

# COMS3008A: Parallel Computing

Hairong Bau

School of Computer Science & Applied Mathematics

# Contents

- 1 Objectives
- 2 Course admin
- 3 Parallel computers
  - Parallel Computing — What is it?
  - Why parallelism?
  - Supercomputers
- 4 Classification of Parallel Computers
  - Control Structure Based Classification
- 5 A Quantitative Look at Parallel Computation
- 6 Limitations of memory system performance

# Outline

- 1 Objectives
- 2 Course admin
- 3 Parallel computers
  - Parallel Computing — What is it?
  - Why parallelism?
  - Supercomputers
- 4 Classification of Parallel Computers
  - Control Structure Based Classification
- 5 A Quantitative Look at Parallel Computation
- 6 Limitations of memory system performance

# Objectives

- Understand the basics of parallel computers, parallel computing, the motivation of parallel computing, and the classification of parallel computers.
- Understand and apply the simple quantitative modelling for parallel program performance.

# Outline

- 1 Objectives
- 2 Course admin
- 3 Parallel computers
  - Parallel Computing — What is it?
  - Why parallelism?
  - Supercomputers
- 4 Classification of Parallel Computers
  - Control Structure Based Classification
- 5 A Quantitative Look at Parallel Computation
- 6 Limitations of memory system performance

- Relevant course information is given in course outline (uploaded on ulwazi course site).
- Course related communications will be announced primarily through the announcement via ulwazi course site. It is important for you to check such announcements regularly to keep informed timely.

# Outline

- 1 Objectives
- 2 Course admin
- 3 **Parallel computers**
  - **Parallel Computing — What is it?**
  - **Why parallelism?**
  - **Supercomputers**
- 4 Classification of Parallel Computers
  - Control Structure Based Classification
- 5 A Quantitative Look at Parallel Computation
- 6 Limitations of memory system performance

# Outline

- 1 Objectives
- 2 Course admin
- 3 **Parallel computers**
  - **Parallel Computing — What is it?**
  - Why parallelism?
  - Supercomputers
- 4 Classification of Parallel Computers
  - Control Structure Based Classification
- 5 A Quantitative Look at Parallel Computation
- 6 Limitations of memory system performance



# Parallel computing

- **Parallel Computer:** A parallel computer is a computer system that uses multiple processing elements simultaneously in a cooperative manner to solve a computational problem.
- **Parallel computing (or processing):** Parallel processing includes techniques and technologies that make it possible to compute in parallel.
  - Hardware, networks, operating systems, parallel libraries, languages, compilers, algorithms, tools etc.
- **Parallel computing is an evolution of serial computing.**
  - Parallelism is natural.
  - Computing problems differ in level or type of parallelism.

# Serial computing

- A problem is broken into a discrete series of instructions;
- Instructions are executed sequentially one after another;
- Executed on a single processor;
- Only one instruction may execute at any moment in time.

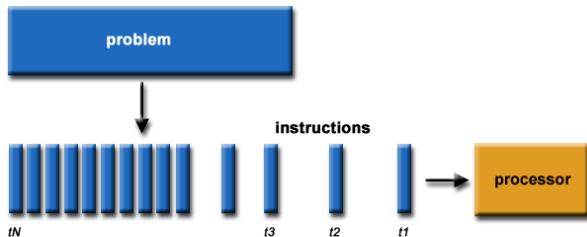


Figure: Serial computing

# Parallel computing

- A problem is broken into discrete parts that can be solved concurrently;
- Each part is further broken down to a series of instructions;
- Instructions from each part execute simultaneously on different processors;
- An overall control/coordination mechanism is employed.

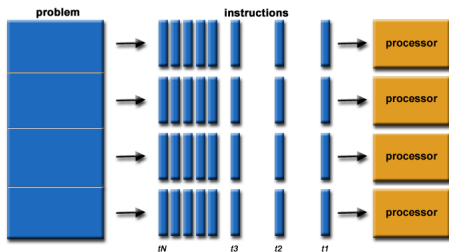


Figure: Parallel computing

# Parallel computing cont.

An example of parallelizing an addition of two vectors is shown in Figure 3.

