# From Amdahl's Law to Big Data:
# *A Story of*
# Computer Sciences and Technology

## Xian-He Sun

Illinois Institute of Technology
sun@iit.edu

# The Journey of Supercomputing

- The Background of Parallel Processing
  - Speedup
  - Sources of overhead

- The Laws of Scalable Computing
  - The Amdahl's law
  - The Gustafson's law
  - The Sun-Ni's law

- Impacts and Discussions

# Performance Evaluation
## *(Improving performance is the goal)*

- **Performance Measurement**
  - Metric, Parameter

- **Performance Prediction**
  - Model, Application-Resource

- **Performance Diagnose/Optimization**
  - Post-execution, Algorithm improvement, Architecture improvement, State-of-the-art, Scheduling, Resource management/Scheduling

# Performance of Parallel Processing

## Models of Speedup

- Speedup

  - Ts = time for the best serial algorithm
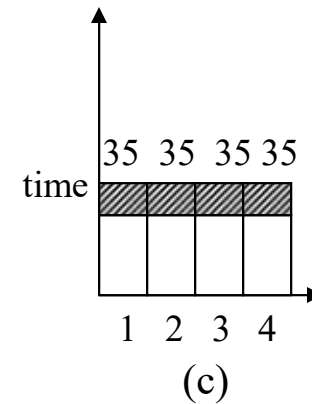
  - Tp= time for parallel algorithm using p processors
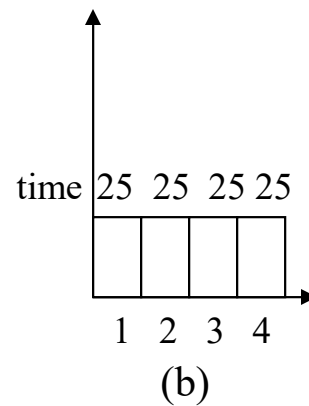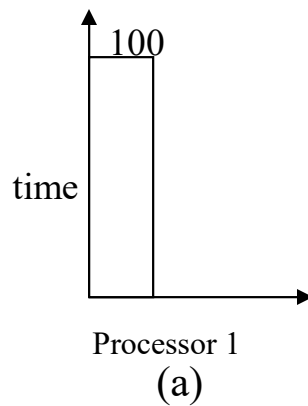
$$S_p = \frac{T_s}{T_p}$$

- Simple enough, but also unexpected complex

$$S_p = \frac{\text{Uniprocessor Execution Time}}{\text{Parallel Execution Time}}$$

# Example
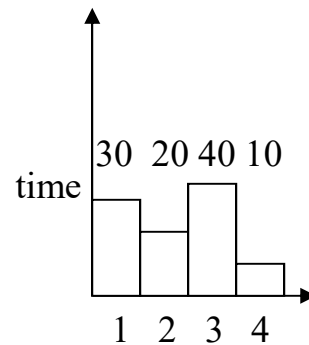


(a)

(b)

(c)

$$S_p = \frac{100}{25} = 4.0,$$

perfect parallelization
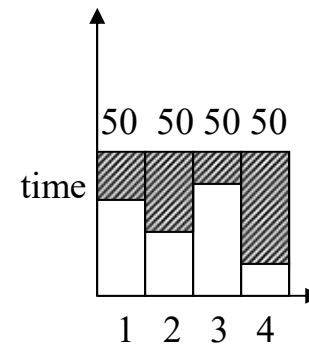
$$S_p = \frac{100}{35} = 2.85,$$

perfect load balancing
but synch cost is 10

# Example (cont.)



(d)

$$S_p = \frac{100}{40} = 2.5,$$

no synch

but load imbalance



(e)

$$S_p = \frac{100}{50} = 2.0,$$

load imbalance

and synch cost

# What Is "Good" Speedup?

- *Linear* speedup:

