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

## Xian-He Sun

Illinois Institute of Technology
sun@iit.edu

## Asian Science Camp, July 30, 2019

# Hot Issues

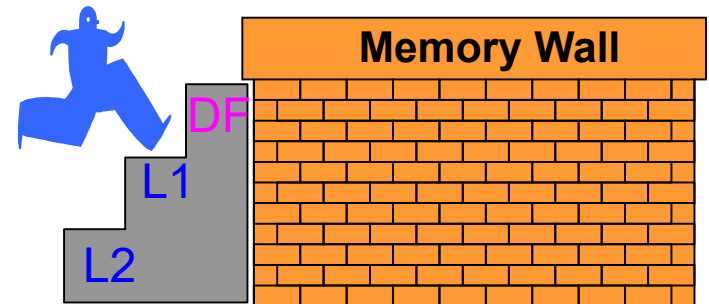- AI and Deep Learning

- Big Data

- *High Performance and Could Computing*

COMPUTING POWER

Memory Wall

DF

L1

L2

# 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

# What is Parallel Processing

- **Parallel Processing**
  - Several working entities work together toward a common goal
- **Parallel Computer**
  - A computer designed for parallel processing
- **Scalable Computing**
  - A parallel computing which can be scaled up to larger size without losing efficiency

- **Supercomputer (high performance computer, high end computer, advanced computer)**
  - A general-purpose computer capable of solving individual problems at extremely high computation speed

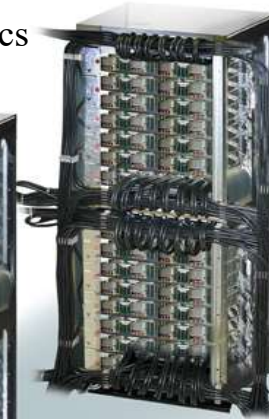# Parallel Processing & Scalable Computing

**Petaflops System**

72 Racks

IBM BG/P

Source: ANL ALCF

**Rack**
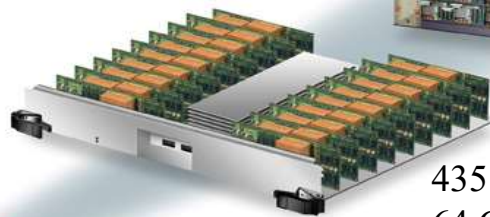
Cabled 8x8x16

32 Node Cards
1024 chips, 4096 procs

1 PF/s
144 TB

**Node Board**

(32 chips  4x4x2)
32 compute, 0-2 IO cards

14 TF/s
2 TB

**Maximum System**

256 racks
3.5 PF/s
512 TB

**Compute Card**

1 chip, 20 DRAMs

435 GF/s
64 GB

**Chip**

4 cores

13.6 GF/s
2.0 GB DDR
Supports 4-way SMP
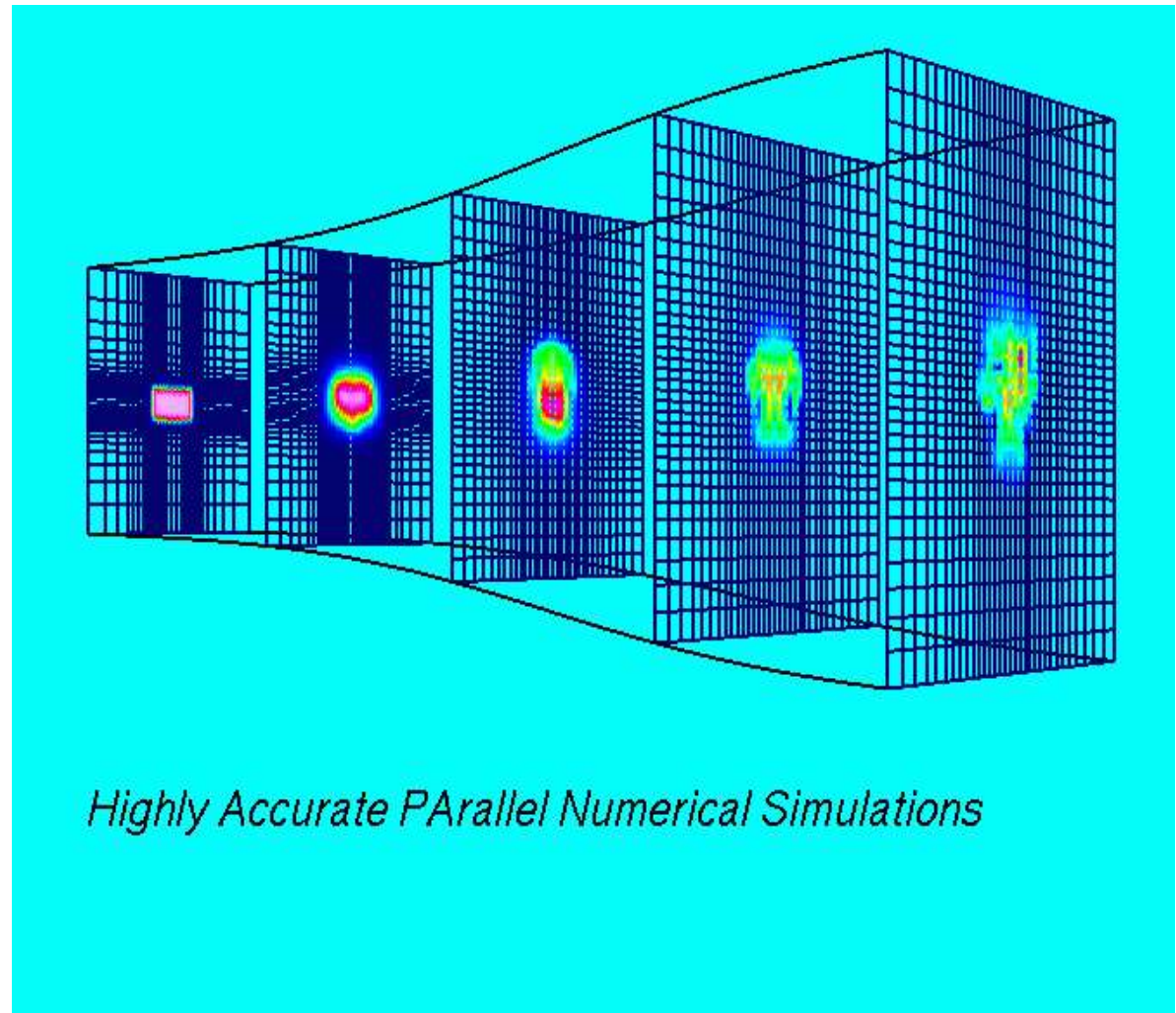
850 MHz
8 MB EDRAM

Front End Node / Service Node

System p Servers

Linux SLES10

HPC SW:

Compilers

GPFS

ESSL

Loadleveler

12/27/2019

# Why Scalable Computing

– Discretization

–Scalable

    More accurate solution

    Sufficient parallelism

    Maintain efficiency

–Efficient in parallel computing

    Load balance

    Communication
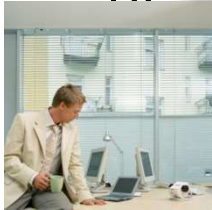
– Mathematically effective

    Adaptive

    Accuracy



Highly Accurate PArallel Numerical Simulations

# Cloud Computing & Big Data

From High Performance Computing to Cloud to Big Data

# 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 of Parallel Processing

## Models of Speedup

- Speedup

  - Ts = time for the best serial algorithm

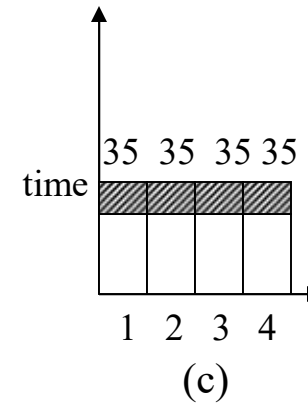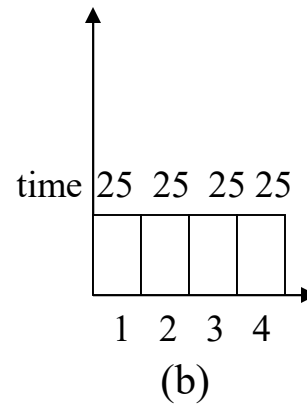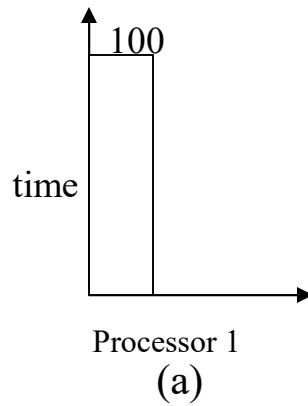  - Tp= time for parallel algorithm using p processors

$$S_p = \frac{T_s}{T_p}$$

- Simple enough, but also unexpected complex

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

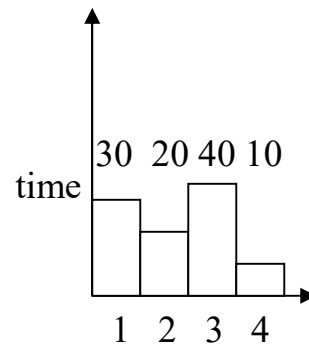# Example



(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)

$$S_p = \frac{100}{40} = 2.5,$$
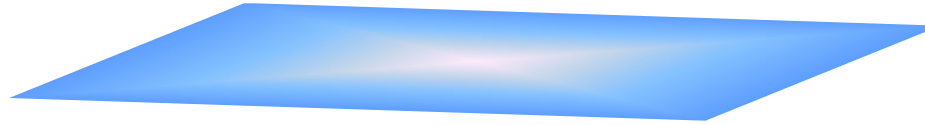
no synch
but load imbalance



(e)

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

load imbalance
and synch cost

# Degradations of Parallel Processing

**Unbalanced Workload**

**Communication Delay**

**Overhead Increases with the Ensemble Size**

**Overheads**

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

# Principals of Architecture Design

- Make common case fast (90/10 Rule)

- Amdahl's Law
  - Law of diminishing returns

- Speedup
  - Achieved performance improvement over original

Gene Amdahl

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

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

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

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

# Amdahl's Law

- Amdahl's law gives a limit on speedup in terms of α

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

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

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

- For example, **if α=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



Amount of Work — $W_1$ $W_1$ $W_1$ $W_1$ $W_1$ / $W_p$ $W_p$ $W_p$ $W_p$ $W_p$ vs. Number of Processors (p): 1 2 3 4 5

Elapsed Time — $T_1$, $T_p$ vs. Number of Processors (p): 1 2 3 4 5

# Amdahl's Law

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

Example: alpha = 20%



Speedup factor = p

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

# Amdahl's Law with Overhead

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

$$Speedup_{FS} = \frac{T_1}{\alpha T_1 + \dfrac{(1-\alpha)T_1}{p} + T_{overhead}} \rightarrow \frac{1}{\alpha + \dfrac{T_{overhead}}{T_1}} \; as \; p \rightarrow \infty$$