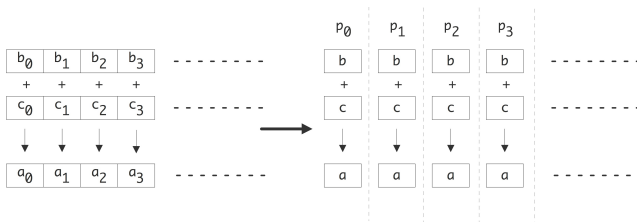


Figure: Parallelization of a vector addition

# Outline

1 Objectives

2 Course admin

3 **Parallel computers**

- Parallel Computing — What is it?
- **Why parallelism?**
- Supercomputers

4 Classification of Parallel Computers

- Control Structure Based Classification

5 A Quantitative Look at Parallel Computation

6 Limitations of memory system performance

# Why parallelism?

- In 2004, Intel changed its course from traditional chip design approach (single core processor) to embrace “dual core” processor structure.



Figure: Intel's shift in chip design to multi-core structure.

# A brief history of processor performance

- What had been happening before multi-processor? – Single-processor machines: they had been getting exponentially faster.

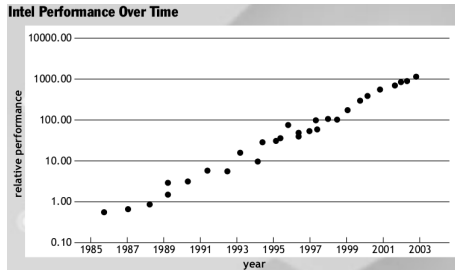


Figure: Intel single processor performance over time

# A brief history of processor performance cont.

- Wider data paths
  - 4 bit → 8 bit → 16 bit → 32 bit → 64 bit
- More efficient pipelining
  - For example, from 3.5 cycles per instruction (CPI) → 1.1 CPI
- Exploiting instruction level parallelism (ILP)
  - “Superscale” processing: e.g., issues up to 4 instructions/cycle
- Faster clock rates
  - e.g., 10MHz → 100MHz → 1GHz → 3GHz

During 80s and 90s, computers had large exponential performance gains.

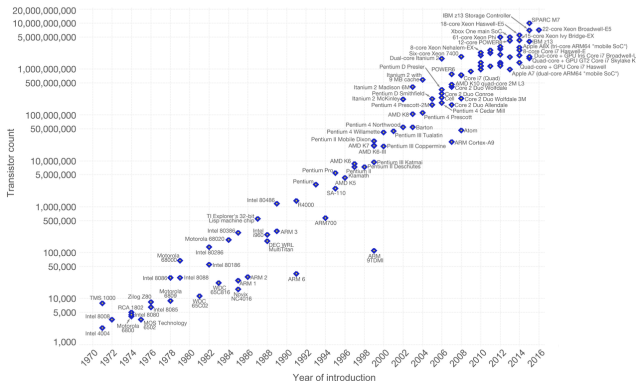


- “The number of transistors on an integrated circuit doubles every two years.” — Gordon E. Moore

## Moore's Law – The number of transistors on integrated circuit chips (1971-2016)

Our World  
in Data

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important as other aspects of technological progress – such as processing speed or the price of electronic products – are strongly linked to Moore's law.



Data source: Wikipedia ([https://en.wikipedia.org/wiki/Transistor\\_count](https://en.wikipedia.org/wiki/Transistor_count))

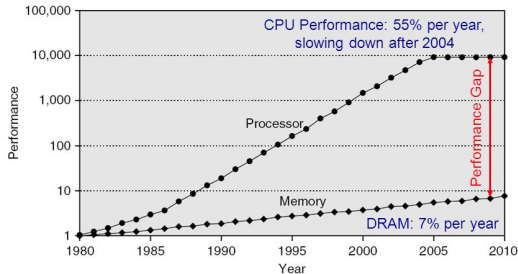
The data visualization is available at [OurWorldinData.org](https://ourworldindata.org). There you find more visualizations and research on this topic.

