

# Speedup

---

Thoai Nam

Faculty of Computer Science and Engineering

HCMC University of Technology



# Outline

---

- ❑ Speedup & Efficiency
- ❑ Amdahl's Law
- ❑ Gustafson's Law
- ❑ Sun & Ni's Law



# Speedup & Efficiency

---

## □ Speedup:

$$S = \text{Time}(\text{the most efficient sequential algorithm}) / \text{Time}(\text{parallel algorithm})$$

## □ Efficiency:

$$E = S / N \quad \text{with } N \text{ is the number of processors}$$



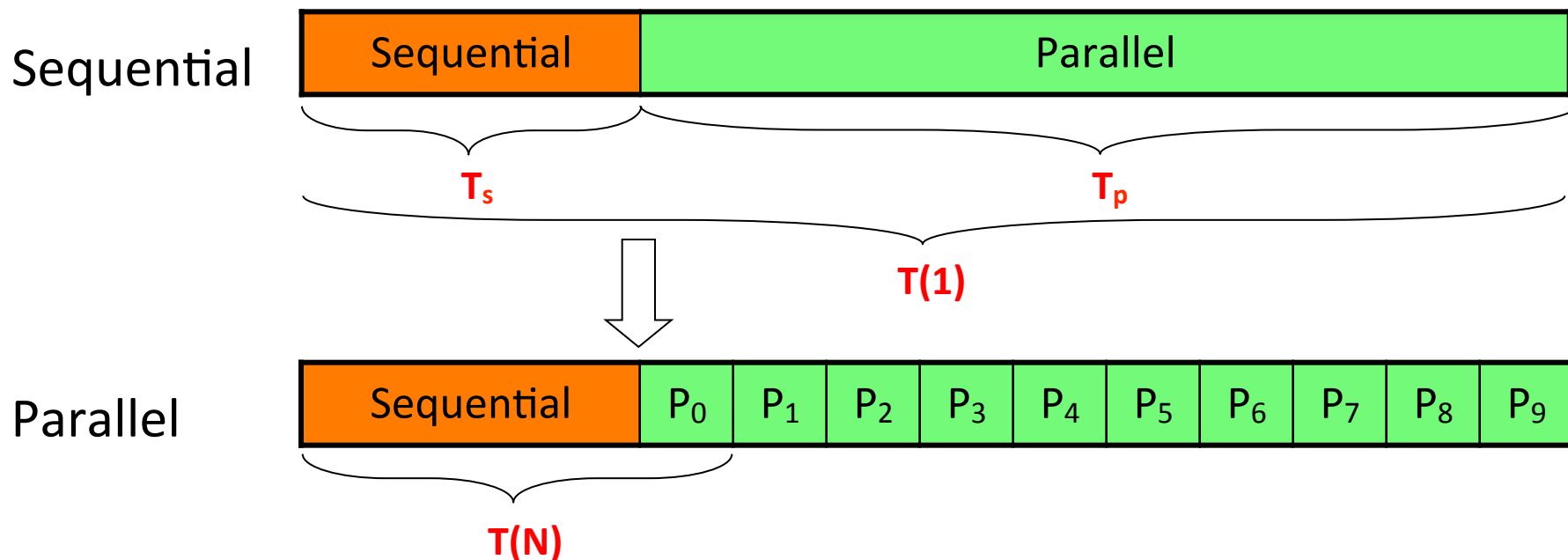
# Amdahl's Law – Fixed Problem Size (1)

---

- ❑ The main objective is to produce the results as soon as possible
  - (ex) video compression, computer graphics, VLSI routing, etc
- ❑ Implications
  - Upper-bound is
  - Make Sequential bottleneck as small as possible
  - Optimize the common case
- ❑ Modified Amdahl's law for **fixed problem size** including the overhead



# Amdahl's Law – Fixed Problem Size (2)



$$T_s = \alpha T(1) \Rightarrow T_p = (1 - \alpha) T(1)$$

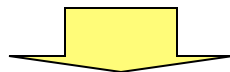
$$T(N) = \alpha T(1) + (1 - \alpha) T(1) / N$$