Amdahl's law: argument against massively parallel systems

# 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

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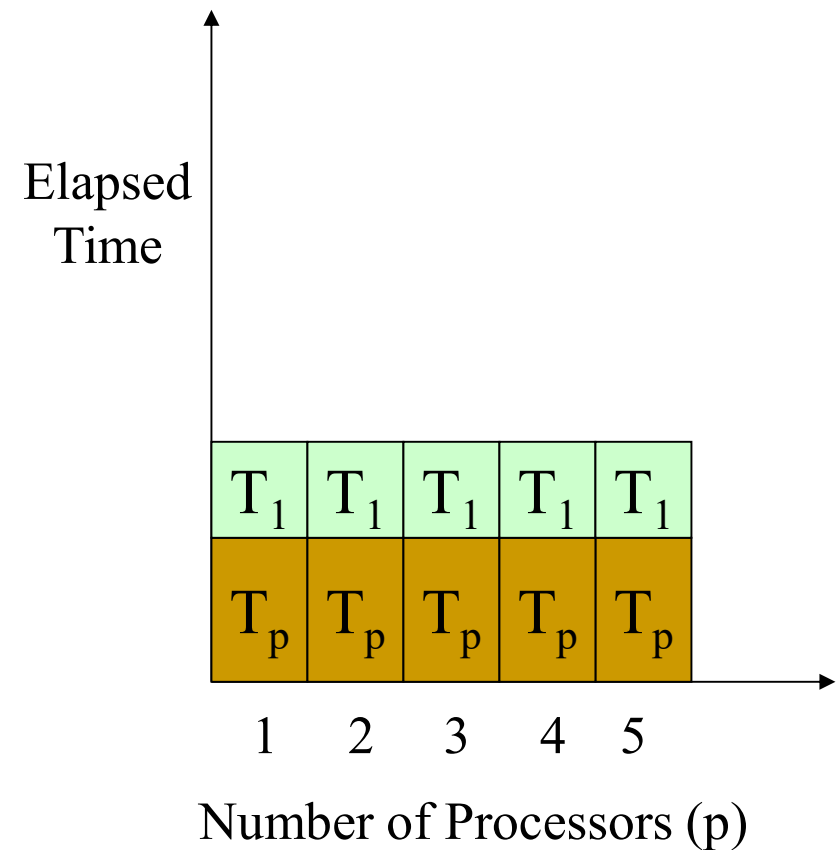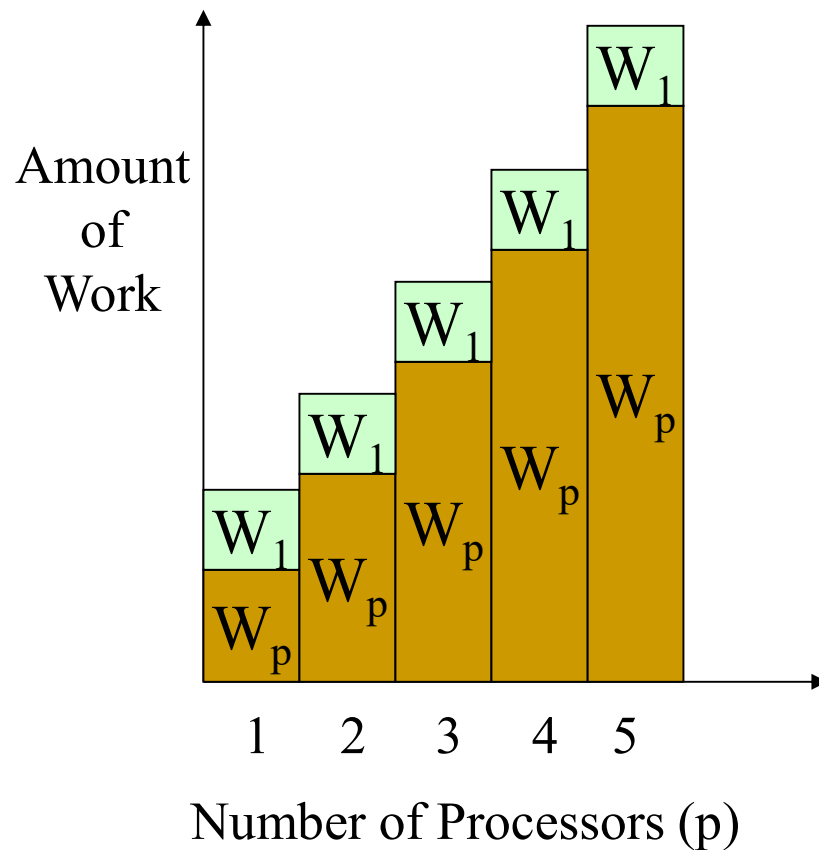
$$S'_p = \frac{\text{Uniprocessor Time of Solving } W'}{\text{Parallel Time of Solving } W'}$$

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

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

# Fixed-Time Speedup (Gustafson)

- Solving a larger application within the time limit



Amount of Work vs. Number of Processors (p); and Elapsed Time vs. Number of Processors (p)

# Reexam Amdahl Law (Fixed-Size Speedup)

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

Amount of Work

| $W_1$ | $W_1$ | $W_1$ | $W_1$ | $W_1$ |
|---|---|---|---|---|
| $W_p$ | $W_p$ | $W_p$ | $W_p$ | $W_p$ |

1   2   3   4   5

Number of Processors (p)

Elapsed Time

$T_1$, $T_p$, $T_1$, $T_p$, $T_1$, $T_p$, $T_1$, $T_p$, $T_1$, $T_p$

1   2   3   4   5

Number of Processors (p)

# 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

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

# 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**

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

X.H. Sun, and L. Ni , *"Scalable Problems and Memory-Bounded Speedup," Journal of Parallel and Distributed Computing*, Vol. 19, pp.27-37, Sept. 1993 (**SC90**).
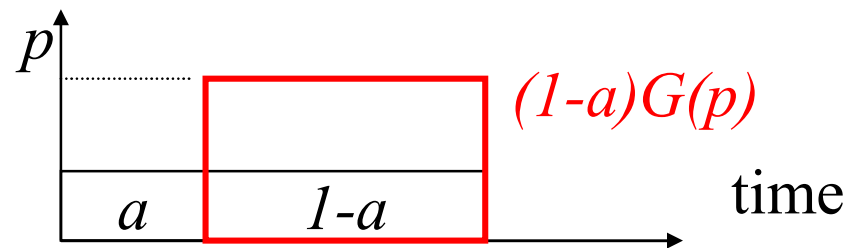
# Sun & Ni's Law

