

# Speedup

Vediamo come e' possibile caratterizzare e studiare le performance di un algoritmo parallelo:

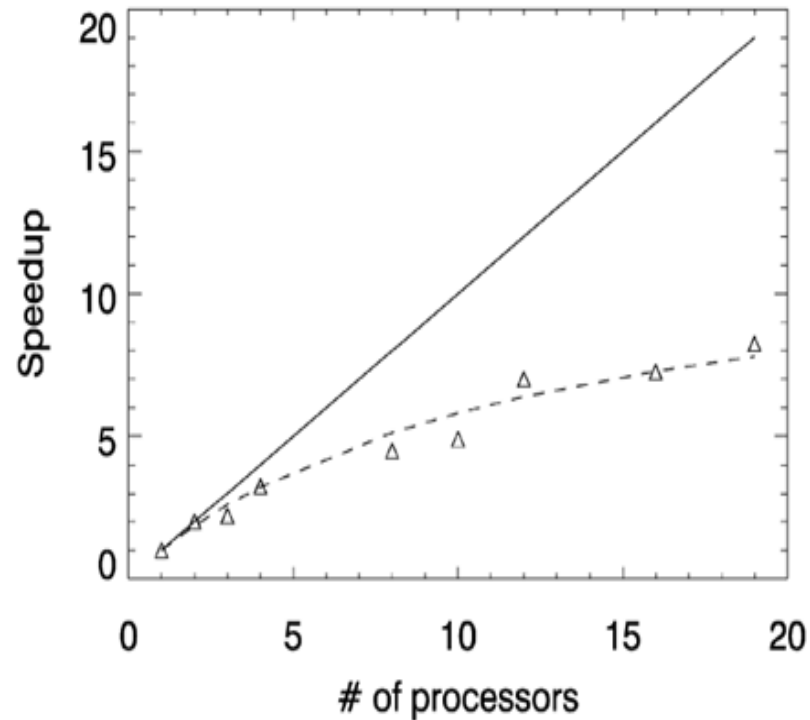
$$S_n = \frac{T_1}{T_p(n)}$$

Dove  $T_1$  e' il tempo impegnato dal miglior algoritmo seriale conosciuto, mentre  $T_p(n)$  e' il tempo di esecuzione dell'algoritmo parallelo su  $n$  processori.

Si definisce anche lo Speedup relativo in cui, invece di usare  $T_1$  si usa  $T_p(1)$ .

# Speedup

Il miglior algoritmo parallelo possibile presenta uno speedup lineare, quindi  $S_n = n$ , in pratica questo limite e' difficilmente raggiungibile.



# Speedup ed Efficienza

**Speedup superlineare:** e' il caso in cui  $S_n > n$ . Casi del genere sono rari, ma ci sono differenti motivazioni che possono portare a risultati del genere. Ad esempio effetti della cache (le dimensioni totali delle cache dei vari processori possono essere tali da contenere tutti i dati relativi al problema).

Oltre allo speedup e' possibile definire l'**efficienza:**

$$E_n = \frac{S_n}{n}$$

il suo valore tipicamente varia tra zero ed uno, ed e' una misura di quanto bene il nostro algoritmo parallelo sfrutta i processori. Anche in questo caso e' possibile parlare di efficienza relativa.

# Scalabilita'

**Costo:** quanto uso di CPU e' richiesto

$$C = n T_p(n) = T_1 / E$$

**Scalabilita':** capacita' di rimanere efficiente su una macchina parallela aumentando il numero di processori (o comunque aumento le dimensione del problema proporzionalmente al numero di processori usati).

# Legge di Amdahl

La legge di Amdahl serve a determinare il miglioramento massimo ottenibile quando solo una parte del codice è stata migliorata:

$$\frac{1}{\sum_{k=0}^n \left( \frac{P_k}{S_k} \right)}$$

- .  $P_k$  è la percentuale di codice che può essere migliorata (o peggiorata)
- .  $S_k$  è lo speedup relativo alla k-esima parte di codice di cui sopra
- . k semplicemente indicizza le varie parti di codice

# Legge di Amdahl

Nel caso specifico del calcolo parallelo la legge puo' essere semplificata come segue:

$$S_n = n / (n * F + (1 - F))$$

dove  $F$  ( $0 < F < 1$ ) e' la frazione di codice che non puo' essere parallelizzata.

Quindi data una frazione di codice che e' stata parallelizzata, ed il numero di processori che possiamo usare, la legge ci dice quale e' il massimo speedup che possiamo aspettarci.

