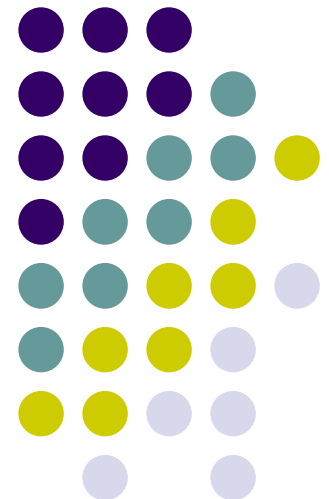


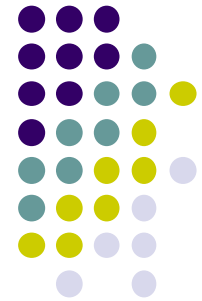
# Computação Paralela

## Optimising performance (MPI)

João Luís Ferreira Sobral  
Departamento de Informática  
Universidade do Minho

Dez 2020





# Performance of parallel applications

## Performance models

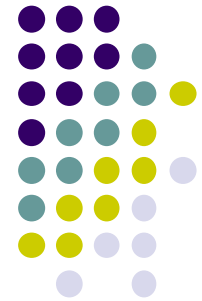
- Make it possible to compare the complexity of algorithms, their scalability and identify bottlenecks before a considerable time is invested in implementation

## What is the definition of performance?

- There are multiple alternatives:
  - **Execution time**, efficiency, scalability, memory requirements, throughput, latency, project costs / development costs, portability, reuse potential
  - The importance of each one depends on the concrete application
- Most common measure in parallel applications is **speed-up**:  $t_{seq}/t_{par}$

## Amdahl law

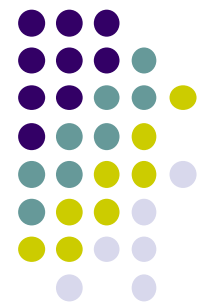
- The sequential component of an application limits the maximum speed-up
  - If  $s$  is the sequential fraction of an algorithm then the maximum possible gain is  $1/s$ .
- Reinforces the idea that we should prefer algorithms suitable for parallel execution: *think parallel*.



# Performance of parallel applications

## Performance models

- Should explain observations and predict behaviour
  - Defined as a function of the problem dimension, number of processing units (PU), number of tasks, etc.
- **Execution time**
  - Time measured since the first process (or thread) starts execution until the last process (thread) terminates (wall time)
- **$T_{exec} = T_{comp} + T_{comm} + T_{free}$** 
  - **Computation time** – time spent in computations
    - Excludes communication/synchronization and free time.
    - The sequential version can be used to estimate  $T_{comp}$ .
  - **Free time** - when a PU becomes starved (without work)
    - Can be complex to measure since it depends on the order of tasks
    - Can be minimized with adequate load distribution and/or overlapping computation and communication