Number of  
processors



# Amdahl's Law – Fixed Problem Size (3)

$$Speedup = \frac{Time(1)}{Time(N)}$$



$$Speedup = \frac{T(1)}{\alpha T(1) + \frac{(1-\alpha)T(1)}{N}} = \frac{1}{\alpha + \frac{(1-\alpha)}{N}} \rightarrow \frac{1}{\alpha} \text{ as } N \rightarrow \infty$$



# Enhanced Amdahl's Law

The overhead includes parallelism  
and interaction overheads

$$Speedup = \frac{T(1)}{\alpha T(1) + \frac{(1-\alpha)T(1)}{N} + T_{overhead}} \rightarrow \frac{1}{\alpha + \frac{T_{overhead}}{T(1)}} \text{ as } N \rightarrow \infty$$



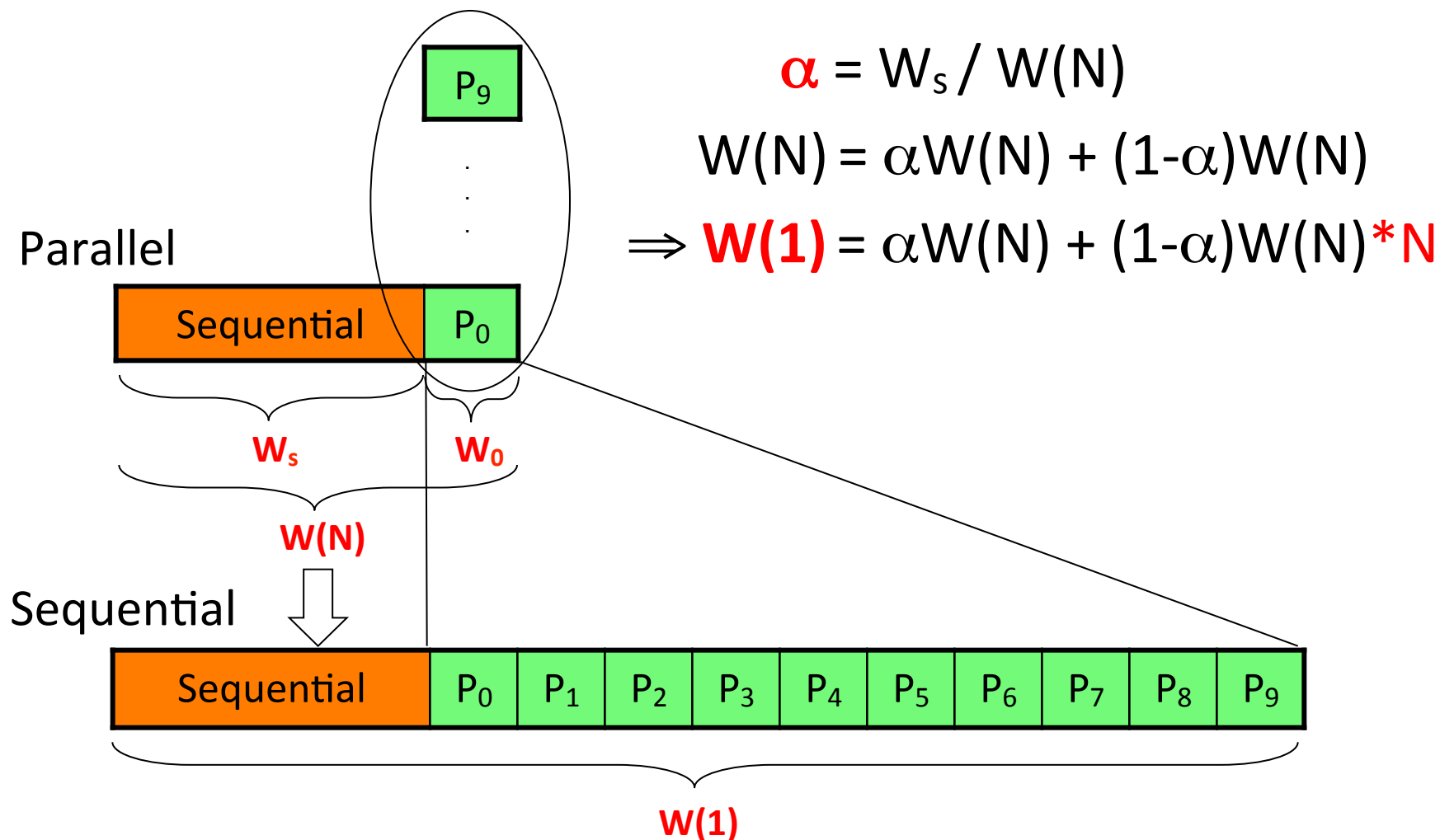
# Gustafson's Law – Fixed Time (1)