## 存储受限理论

Xian-He Sun

Lionel M. Ni



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

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

$$S_p^* = \frac{\text{Uniprocessor Time of Solving } W^*}{\text{Parallel Time of Solving } W^*}$$
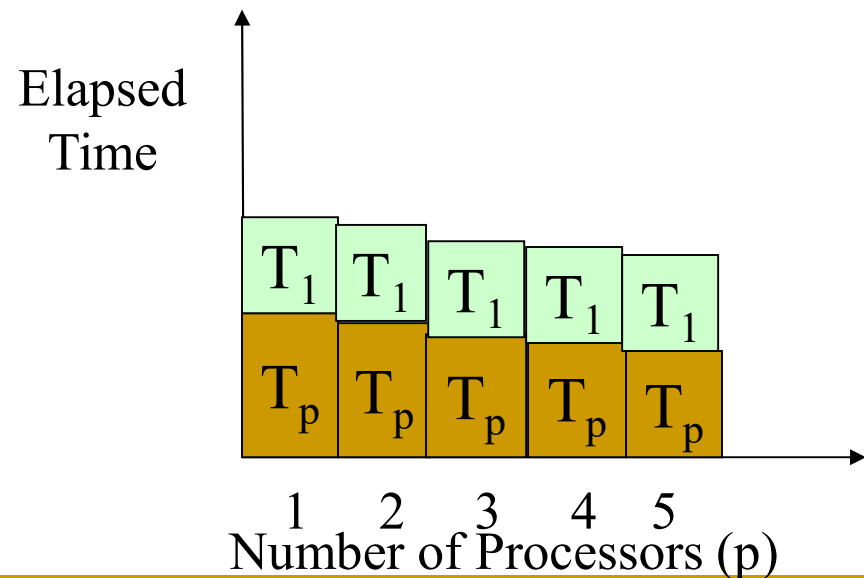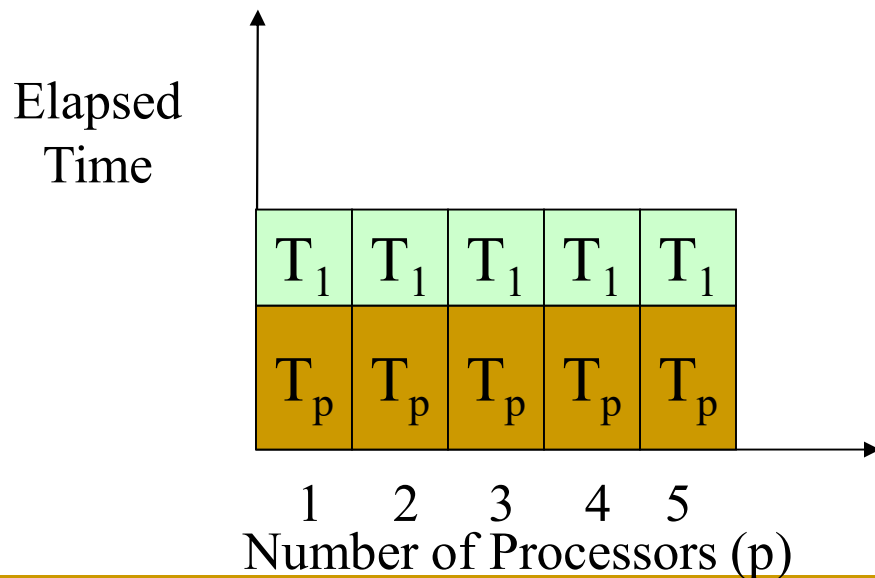
# Memory-Bounded Speedup (Sun & Ni)

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

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

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

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)

# 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

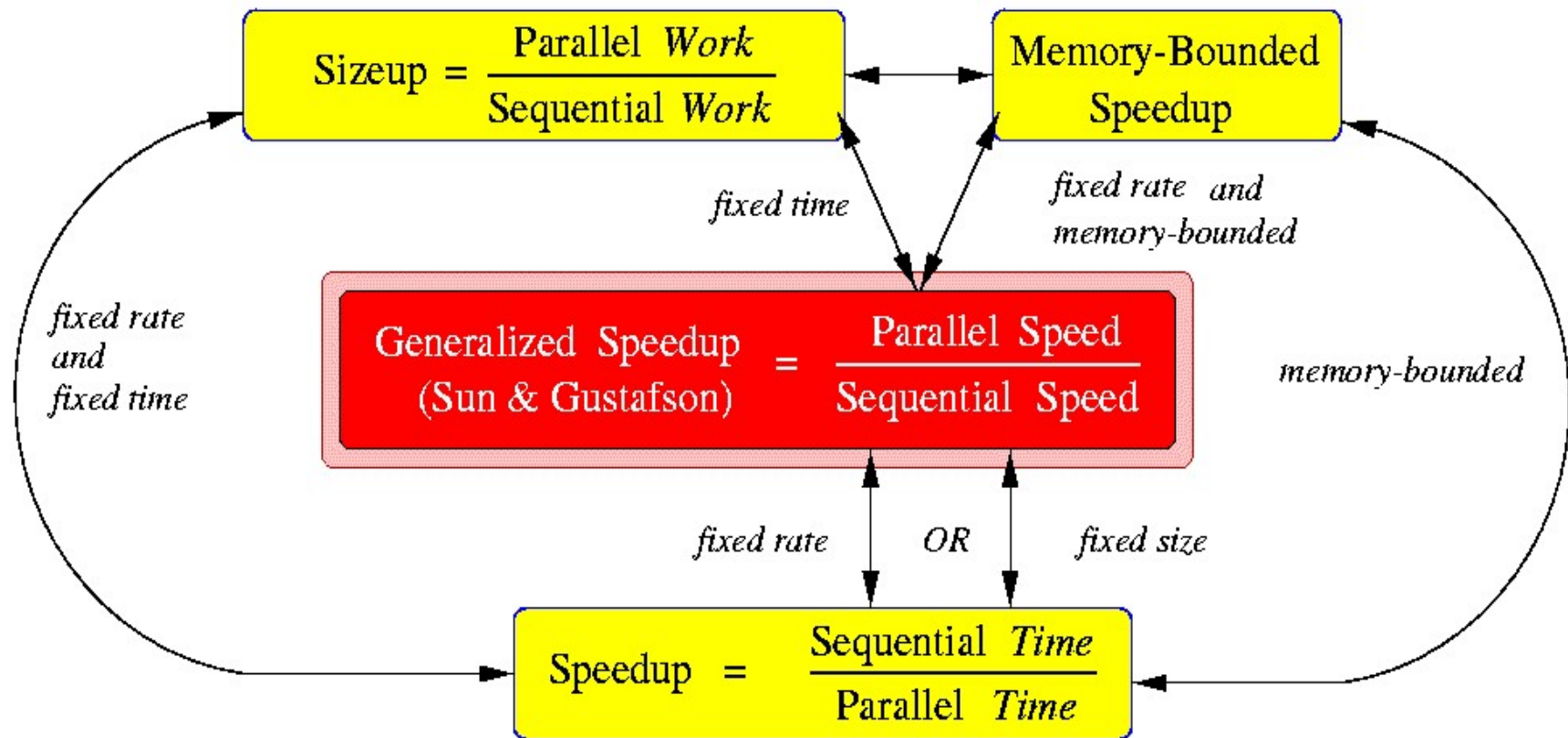$$S_p = \frac{Parallel\ Speed}{Sequential\ Speed}$$