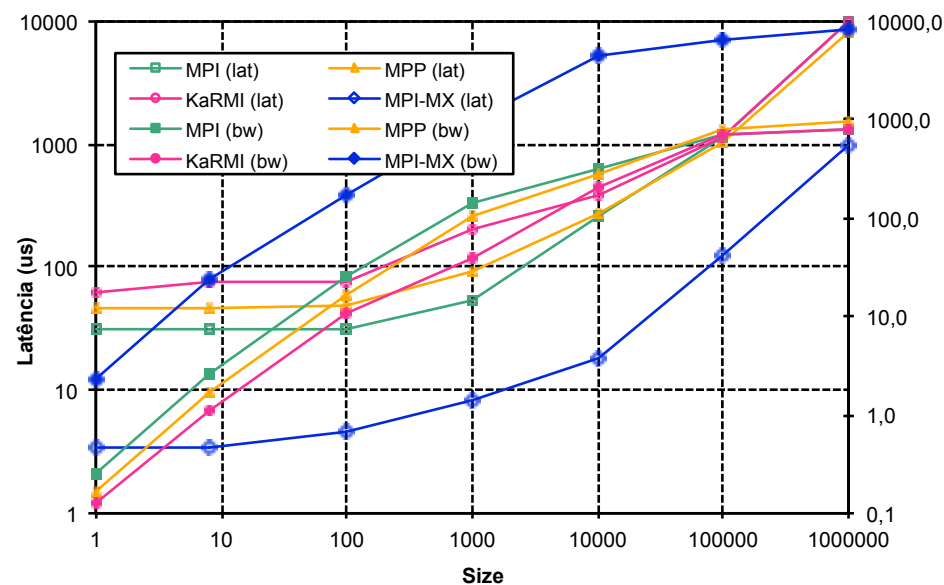
# Performance of parallel applications

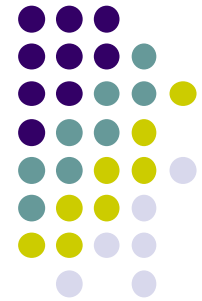
## Performance models(cont.)

- **Communication time** – time that PU spend sending/receiving data.
  - Computed using communication latency (ts) and throughput (1/tw):
    - **$T_{mens} = ts + twL$**

ts and tw can be obtained experimentally, by a ping-pong test and a linear regression.

Tamnhos	MPI (lat)	MPI (bw)	MPI-MX (lat)	MPI-MX (bw)	MPP (lat)	MPP (bw)	KaRMI (lat)	KaRMI (bw)
1	31	0,3	3,4	2,3	46	0,2	63	0,1
8	31	2,6	3,4	23,7	46	1,7	75	1,1
100	31	25,6	4,7	168,7	48	16,6	75	10,7
1000	55	146,3	8,1	983,9	94	106,4	200	40,0
10000	258	310,3	18,4	4355,9	279	286,0	387	206,0
100000	1136	703,8	125,5	6373,0	1017	786,5	1137	703,0
1E+06	9859	811,4	970,2	8246,0	8282	953,2	9787	817,0
ts (us)	31		3		46		63	
tw (us)	0,010		0,001		0,008		0,010	





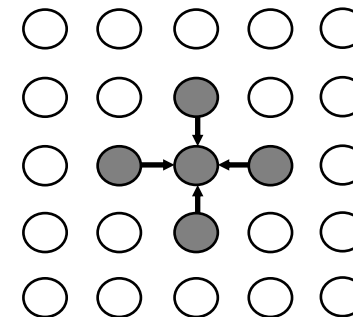
# Performance of parallel applications

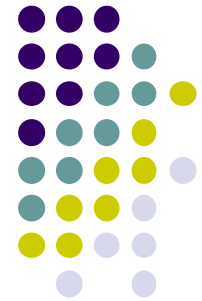
## Performance models – Example Jacobi Method

- Iterative method, at each iteration the new matrix value ( $X^{t+1}$ ) is computed as the average of neighbour values from previous iteration ( $X^t$ )

$$X_{i,j}^{(t+1)} = aX_{i,j}^{(t)} + b(X_{i-1,j}^{(t)} + X_{i,j-1}^{(t)} + X_{i+1,j}^{(t)} + X_{i,j+1}^{(t)})$$

```
for(int t=0; t<Niter; t++) {  
    for(int i=1; i<N-1; i++)  
        for(int j=1; j<N-1; j++)  
            r[i][j] = a*x[i][j]+ b*(x[i-1][j] + x[i+1][j] + x[i][j-1] + x[i][j+1]);  
    // x = r on the next iteration  
}
```





# Performance of parallel applications

## Performance models – Example Jacobi Method

- **Execution time of one iteration** for a  $N \times N$  matrix, on  $P$  processors, using a partition by rows with  $N/P$  rows per processor.