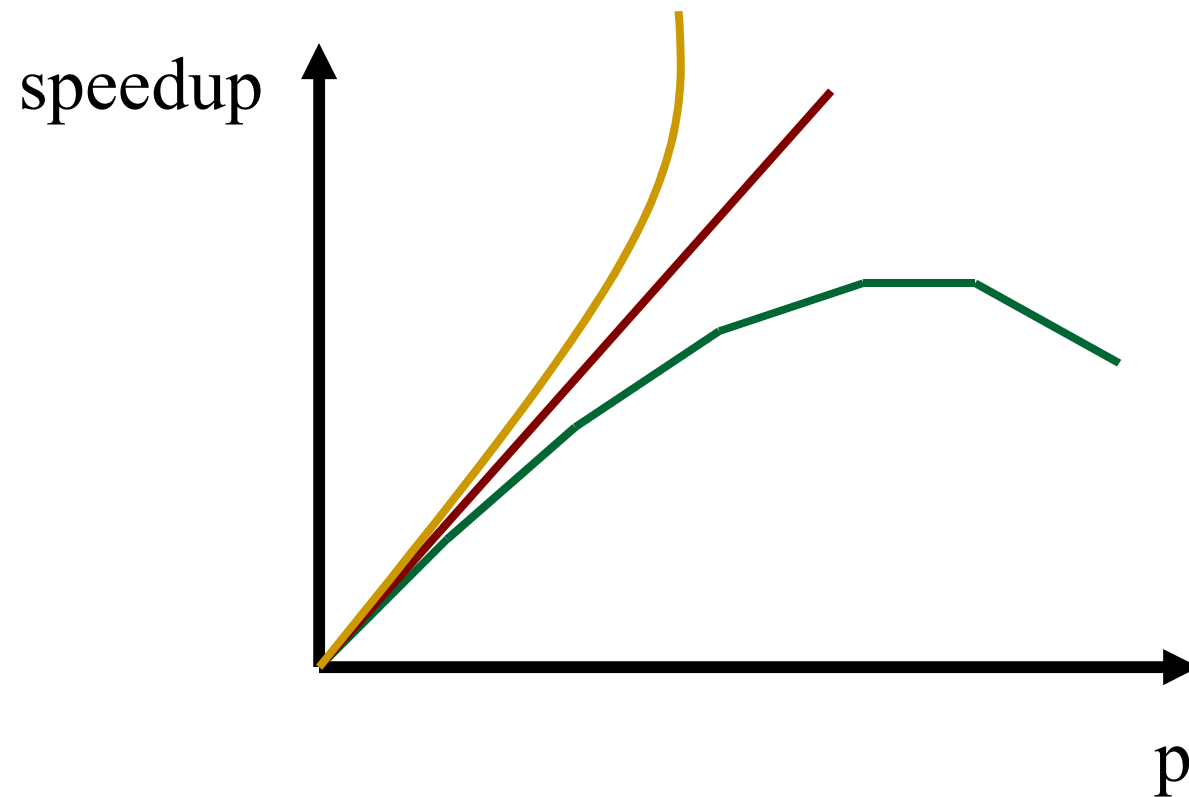$$S_p = p$$

- *Superlinear* speedup

$$S_p > p$$

- *Sub-linear* speedup:

$$S_p < p$$

# Speedup

# Sources of Parallel Overheads

- Interprocessor communication

- Load imbalance

- Synchronization
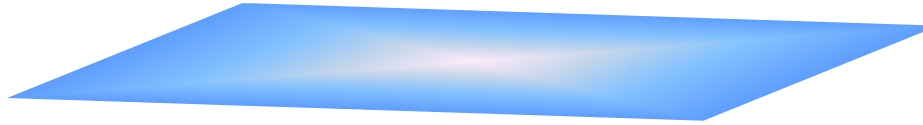
- Extra computation

# Causes of Superlinear Speedup

- Cache size increased

- Overhead reduced

- Latency hidden

- Randomized algorithms

- Mathematical inefficiency of the serial algorithm

- Higher memory access cost in sequential processing
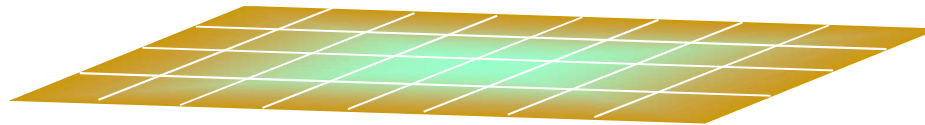
• X.H. Sun, and J. Zhu, "Performance Considerations of Shared Virtual Memory Machines," *IEEE Trans. on Parallel and Distributed Systems*, Nov. 1995

# Degradations of Parallel Processing
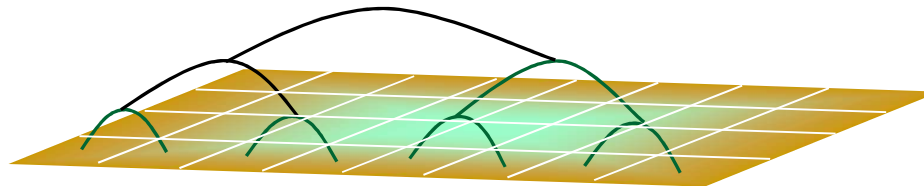
**Unbalanced Workload**

**Communication Delay**

**Overhead Increases with the Ensemble Size**

**Overheads**

- *communication*
- *Load imbalance*
- *Synchronization*
- *Extra computation*

# Degradations of Distributed Computing

Unbalanced Computing Power and Workload

Shared Computing and Communication Resource

Uncertainty, Heterogeneity, and Overhead Increases
with the Ensemble Size

# Principals of Architecture Design


Gene Amdahl

- Make common case fast (90/10 Rule)
- Amdahl's Law
  - Law of diminishing returns
- Speedup
  - Achieved performance improvement over original

$$\text{Speedup Overall} = \frac{\text{speed new}}{\text{speed old}} = \frac{\text{execution time old}}{\text{execution time new}}$$

Here performance is measured in **Speed**

# Amdahl's Law

Execution time of any code has two portions

Portion I:  not affected by enhancement
Portion II: affected by enhancement

$$\text{execution time}_{old} = \text{execution time}_{p1} + \text{execution time}_{p2}$$

$\alpha$ is % of original code that cannot benefit from enhancement

As p -> infinity, execution time$_{new}$ -> $\alpha$ * execution time$_{old}$

$$\text{execution time}_{new} = (\alpha) * \text{execution time}_{old} + (1-\alpha) * \frac{\text{execution time}_{old}}{p}$$

Execution time$_{p1}$          Execution time$_{p2}$

$p$ is speedup factor of old/new
execution times for portion II

# Amdahl's Law for Parallel Processing (1967)

- Let $\alpha$ = fraction of program (algorithm) that is <u>serial</u> and <u>cannot be parallelized</u>. For instance:
  - ☐ Loop initialization
  - ☐ Reading/writing to a single disk
  - ☐ Procedure call overhead
- Parallel run time is given by

$$\text{execution time}_{new} = (\alpha) * \text{execution time}_{old} + (1-\alpha) * \frac{\text{execution time}_{old}}{p}$$

$$T_p = (\alpha + \frac{1-\alpha}{p}) \bullet T_s$$

*Gene M Amdahl, "Validity of the single processor approach to achieving large scale computing capabilities," AFIPS spring joint computer conference, 1967*

# Amdahl's Law

- Amdahl's law gives a limit on speedup in terms of $\alpha$

$$S_p = \frac{T_s}{T_p} = \frac{T_s}{\alpha T_s + \dfrac{(1-\alpha)T_s}{p}} = \frac{1}{\alpha + \dfrac{1-\alpha}{p}}$$

- If we assume that the serial fraction is fixed, then the speedup for infinite processors is limited by $1/\alpha$

$$\lim_{p->\infty} S_p = \frac{1}{\alpha}$$

- For example, **if $\alpha$=10%,** then the maximum speedup is **10,** even if we use an infinite number of processors