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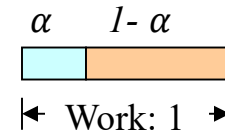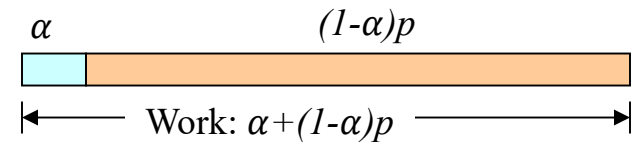
# Models of Speedup

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

Memory-Bounded Speedup

*fixed time*

*fixed rate and memory-bounded*

*fixed rate and fixed time*

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

*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
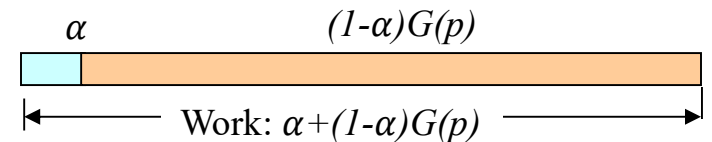  - Speedup emphasizes on time reduction

$\alpha$  $1-\alpha$

Work: 1

- Gustafson's Law, 1988
  - Fixed-time speedup model

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

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

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

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

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

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.
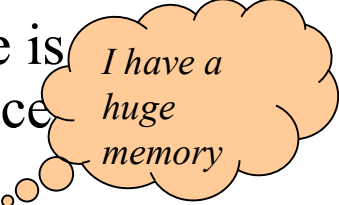
# *The Three Laws:* and their impact

*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, except 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)

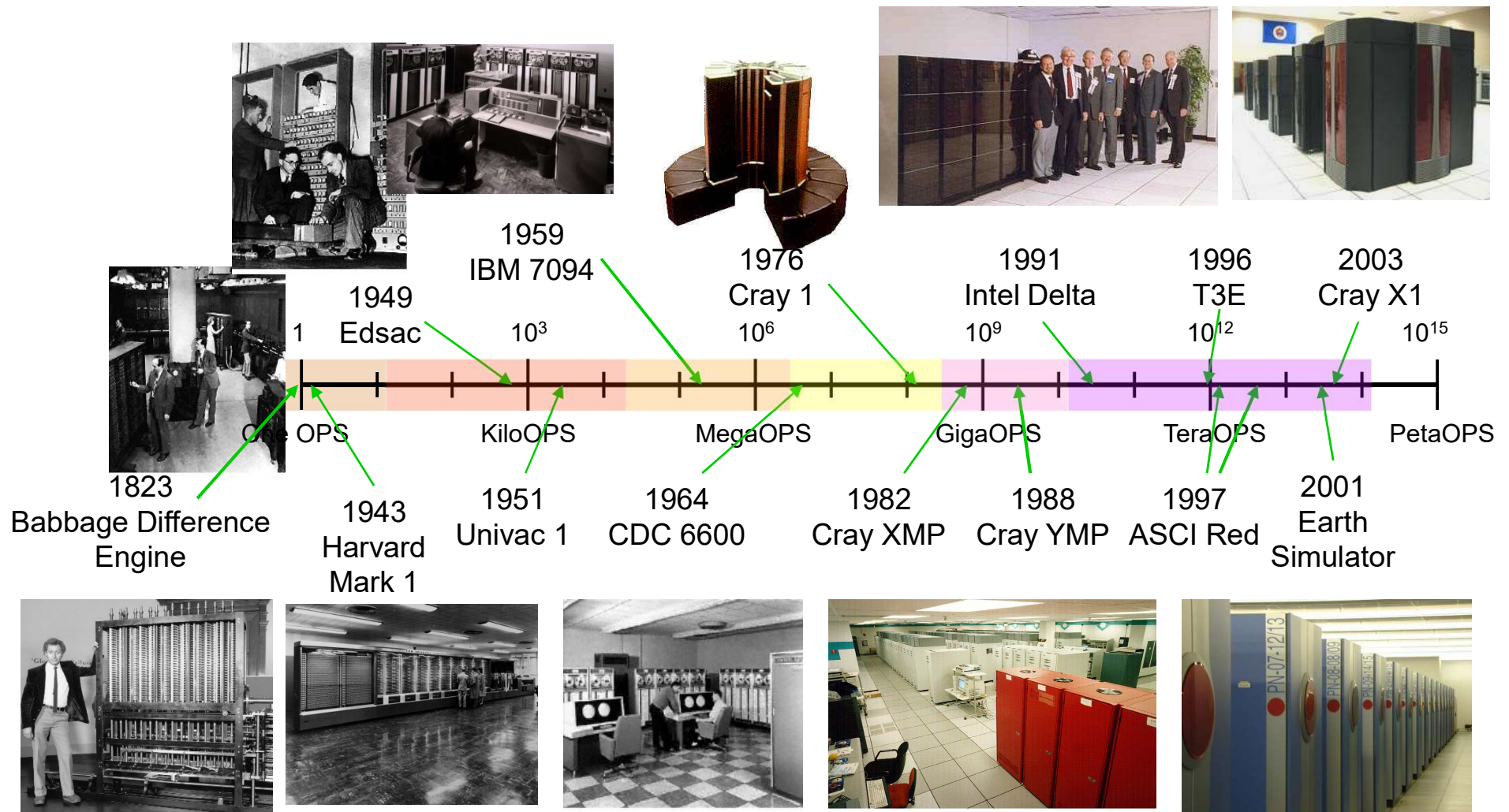- The **Memory-Wall Problem** (1994) shows memory-bound is a general performance issue for computing, not just for parallel computing