**$T_{comp}$**  = operations per element  $\times$  n° elements per processor  $\times t_c$

$$= 6 \times (N \times N/P) \times t_c \quad (t_c = \text{time for a arithmetic operation})$$

$$= 6t_c N^2/P$$

**$T_{comm}$**  = messages per processor  $\times$  time required for each message

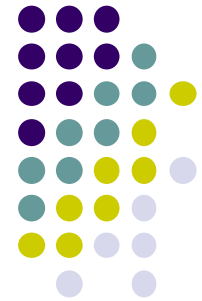
$$= 2 \times (t_s + t_w N)$$

**$T_{free}$**  = 0 , since in this problem the workload is well distributed

$$\mathbf{T_{exec} = T_{comp} + T_{comm} + T_{free}}$$

$$= 6t_c N^2/P + 2t_s + 2t_w N$$

$$= O(N^2/P + N)$$



# Performance of parallel applications

## Performance models – Example (cont)

- In certain cases, execution time may not be the most adequate performance measure.
- Speed-up and efficiency are two related metrics.
- **Speed-up (G)** indicates the reduction in execution time attained in P processors
  - Ratio between the *best sequential algorithm* and the execution time of the parallel version

$$\text{speed-up} = T_{\text{seq}} / T_{\text{par}},$$

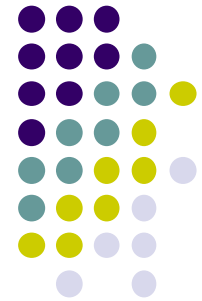
- **Efficiency (E)** gives the fraction of time that processors perform useful work:

$$E = T_{\text{seq}} / (P \times T_{\text{par}})$$

- Jacobi case:

$$G = \frac{6t_c N^2 P}{6t_c N^2 + 2Pt_s + 2Pt_w N}$$

$$E = \frac{6t_c N^2}{6t_c N^2 + 2Pt_s + 2Pt_w N}$$

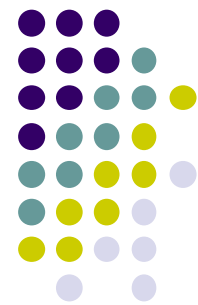


# Performance of parallel applications

## Scalability analysis

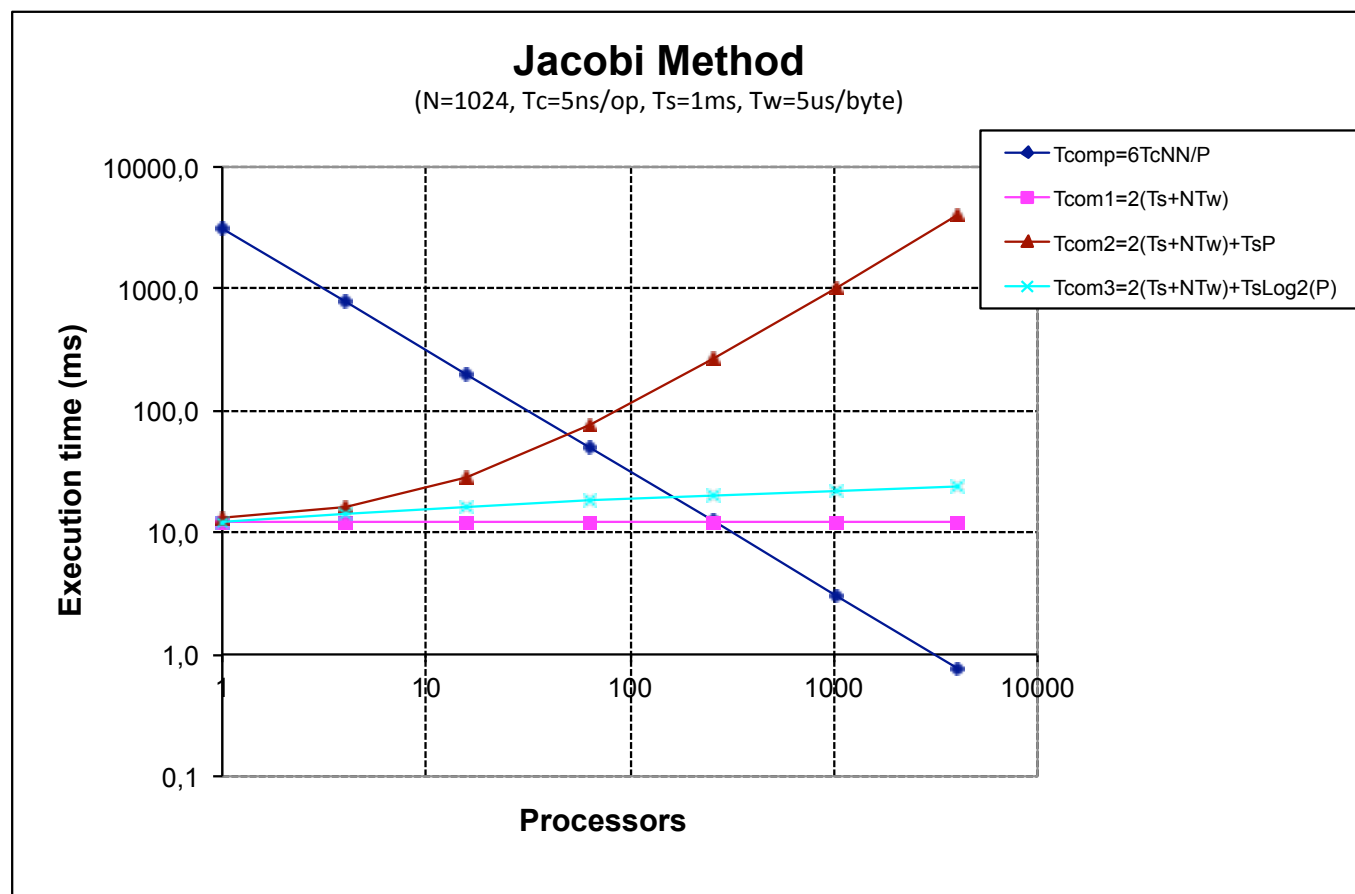
- Execution time, speed-up and efficiency can be used for quantitative analysis of performance
  - Jacobi example:
    - Execution time decreases when  $P$  increases, but it is limited by the time to exchange two lines
    - Execution time increases with  $N$ ,  $t_c$ ,  $t_s$  e  $t_w$
    - Efficiency decreases when  $P$ ,  $t_s$  e  $t_w$  increase
    - Efficiency increases with  $N$  and  $T_c$ ;
- $$T_{\text{exe}} = 6t_c N^2 / P + 2t_s + 2t_w N$$
- $$E = \frac{6t_c N^2}{6t_c N^2 + 2Pt_s + 2Pt_w N}$$
- *Scalability for problems with fixed size (strong scalability)*
    - Analysis of  $T_{\text{exec}}$  and  $E$  when  $P$  increases
    - In general,  $E$  decreases.  $T_{\text{exec}}$  can increase if it has a positive power of  $P$ .
  - *Scalability for problems with variable size (weak scalability)*
    - In some cases, more processors are used to solve larger problems, keeping the same efficiency levels
    - **Isoefficiency** indicates what is the required increase in the problem dimension, to keep the same efficiency, when the number of processors increases