# Amdahl Law

- The sequential part becomes the dominate factor quickly

# Amdahl's Law

$$\text{execution time}_{new} = (\alpha) * \text{execution time}_{old} + (1-\alpha) * \frac{\text{execution time}_{old}}{p}$$

Example: alpha = 20%



Speedup factor = p

$$\text{Speedup}_{overall} = \frac{\text{execution time}_{old}}{\text{execution time}_{new}} = \frac{1}{(\alpha) + \frac{1-\alpha}{p}}$$

# Amdahl's Law with Overhead

- To include overhead will be even worse
- The overhead includes parallelism and interaction overheads

$$Speedup_{FS} = \frac{T1}{\alpha T1 + \dfrac{(1-\alpha)T1}{p} + T_{overhead}} \rightarrow \frac{1}{\alpha + \dfrac{T_{overhead}}{T1}} \; as \; p \rightarrow \infty$$

Amdahl's law: argument against massively parallel systems

# Example

- Enhanced mode is used 50% of resulting execution time$_{new}$
- Portions of code using enhanced mode improve performance by factor 10x
- What is speedup for fast mode?

First find 1-α (% of code affected by enhancement)

$$\text{execution time}_{new} = \underbrace{(.5) * \text{execution time}_{new}}_{\text{not enhanced}} + \underbrace{(.5) * \text{execution time}_{new}}_{\text{enhanced}}$$

$$(\alpha) * \text{execution time}_{old} = \frac{(1-\alpha) * \text{execution time}_{old}}{10}$$

$$1 - \alpha = .909091 \qquad \text{Speedup}_{overall} = \cfrac{1}{(\alpha) + \cfrac{1-\alpha}{n}} = \cfrac{1}{(1-.9) + \cfrac{.9}{10}} = 5.3$$

# Example

- Enhancement: Parallel processing with 20 computing nodes
- Portions of code containing computations run 20x faster in parallel processing mode.
- What % of original code must be parallelizable to achieve $\text{speedup}_{overall} = 2$?

$$\text{Speedup}_{overall} = \frac{\text{execution time}_{old}}{\text{execution time}_{new}} = \frac{1}{(\alpha) + \frac{1-\alpha}{n}}$$

$$2 = \frac{1}{(\alpha) + \frac{1-\alpha}{20}} \qquad 1 - \alpha = .5263$$

# History back to 1988


IBM 7030 Stretch


IBM 7950 Harvest

All have up to 8 processors, citing Amdahl's law,

$$\lim_{p \to \infty} Speedup_{Amdahl} = \frac{1}{\alpha}$$


Cray X-MP
Fastest computer 1983-1985


Cray Y-MP


Gene Amdahl

# Summit: the World Fastest Computer



- ➢ 148.6 petaflops (187.66 petaflop theoretical peak)
- ➢ 2,282,544 IBM Power 9 core
- ➢ 2,090,880 Nvidia Volta GV100 core
- ➢ Power efficiency 11.324gigaflop/watt

# Bombshell: *Gustafson, etc. Got Speedup of more than 1,000 on Three Applications*

- On a 1024-processor nCUBE parallel computer
- For three applications: wave mechanics, fluid dynamics, and structural analysis.
- Introduced the concept of **Scalable Computing,** *problem size increases with the machine size*

*John L. Gustafson, Gary R. Montry, and Robert E. Benner, "Development of Parallel Methods for a 1024-Processor Hypercube," SIAM Journal on Scientific and Statistical Computing, Vol. 9, No.4, 1988* (submitted 3/10/1988, accepted 3/25/1988, appeared April 1988)

*John Gustafson, "Reevaluation of Amdahl's Law," Communications of the ACM, Vol. 31, No. 5, May 1988.*

# Reevaluate Amdahl's Law

- **Amdahl's Law** is designed for technology improvement, but has been widely used to against parallel processing in terms of reducing execution time