---

- ❑ User wants more accurate results within a time limit
    - Execution time is fixed as system scales
    - (ex) FEM (Finite element method) for structural analysis, FDM (Finite difference method) for fluid dynamics
  - ❑ Properties of a work metric
    - Easy to measure
    - Architecture independent
    - Easy to model with an analytical expression
    - No additional experiment to measure the work
    - The measure of work should scale linearly with sequential time complexity of the algorithm
  - ❑ Time constrained seems to be most generally viable model!
-



# Gustafson's Law – Fixed Time (2)





# Gustafson's Law – Fixed Time without overhead

---

Time = Work \* k

$W(N) = W$

$$Speedup = \frac{T(1)}{T(N)} = \frac{W(1) * k}{W(N) * k} = \frac{\alpha W + (1 - \alpha)NW}{W} = \alpha + (1 - \alpha)N$$



# Gustafson's Law – Fixed Time with overhead

---

$$W(N) = W + W_0$$

$$\text{Speedup} = \frac{T(1)}{T(N)} = \frac{W(1) * k}{W(N) * k} = \frac{\alpha W + (1 - \alpha)NW}{W + W_0} = \frac{\alpha + (1 - \alpha)N}{1 + \frac{W_0}{W}}$$



# Sun and Ni's Law – Fixed Memory (1)

---

- ❑ Scale the largest possible solution limited by the memory space. Or, fix memory usage per processor
- ❑ Speedup
  - $\text{Time}(1)/\text{Time}(N)$  for scaled up problem is not appropriate
  - For simple profile, and  $G(N)$  is the increase of parallel workload as the memory capacity increases  $N$  times



## Sun and Ni's Law – Fixed Memory (2)

- $W = \alpha W + (1 - \alpha)W$
- Let  $M$  be the memory capacity of a single node
- $N$  nodes:
  - the increased memory  $N * M$
  - The scaled work:  $W = \alpha W + (1 - \alpha)W * G(N)$

$$Speedup_{MC} = \frac{\alpha + (1 - \alpha)G(N)}{\alpha + (1 - \alpha)\frac{G(N)}{N}}$$



# Sun and Ni's Law – Fixed Memory (3)

## □ Definition:

A function  $g$  is homomorphism if there exists a function such that  $\bar{g}$  for any real number  $c$  and variable  $x$ ,

$$g(cx) = \bar{g}(c) * g(x)$$

## □ Theorem:

If  $W = g(M)$  for some homomorphism function  $g$ , then with all data being shared by all available processors, the simplified memory-bounded speedup is

$$S_N^* = \frac{W_1 + \bar{g}(N)W_N}{W_1 + \frac{\bar{g}(N)}{N}W_N} = \frac{\alpha + (1 - \alpha)G(N)}{\alpha + (1 - \alpha)\frac{G(N)}{N}}$$



# Sun and Ni's Law – Fixed Memory (4)

Proof:

Let the memory requirement of  $W_n$  be  $M$ ,  $W_n = g(M)$ .

$M$  is the memory requirement when 1 node is available.

With  $N$  nodes available, the memory capacity will increase to  $N * M$ .

Using all of the available memory, for the scaled parallel portion  $W_N^*$ :  $W_N^* = g(N * M) = \bar{g}(N) * g(M) = \bar{g}(N) * W_N$

$$S_N^* = \frac{W_1^* + W_N^*}{W_1^* + \frac{W_N^*}{N}} = \frac{W_1 + \bar{g}(N)W_N}{W_1 + \frac{\bar{g}(N)}{N}W_N}$$