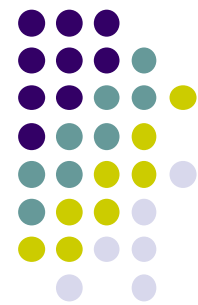




# Performance of parallel applications

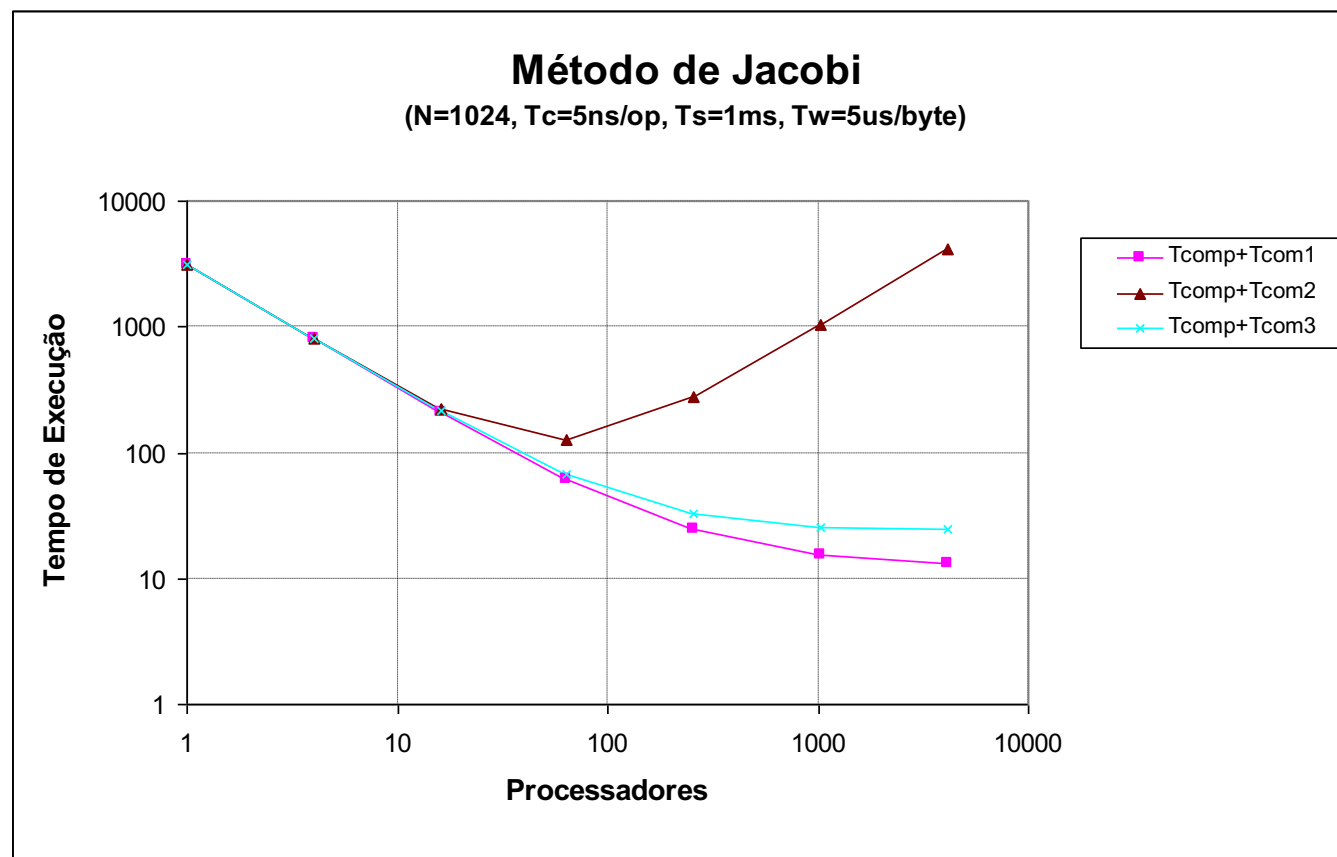
## Scalability analysis (cont)