- **But**: large computers are not (only) designed for solving existing problem faster, they are designed for solving otherwise unsolvable large problems

- The introduction of **scalable computing,** where *problem size increases with the machine size*

- Fixed-Time Speedup (Gustafson, 88)

  o    Emphasis on work finished in a fixed time

  o    Problem size is **scaled** from $W$ to $W'$

  o    $W'$: Work finished within the fixed time with

  parallel processing

John L. Gustafson

$$S'_p = \frac{\text{Uniprocessor Time of Solving } W'}{\text{Parallel Time of Solving } W'}$$

$$= \frac{\text{Uniprocessor Time of Solving } W'}{\text{Uniprocessor Time of Solving } W}$$

$$= \frac{W'}{W}$$

# Fixed-Time Speedup (Gustafson)

■ Solving a larger application within the time limit



Amount of Work

Number of Processors (p)

Elapsed Time

Number of Processors (p)

# Reexam Amdahl Law (Fixed-Size Speedup)

- It is on time reduction for solving a fixed problem (size)

- Amdahl's law (Fixed-Size Speedup)
  - Emphasis on turnaround time
  - Problem size, $W$, is fixed

$$S_p = \frac{\text{Uniprocessor Execution Time}}{\text{Parallel Execution Time}}$$

$$S_p = \frac{\text{Uniprocessor Time of Solving } W}{\text{Parallel Time of Solving } W}$$

# Gustafson's Law (Without Overhead)

- Under **Gustafson's Law** the parallel processing part is changing with the number of processors, *p,* and problem size
- Linear speedup



$$\alpha = \frac{t_s}{t_s + t_p}$$

$$Speedup_{FT} = \frac{Work(p)}{Work(1)} = \frac{\alpha W + (1 - \alpha)pW}{W} = \alpha + (1 - \alpha)p$$

If α=0.1

$$Speedup_{FT} = \alpha + (1 - \alpha)p = 0.1 + 0.9p$$

# Gustafson's Law (With Overhead)

- In practice, the **overhead** can increase with P, and **limit** the scalability

$$Speedup_{FT} = \frac{Work(p)}{Work(1)} = \frac{\alpha + (1-\alpha)p}{1 + T_{overhead}/T_1}$$

$$Speedup_{FT} = \frac{Work(p)}{Work(1)} = \alpha + (1 - \alpha - \frac{T_{overhead}}{T_1})p$$

# But: *Gustafson's Applications are not Scalable*

- Most applications cannot get more than 1,000 speedup on a 1024-processor nCUBE parallel computer

  *Parallel Processing overhead*

- Even the three applications are not **Scalable** (increase *problem size further does not help*)

  *Why?*

# Memory Constrained Scaling:
## *Sun and Ni's Law*

- **Scaling is limited by memory space** (disk will increase overhead significantly), e.g. fixed memory capacity/usage per processor
  - (ex) N-body problem
- Problem size is scaled from W to W*, W* is the work executed under memory limitation
- The relation between memory & computing requirement is determined by the underlying algorithm/program
- **Memory-scaling function**

$$W* = G(p * M)$$

# Sun & Ni's Law

## 存储受限理论

Xian-He Sun

Lionel M. Ni



$(1-a)G(p)$

$$Speedup_{MB} = \frac{Work(p)/Time(p)}{Work(1)/Time(1)} = \frac{\alpha + (1-\alpha)G(p)}{\alpha + (1-\alpha)G(p)/p}$$

Assuming $\alpha = 0.1$, the problem needs $2n^3$ computation and $3n^2$ memory

Then $G(p) = G(\text{p}) = p^{\frac{3}{2}}$, and

$$Speedup_{MB} = \left(0.1 + 0.9 \times p^{\frac{3}{2}}\right) \Big/ \left(0.1 + (0.9 \times p^{\frac{3}{2}})/p\right)$$