*William Wulf, Sally Mckee, "Hitting the memory wall: implications of the obvious," ACM SIGARCH Computer Architecture News Homepage archive, Vol. 23 Issue 1, March 1995*
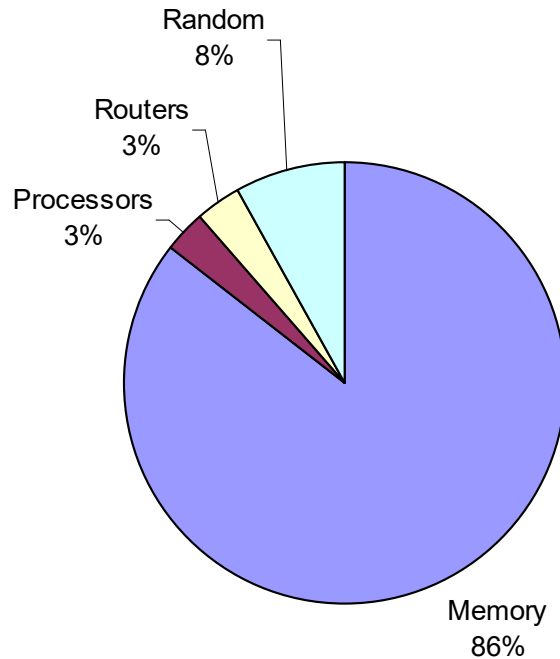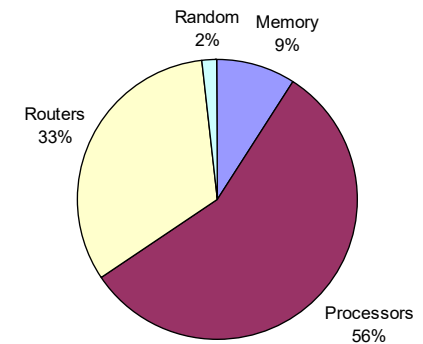
# Impact of Scalable Computing



1949
Edsac

1959
IBM 7094

1976
Cray 1

1991
Intel Delta

1996
T3E

2003
Cray X1

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

One OPS    KiloOPS    MegaOPS    GigaOPS    TeraOPS    PetaOPS

1823
Babbage Difference
Engine

1943
Harvard
Mark 1

1951
Univac 1

1964
CDC 6600

1982
Cray XMP

1988
Cray YMP

1997
ASCI Red

2001
Earth
Simulator

# Impact: Computing/Memory Trade-off

## Silicon Area Distribution

Random
8%

Routers
3%

Processors
3%

Memory
86%

## Power Distribution

Random
2%

Memory
9%

Routers
33%

Processors
56%

Modern microprocessors such as the Pentium Pro, Alpha 21164, Strong Arm SA110, and Longson-3A use 80% or more of their transistors for the on-chip cache

# Impact of Memory-Bounded Speedup
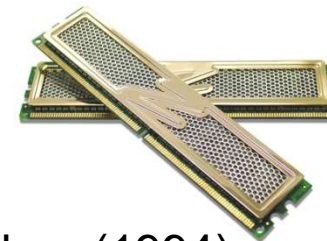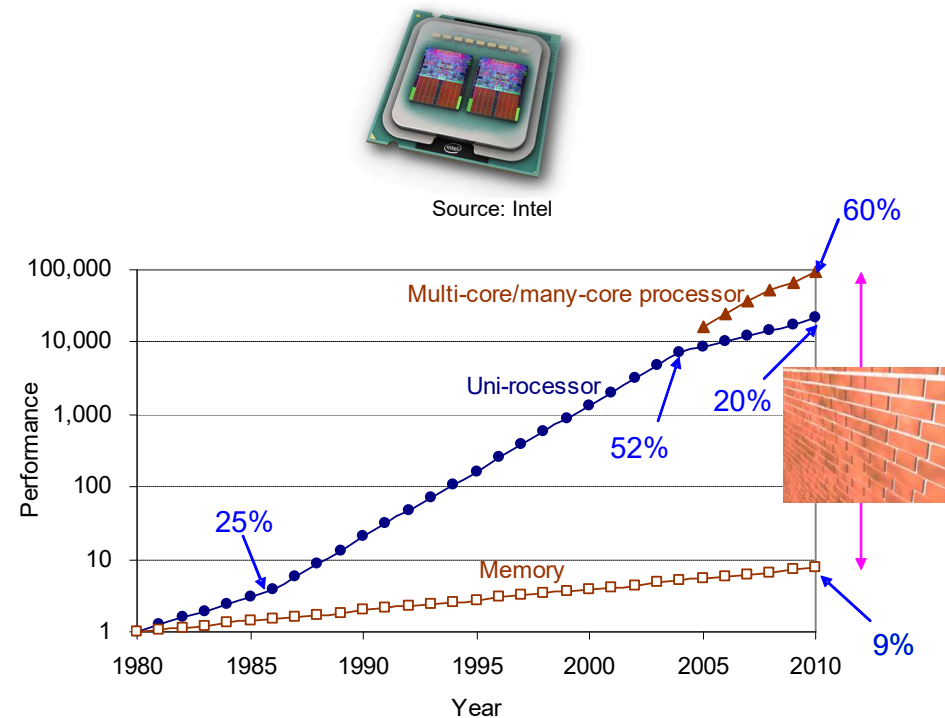
- **W = G(M)** shows the trade-off between computing & memory
  - W, the work in floating point operation
  - M, the memory requirement
  - G, the data reuse rate