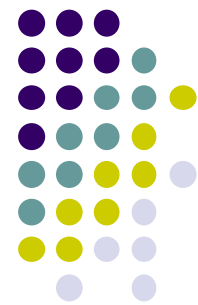




# Performance of parallel applications

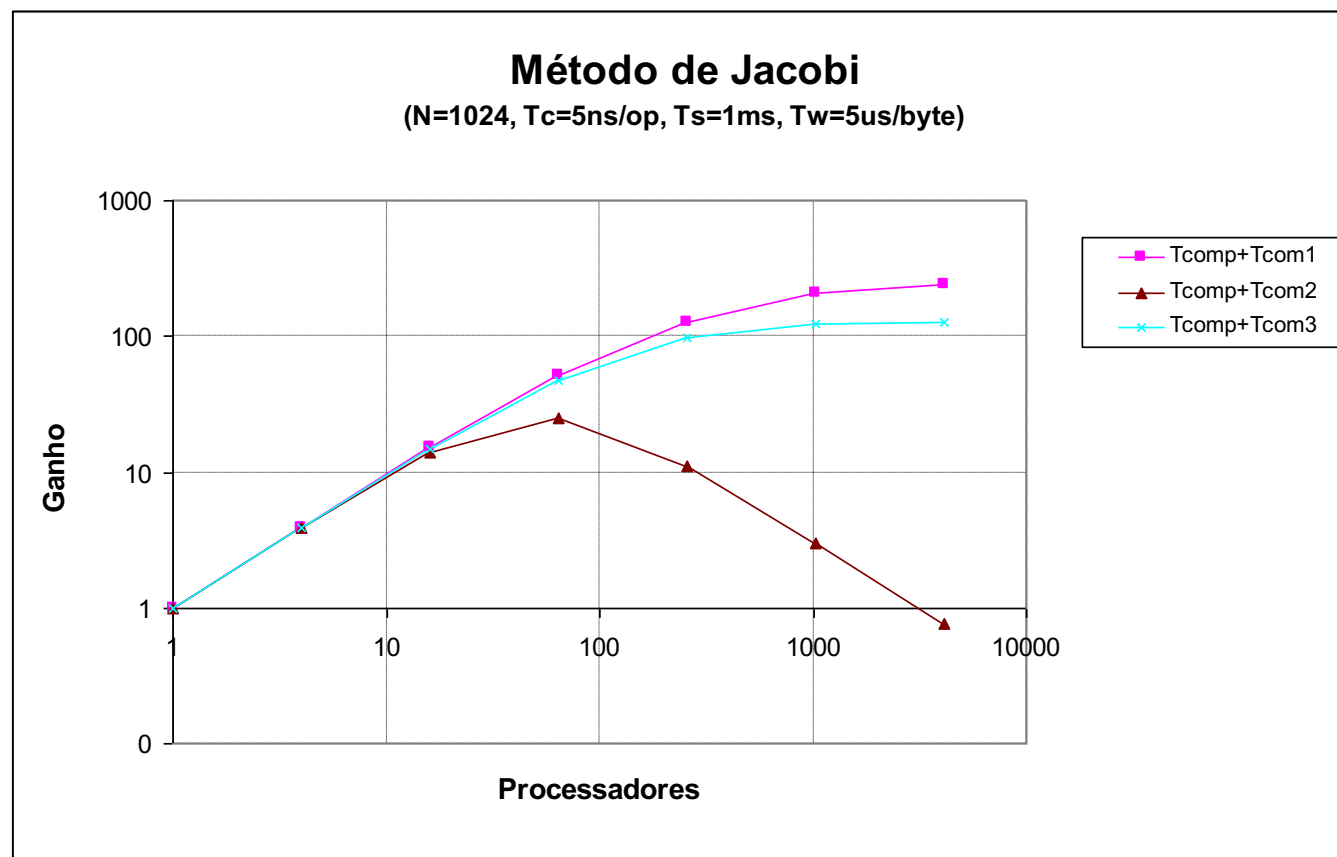
## Scalability analysis (cont)