# Memory-Bounded Speedup 存储受限理论

## (Sun & Ni, 90)

• Emphasis on work finished under current physical limitation

  ° Problem size is scaled from $W$ to $W^*$

  ° $W^*$: Work executed under memory limitation with parallel processing

$$S_p^* = \frac{\text{Uniprocessor Time of Solving } W^*}{\text{Parallel Time of Solving } W^*}$$

# Memory-Bounded Speedup (Sun & Ni)

- In practice, memory-bounded performs better than fixed-time but both hard to achieve linear speedup

$$Speedup_{MB} = \frac{Work(p)/Time(p)}{Work(1)/Time(1)} = \frac{\alpha + (1-\alpha)G(p)}{\alpha + \dfrac{(1-\alpha)G(p)}{p} + overhead(p, G(p))}$$

$$Speedup_{FT} = \frac{Work(p)}{Work(1)} = \alpha + (1 - \alpha - \frac{T_{overhead}}{T_1})p$$

Elapsed Time

| $T_1$ | $T_1$ | $T_1$ | $T_1$ | $T_1$ |
| $T_p$ | $T_p$ | $T_p$ | $T_p$ | $T_p$ |

1   2   3   4   5
Number of Processors (p)

Elapsed Time

| $T_1$ | $T_1$ | $T_1$ | $T_1$ | $T_1$ |
| $T_p$ | $T_p$ | $T_p$ | $T_p$ | $T_p$ |

1   2   3   4   5
Number of Processors (p)

# Example of calculating Sun & Ni's Law
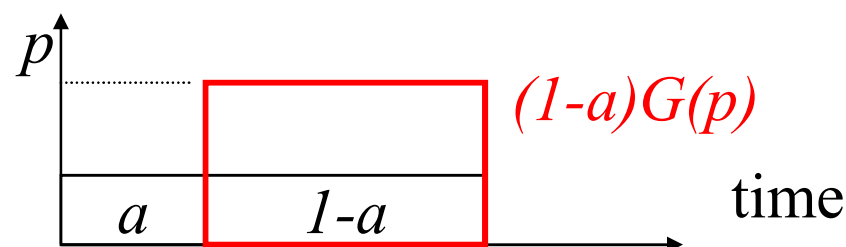
$$Speedup_{MB} = \frac{Work(p)/Time(p)}{Work(1)/Time(1)} = \frac{\alpha + (1-\alpha)G(p)}{\alpha + (1-\alpha)G(p)/p}$$

- Example:
  - Assuming that the problem needs $2n^3$ computation and $3n^2$ memory requirement
  - $\alpha = 0.6$
  - $p = 8$
  - $G(p) = G(8) = 8^{\frac{3}{2}} \approx 22.63$
  - $Speedup_{MB} = (0.6 + 0.4 \times 22.63)/(0.6 + (0.4 \times 22.63)/8) \approx 5.576$

# Rethinking of Speedup

- Speedup

$$S_p = \frac{Uniprocessor\ ExecutionTime}{Parallel\ ExecutionTime}$$

- It is only the true speedup if problem size is fixed, but now we have scalable computing
- Generalized speedup

$$S_p = \frac{Parallel\ Speed}{Sequential\ Speed}$$

**X.H. Sun, and J. Gustafson**, *"Toward A Better Parallel Performance Metric," Parallel Computing*, Vol. 17, pp.1093-1109, Dec. 1991.

