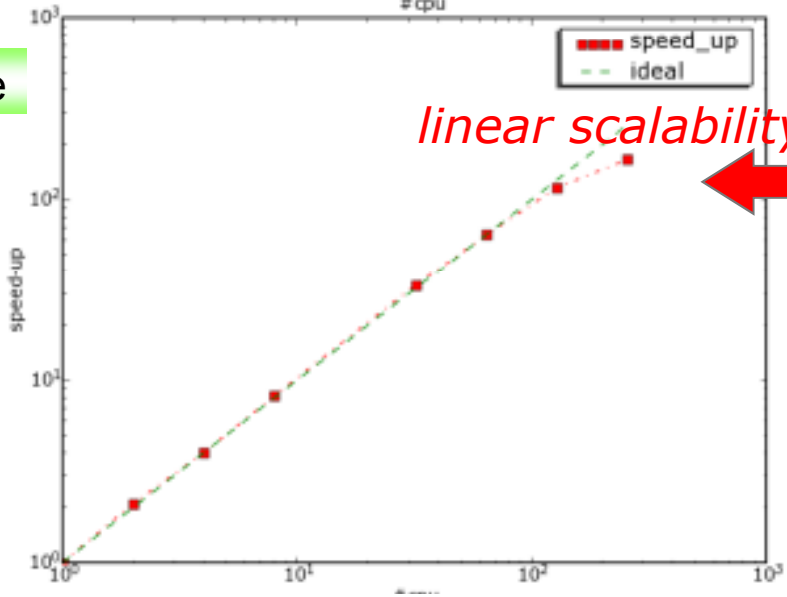


linear scalability till 16 cpus



Large

L3



linear scalability till 128 cpus