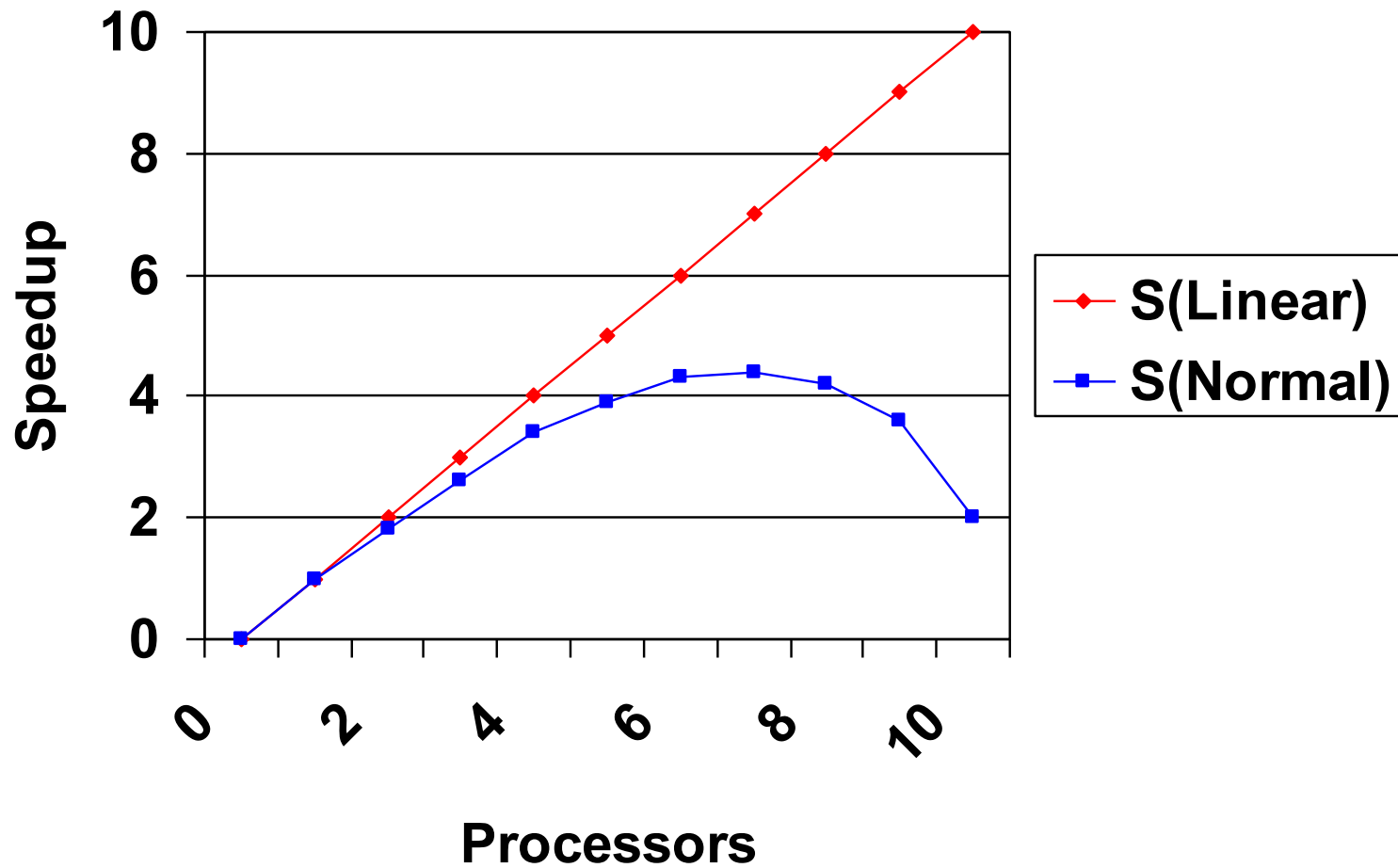
# Speedup

$$S_N^* = \frac{W_1 + G(N)W_N}{W_1 + \frac{G(N)}{N}W_N}$$

- When the problem size is independent of the system, the problem size is fixed,  $G(N)=1 \Rightarrow$  **Amdahl's Law**.
- When memory is increased N times, the workload also increases N times,  $G(N)=N \Rightarrow$  **Gustafson's Law**
- For most of the scientific and engineering applications, the computation requirement increases faster than the memory requirement,  $G(N)>N$ .