- **W = G(M)** unifies the models
  - G(p) = 1, Amdahl's law
  - G(p) = p, Gustafson's law

- Reveal memory is the performance bottleneck
  - Memory-bounded algorithms and analysis in

    Dynamic programming, distributed optimization, search, convolution, regression, etc.
  - The Memory-Wall problem (1994)

# Impact: The Memory-wall Problem

- **Processor performance increases rapidly**
  - ❑ Uni-processor: ~52% until 2004
  - ❑ Aggregate multi-core/many-core processor performance even higher since 2004
- **Memory: ~9% per year**
  - ❑ Storage: ~6% per year
- **Processor-memory speed gap keeps increasing**



Source: Intel



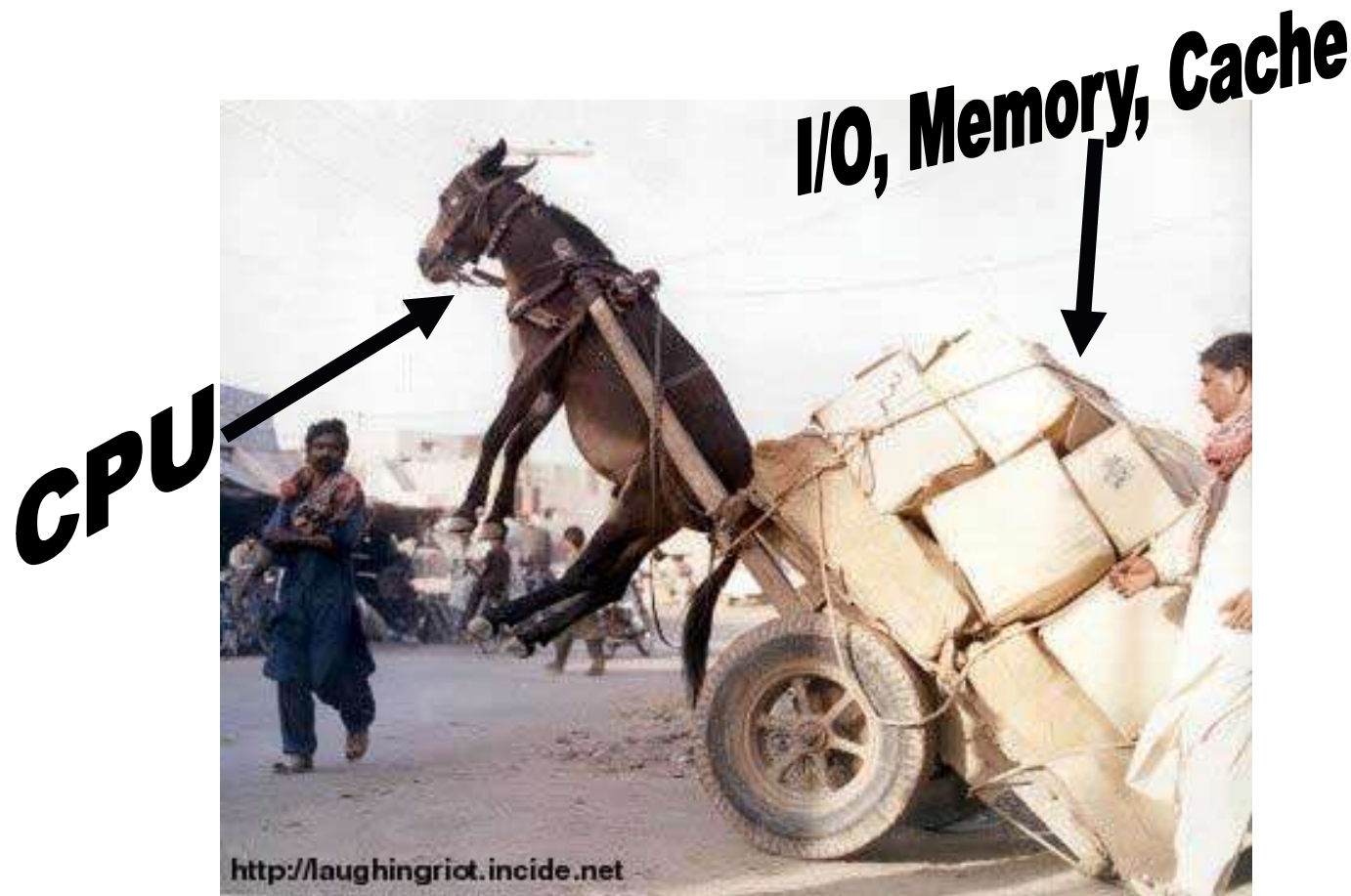Source: OCZ

Memory-bounded speedup (1990), Memory wall problem (1994)

# *The Beauty of Mathematics*



- The ability of abstract
- In depth understanding
  of the engineering issues
- Creative thinking

- Complex Specificity, Simple Genericity
- Abstract the complex specificity into simple genericity
- Engineering, mathematics, philosophy
- Everybody understand something, at a different level
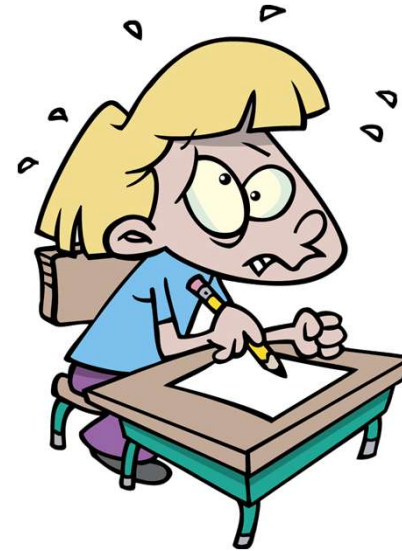- Your understanding determine your ability to apply it