# Performance Example

- Performance is measured in **Speed**
- But since they work on the same of work, Speedup equals Time Reduction
- My car (X) travels a distance of 1 in one hour
  - Time between start and completion of event is 1 hour
    - Execution time
- Your car (Y) travels a distance of 1 in two hours
- Intuition: my car is twice as fast as your car

# Performance Example

- My car (X) travels a distance of 1 in one hour
    - Time between start and completion of event is 1 hour
        - Execution time
- Your car (Y) travels a distance of 1 in two hours
- Intuition: my car is twice as fast as your car
- Intuition assumes performance = speed

# Is intuition correct?

$$\frac{\text{execution time of your car (Y)}}{\text{execution time of my car (X)}} = n$$

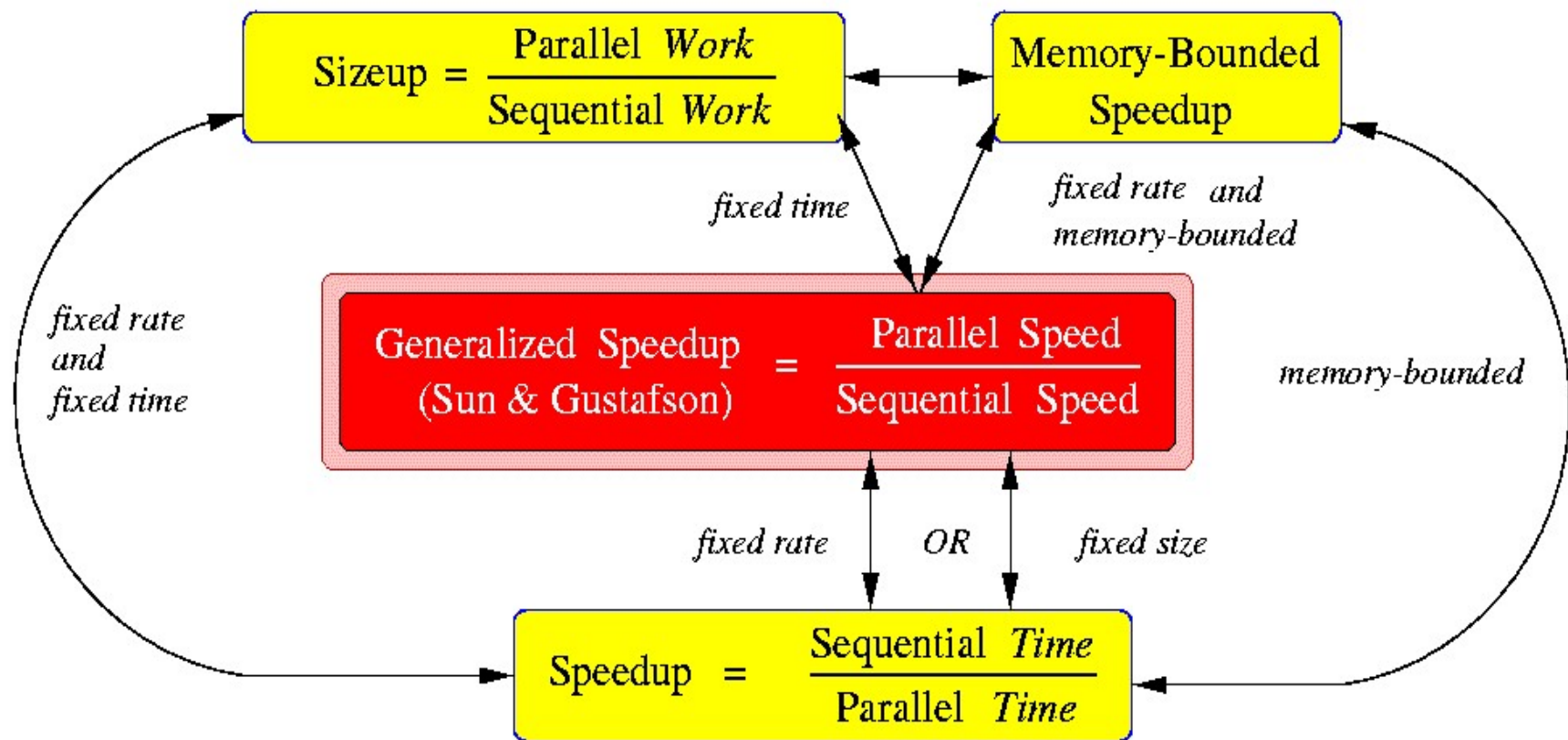Implies X performs n times better than Y

$$\text{performance} = \text{speed} = \frac{1}{\text{execution time}}$$

$$n = \frac{\text{execution time of Y}}{\text{execution time of X}} = \frac{1/\text{speed of Y}}{1/\text{speed of X}} = \frac{\text{speed of X}}{\text{speed of Y}} = \frac{\text{performance of X}}{\text{performance of Y}}$$