# Examples





# Scalability

---

- ❑ Parallelizing a code does not always result in a speedup; sometimes it actually slows the code down! This can be due to a poor choice of algorithm or to poor coding
- ❑ The best possible speedup is **linear**, i.e. it is proportional to the number of processors:  $T(N) = T(1)/N$  where **N = number of processors**, **T(1) = time for serial run**.
- ❑ A code that continues to speed up reasonably close to linearly as the number of processors increases is said to be **scalable**. Many codes scale up to some number of processors but adding more processors then brings no improvement. Very few, if any, codes are indefinitely scalable.



# Factors That Limit Speedup

---

## ☐ Software overhead

Even with a completely equivalent algorithm, software overhead arises in the concurrent implementation. (e.g. there may be additional index calculations necessitated by the manner in which data are "split up" among processors.) i.e. there is generally more lines of code to be executed in the parallel program than the sequential program.

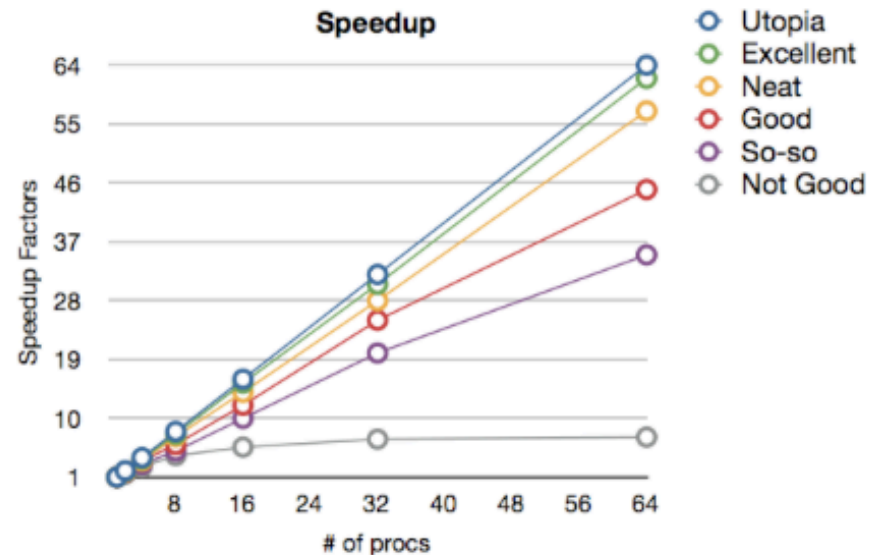
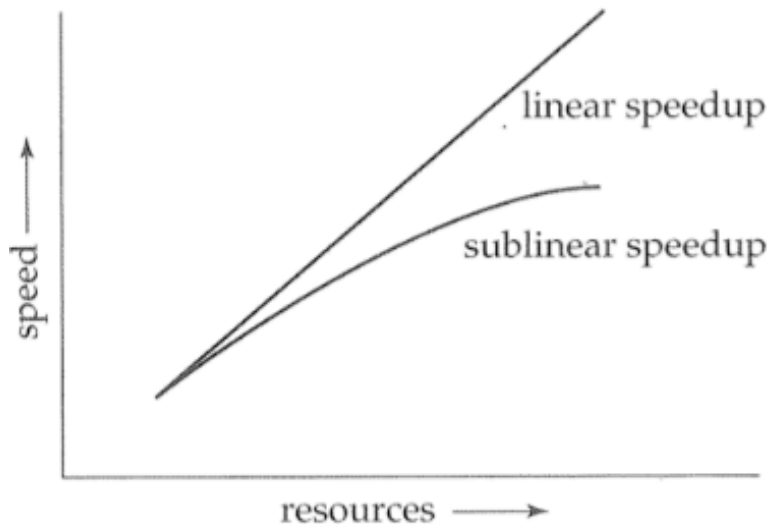
## ☐ Load balancing

## ☐ Communication overhead

# Speedup

- The fundamental concept in parallelism

- $T(1)$  = time to execute task on a single resource
- $T(n)$  = time to execute task on  $n$  resources
- $\text{Speedup} = T(1)/T(n)$



[http://web.eecs.utk.edu/~huangj/hpc/hpc\\_intro.php](http://web.eecs.utk.edu/~huangj/hpc/hpc_intro.php)

From Silberschatz, Korth, and Sudarshan. *Database Systems Concepts*, 4<sup>th</sup> Ed.