# Legge di Amdahl

Vediamo come si deduce tale legge. Se  $t_s$  e' il tempo impiegato dal singolo processore per il completamento del task, nel caso parallelo avremo un tempo pari a  $F*t_s + [(1-F)*t_s / n]$  e dunque la legge di Amdahl dice che:

$$S_n = t_s / (F*t_s + [(1-F)*t_s / n])$$

$$S_n = 1 / (F + [(1-F) / n])$$

$$S_n = n / (n*F + (1-F))$$

$(1 - F)$  rappresenta chiaramente la frazione di codice parallelo.

# Legge di Amdahl

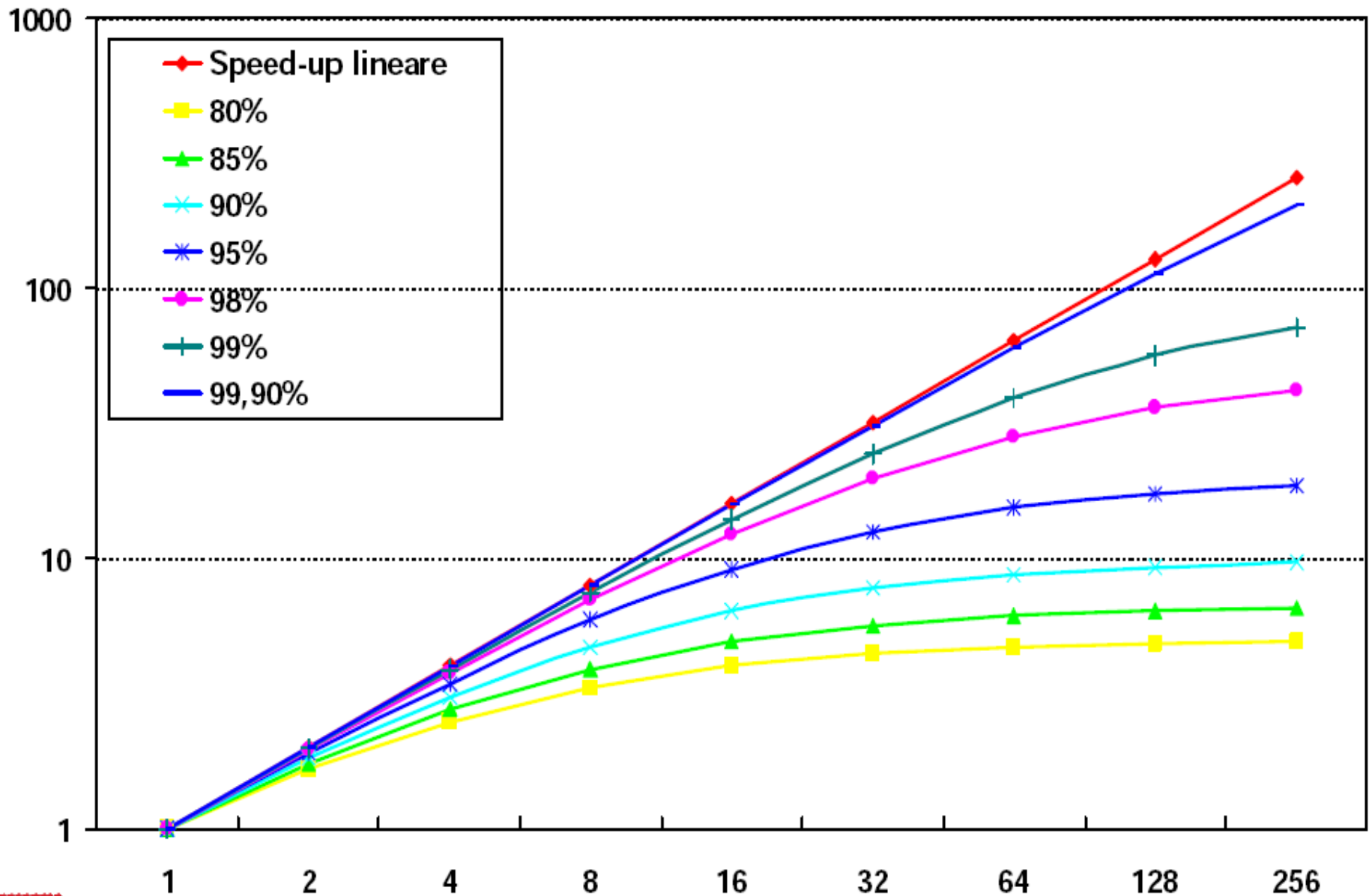
Quando  $n$  tende all'infinito otteniamo:

$$S_n = 1 / F$$

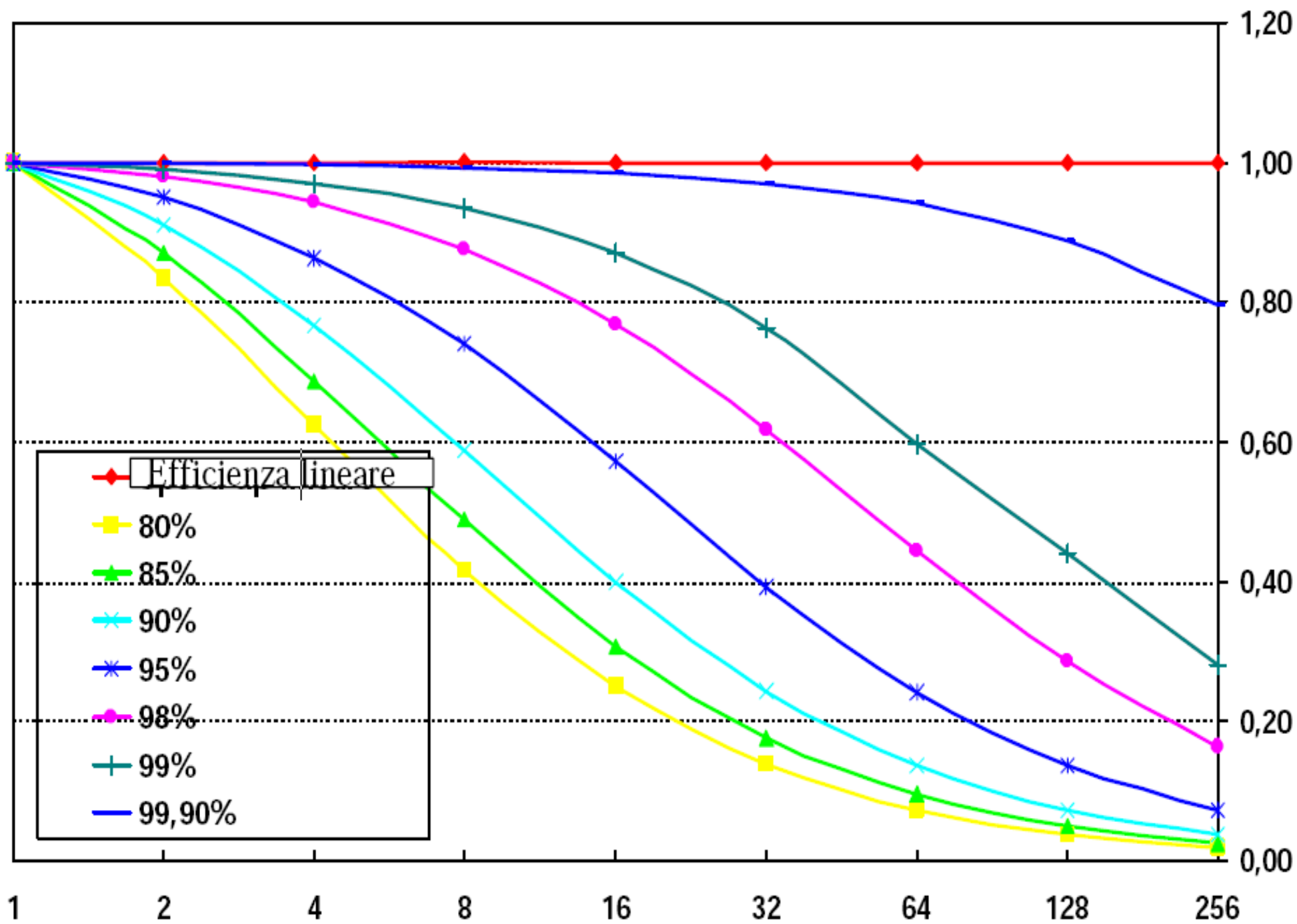
che e' il massimo speedup che possiamo aspettarci, a prescindere dal numero di processori nel caso in cui riusciamo a parallelizzare una frazione  $(1 - F)$  del codice in questione.