Speed is one measure of performance, throughput is another

# Models of Speedup

$$\text{Sizeup} = \frac{\text{Parallel } Work}{\text{Sequential } Work}$$

Memory-Bounded Speedup

fixed time

fixed rate and memory-bounded

$$\text{Generalized Speedup} \atop (\text{Sun \& Gustafson}) = \frac{\text{Parallel Speed}}{\text{Sequential Speed}}$$

fixed rate and fixed time

memory-bounded

fixed rate       OR       fixed size
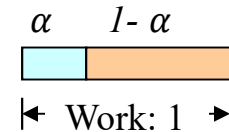
$$\text{Speedup} = \frac{\text{Sequential } Time}{\text{Parallel } Time}$$

# The Three Laws

- Tacit assumption of Amdahl's law

  - Problem size is fixed
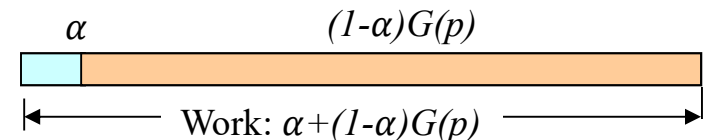
  - Speedup emphasizes on time reduction

$\alpha \quad 1-\alpha$

Work: 1

- Gustafson's Law, 1988

  - Fixed-time speedup model

$\alpha \qquad (1-\alpha)p$

Work: $\alpha+(1-\alpha)p$

$$Speedup_{fixed-time} = \frac{Sequential\ Time\ of\ Solving\ Scaled\ Workload}{Parallel\ Time\ of\ Solving\ Scaled\ Workload}$$

$$= \alpha+(1-\alpha)p$$

$\alpha \qquad (1-\alpha)G(p)$

Work: $\alpha+(1-\alpha)G(p)$

- Sun and Ni's law, 1990

  - Memory-bounded speedup model

$$Speedup_{memory-bound} = \frac{Sequential\ Time\ of\ Solving\ Scaled\ Workload}{Parallel\ Time\ of\ Solving\ Scaled\ Workload}$$

$$= \frac{\alpha + (1-\alpha)G(p)}{\alpha + (1-\alpha)\,G(p)/p}$$

*X.-H. Sun, and L. Ni, "Another View of Parallel Speedup," Proc. of IEEE Supercomputing'90, NY, NY, Nov.12--Nov.16, 1990.*

# *The Three Laws:* and their impact

*I can improve Amdahl's law*

- **Amdahl's law** (1967) shows the inherent limitation of parallel processing

- **Gustafson's law** (scalable computing, 1988) shows there is no inherent limitation for scalable parallel computing, exce engineering issues

*I have a huge memory*

- **Sun-Ni's law** (memory-bounded, 1990) shows memory (data) is the constraint of scalable computing (**the** engineering issue)

- The **Memory-Wall Problem** (1994) shows memory-bound is a general performance issue for computing, not just for parallel computing

# Impact of Scalable Computing



1823
Babbage Difference Engine

1943
Harvard Mark 1

1949
Edsac

1959
IBM 7094

1951
Univac 1

1976
Cray 1

1964
CDC 6600

1982
Cray XMP

1991
Intel Delta

1988
Cray YMP

1996
T3E

1997
ASCI Red

2003
Cray X1

2001
Earth Simulator

1    $10^3$    $10^6$    $10^9$    $10^{12}$    $10^{15}$

One OPS    KiloOPS    MegaOPS    GigaOPS    TeraOPS    PetaOPS