- 厚积薄发，可遇不可求

# Big Data Makes Memory-Bound Even Worse



CPU

I/O, Memory, Cache

http://laughingriot.incide.net

- Source: Bob Colwell keynote ISCA'29 2002 http://systems.cs.colorado.edu/ISCA2002/Colwell-ISCA-KEYNOTE-2002-final.ppt

# How do we solve the memory-bound constraint or the memory-wall problem

# Solution: Memory Hierarchy

SRAM

DRAM

| Reg File | L1 Data cache | L2 Cache | Main Memory | DISK |

L1 Inst cache

# More on Memory Hierarchy & Concurrency

Multi-core
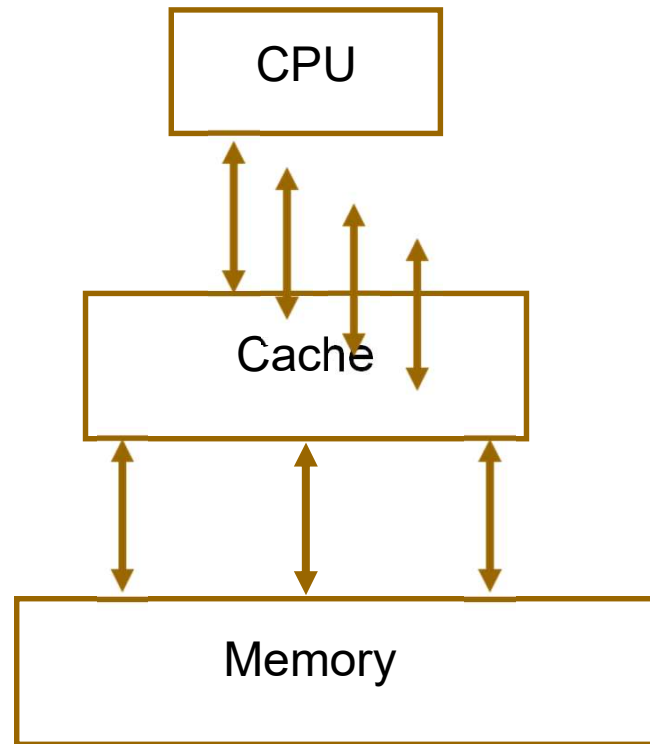Multi-threading
Multi-issue

Out-of-order Execution
Speculative Execution
Runahead Execution

**CPU**

Multi-banked Cache
Multi-level Cache

Pipelined Cache
Non-blocking Cache
Data Prefetching
Write buffer

**Cache**

Multi-channel
Multi-rank
Multi-bank

Pipeline
Non-blocking
Prefetching
Write buffer

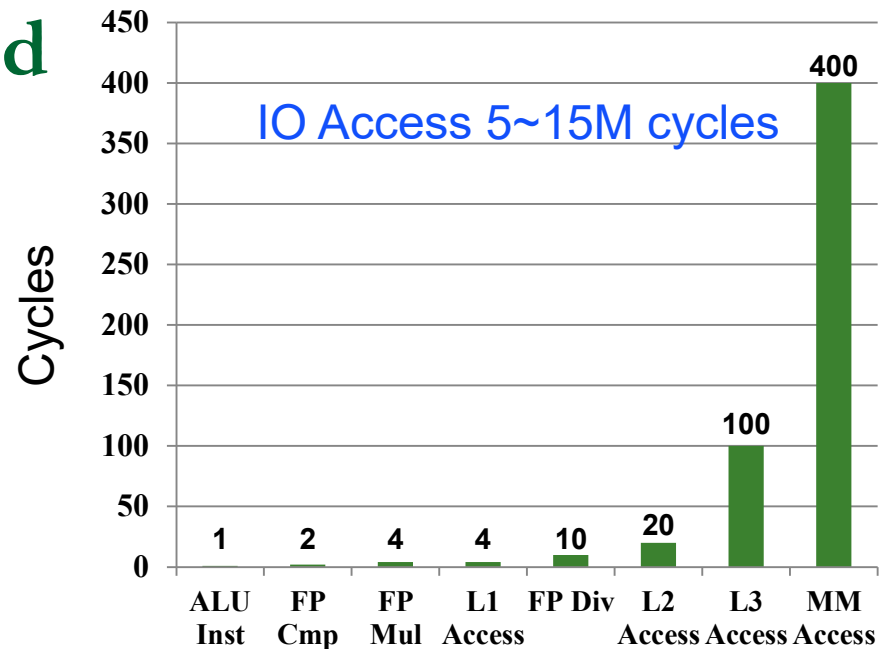**Memory**

## Input-Output (I/O)

*Parallel File System*

## Disks

# Assumption of Current Solutions

❑ Memory Hierarchy: Locality
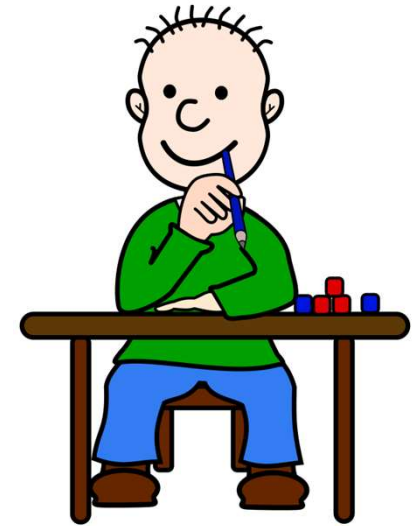❑ Concurrence: Data access pattern
  ○ Data stream

**Extremely Unbalanced Operation Latency**

**Performances vary largely**

IO Access 5~15M cycles

Cycles

| Label | Value |
|-------|-------|
| ALU Inst | 1 |
| FP Cmp | 2 |
| FP Mul | 4 |
| L1 Access | 4 |
| FP Div | 10 |
| L2 Access | 20 |
| L3 Access | 100 |
| MM Access | 400 |

*How do we further solve the memory-bound constraint or the memory-wall problem*

SEE YOU NEXT TIME
且听下回分解

# Welcome to my Research Team

*How can we produce classical research results?*