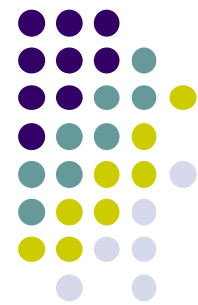




# Performance of parallel applications

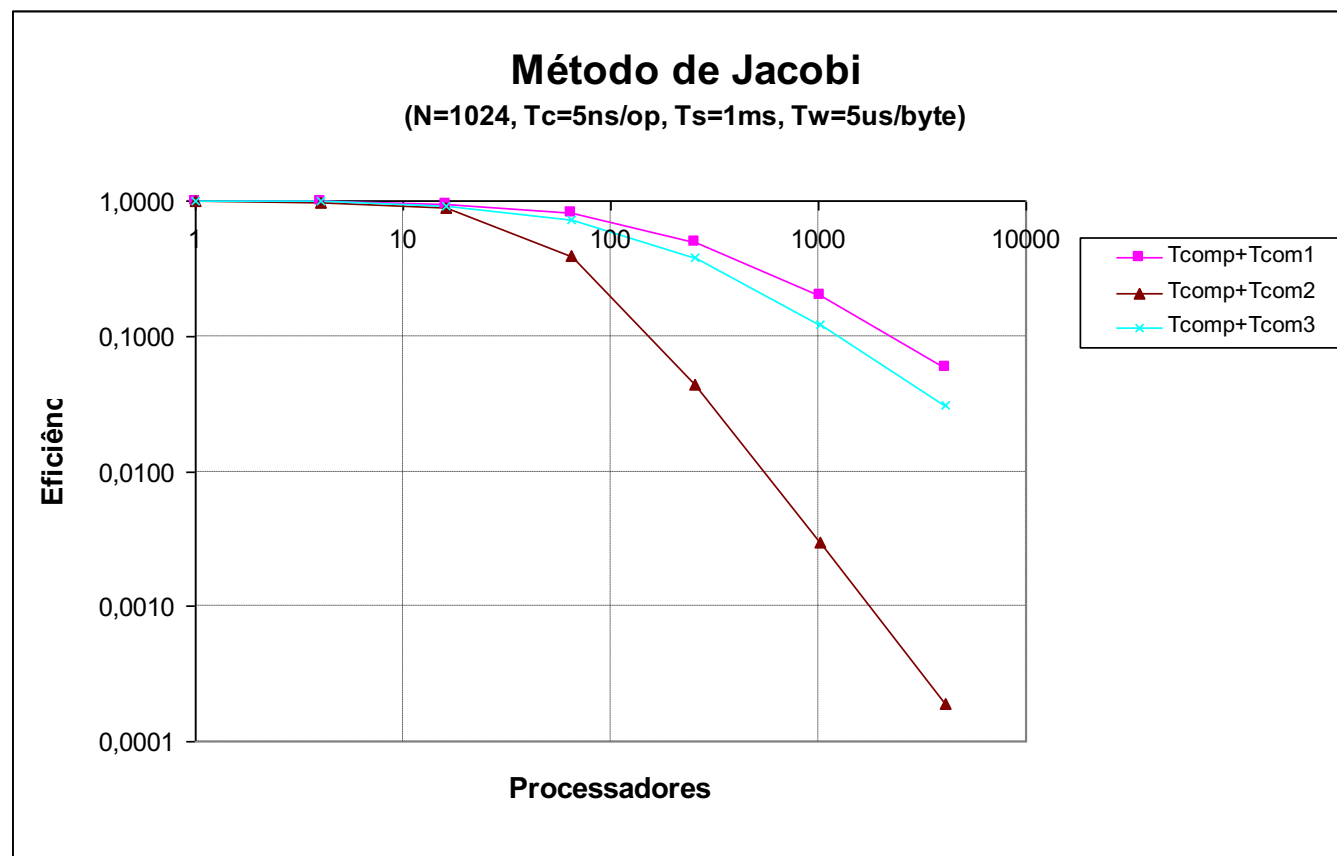
## Scalability analysis (cont)

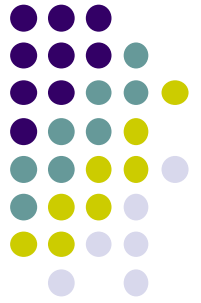




# Performance of parallel applications

## Scalability analysis (cont)

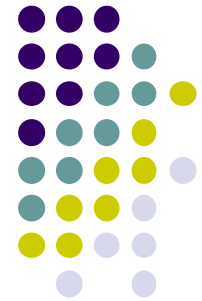




# Performance of parallel applications

## Measuring time in MPI

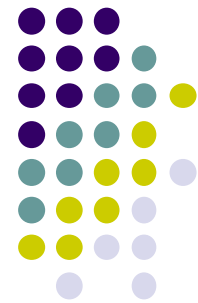
- **Time functions in MPI**
  - `double MPI_Wtime()` – returns the wall time (high resolution)
  - `double MPI_Wtick()` – returns the clock resolution (in seconds)
- **Wall time can differ from process to process**
  - There is no notion of “global time”
    - Each machine provides a local wall time
  - Application execution time should be wall time of the slowest process
  - Note: in some parallel algorithms process termination is not trivial