Se la percentuale di codice seriale e' il 10% ( $F = 0.10$ ) il massimo speedup ottenibile sara' 10 a prescindere dal numero di processori utilizzato. Dunque in generale se ne deduce che il parallelismo e' utile o quando si usa un piccolo numero di processori, o per quei problemi cosi' detti embarrassingly parallel.





Chiaramente se  $F$  tende a zero  $S_n$  tende ad  $n$



**Altri Costi**

# Load Balancing

Se dividiamo un lavoro in tanti "task" piu' piccoli e li eseguiamo in parallelo dobbiamo alla fine aspettare fino che tutti i "task" non sono stati completati.

La velocita' del lavoro e' dunque determinata da quello del "task" piu' lento. Tutti i task devono essere disegnati di modo da avere piu' o meno lo stesso tempo di esecuzione. Se cosi' non fosse, potremmo trovarci nella situazione di avere qualche processore inutilmente in "idle".

Ottenere un buon "load balancing" non e' sempre semplice.

# Sincronizzazione

In tutti i sistemi shared memory lo scambio di dati tra i processori avviene mediante variabili memorizzate nella memoria condivisa.

Chiaramente spesso si pone il problema della sincronizzazione, cioè dell'accesso simultaneo in lettura e scrittura di tali variabili, da parte di più processori.

Questo problema di sincronizzazione è tipicamente gestito mediante meccanismi di "lock-unlock". Questo comporta in generale una certa perdita di tempo.

Altri problemi sorgono a causa della necessità di mantenere la coerenza della "cache" (Questo spesso comporta una certa perdita dei vantaggi dati dalla presenza stessa della cache)

# Tempi di comunicazione (memoria distribuita)

Rispetto alle macchine seriali i tempi di comunicazione sono fondamentali, in particolare per i cluster:

<b>Nome</b>	<b>Vendore</b>	<b>CPU</b>	<b>N. CPU</b>	<b>Peak (GF)</b>	<b>Sustained (GF)</b>
K.I.S.T	IBM	Intel Xeon (2.4 Ghz)	1024	4915	3067
TotalFinalElf	IBM	Intel Xeon (2.4 Ghz)	1024	4915	1755

Dove e' la differenza ? Esclusivamente nelle rete di interconnessione:

1. Myrinet
2. Gigabit Ethernet

Parametri importanti: latenza, banda,

# Limiti della legge di Amdahl

La legge di Amdahl (1967) pare mettere un limite considerevole allo sviluppo del calcolo parallelo. Tuttavia da quando e' stata formulata i fatti hanno mostrato che anche altri fattori vanno considerati:

La legge di Amdahl si basa su alcuni presupposti che non sono sempre del tutto veri:

- .basta pensare a casi di speedup superscalari (cumulative cache size)
- .la legge di Amdahl parte dal presupposto che l'argoritmo seriale sia sempre e comunque la miglior soluzione al problema. Alcuni problemi si prestano meglio ad essere risolti con algoritmi paralleli
- .**la legge di Amdahl assume che la dimensione del problema rimanga la stessa al crescere del numero di processori. Nella maggior parte dei casi piu' processori implicano anche la possibilita' di affrontare problemi piu' grandi.**

# La legge di Gustafson

Questa legge e' datata 1988. Indichiamo con  $s$  il tempo impiegato all'esecuzione del codice sequenziale, mentre con  $p$  il tempo impegnato nella parte parallela. Se usiamo  $n$  processori immaginiamo anche di crescere la dimensione del problema di  $n$  volte. Al crescere di  $n$  supponiamo che  $s$  rimanga fisso (ad esempio  $s$  e' il tempo di inizializzazione del task), allora avro':

$$\begin{aligned}T_1 &= s + n*p \\ T_p(n) &= s + p\end{aligned}$$

dunque:

$$S_n = (s + n*p) / (s + p)$$

Dove  $S_n$  e' lo speedup per  $n$  processori. (Se usiamo  $F = s / (s + n * p)$  riotteniamo la legge di Amdahl)



# La legge di Gustafson

Quando  $n$  cresce  $S_n$  cresce linearmente. Se  $s$  tende a zero allora otteniamo:

$$S_n = n$$

Gustafson's law argues that even using massively parallel computer systems does not influence the serial part and regards this part as a constant one.

In comparison to that, the hypothesis of Amdahl's law results from the idea that the influence of the serial part grows with the number of processes.

(vedi **gust**)

# Esercizio

Provate a determinare lo speedup del primo programma MPI che abbiamo visto assieme.