# Performance of parallel applications

## Experimental study and evaluation of implementations

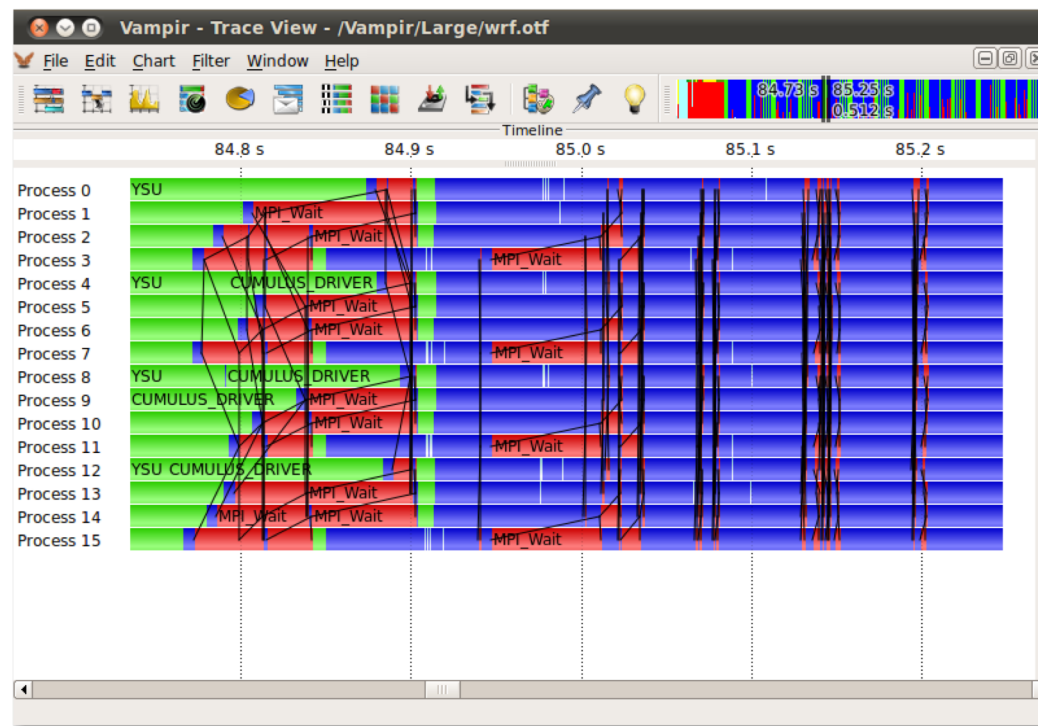
- Parallel computing has a strong experimental component
  - Many problems are too complex for a realization only based on models
  - Performance model can be calibrated with experimental data (e.g.,  $T_c$ )
- How to ensure that results are precise and reproducible?
  - Perform multiple experiments and verify clock resolution
  - Results should not change among in small difference: less than 2-3%
- Execution profile:
  - Gather several performance data: number of messages, data volume transmitted
  - Can be implemented by specific tools or by directly instrumenting the code
    - There is always an overhead introduced in the base application
- *Speed-up anomalies*
  - superlinear (superior to the number of processors) – in most cases it is due the cache effect

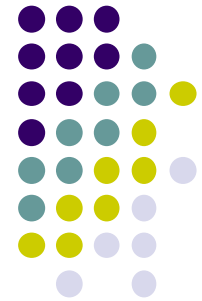


# Performance of parallel applications

## Technique to measure the application time-profile (*profiling*)

- **Polling:** the application is periodically interrupted to collect performance data
- **Instrumentation:** code is introduced (by the programmer or by tools) to collect performance data about useful events
- Instrumentation tends to produce better results but also produces more interference (e.g., overhead)
- Exemple: vampir





# Performance of parallel applications

## Distributed memory (MPI) vs Shared memory (OpenMP) optimisation

- **Distributed memory (vs shared memory)**
  - Data placement is explicit (vs implicit)
  - Static scheduling is preferred (vs dynamic)
  - Synchronization is costly (only performed by global barriers & message send)
- How to improve scalability on distributed memory?
  - Minimise communication among processes
    - Eventually duplicating computation
  - Minimise idle (free) time with a good load distribution
- Practical advise
  - Measure communication overhead
  - Measure load balance
  - Avoid centralised control