Licensed under [CC-BY-SA](#) by the author Max Roser.

Figure: Moore's Law – transistor count 1971 - 2016

# A brief history of processor performance cont.

## Processor-Memory Performance Gap



- ❖ 1980 – No cache in microprocessor
- ❖ 1995 – Two-level cache on microprocessor

Figure: Processor-memory performance gap

# A brief history of processor performance cont.

- For example, Intel Itanium II
  - 6-way integer unit < 2% die area;
  - Cache logic > 50% die area.
- Most of chip there to keep these 6 integer units at 'peak' rate.
- Main issue is external DRAM latency (50ns) to internal clock rate (0.25ns) ratio is 200:1.



**Figure:** Illustration of the die area for integer unit and cache logic for Intel Itanium II.

# A brief history of parallel computing

- Greater clock frequency, greater electrical power
- Intel VP Patrick Gelsinger (ISSCC 2001): “If scaling continues at present pace, by 2005, high speed processors would have power density of nuclear reactor, by 2010, a rocket nozzle, and by 2015, surface of sun.”

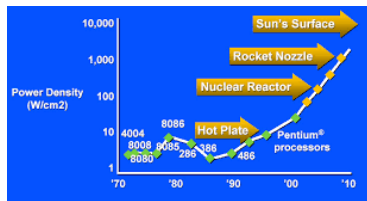
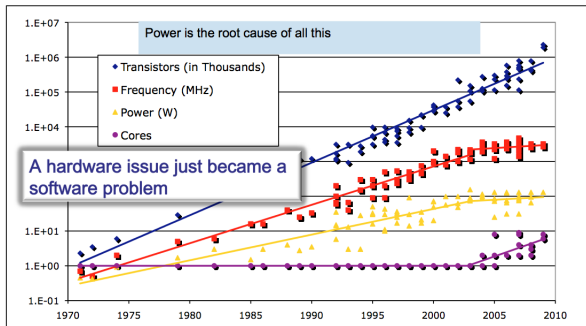


Figure: Intel VP Patrick Gelsinger (ISSCC 2001)

# A brief history of parallel computing cont.

- Add multiple cores to add performance, keep clock frequency the same or reduced.



Data from Kunle Olukotun, Lance Hammond, Herb Sutter,  
Burton Smith, Chris Batten, and Krste Asanović  
Slide from Kathy Yelick

Figure: Growth of transistors, frequency, cores, and power consumption

## A brief history of parallel computing cont.

- What we can do? We can still pack more and more transistors on to a die, but we cannot make the individual transistor faster and faster.
- We have to make better use of those more and more transistors to increase the performance instead of making only the clock speed faster. One solution is to use parallelism.

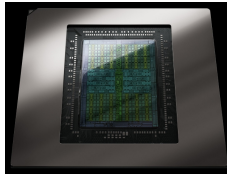
# Why parallelism – summary

- Serial machines have inherent limitations:
  - Processor speed, memory bottlenecks, . . .
- Parallelism has become the future of computing.
- Two primary benefits of parallel computing:
  - Solve fixed size problem in shorter time
  - Solve larger problems in a fixed time
- Other factors motivate parallel processing:
  - Effective use of machine resources;
  - Cost efficiencies;
  - Overcoming memory constraints.
- Performance is still the driving concern.
- Technology push
- Application pull
  - Application performance demands hardware advances;
  - Hardware advances generate new applications;
  - New applications have greater performance demands.

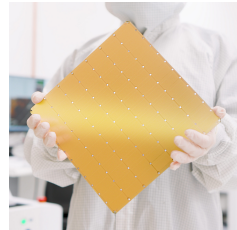
# Modern processors



(a)



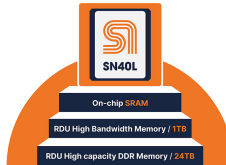
(b)



(c)



(d)



(e)



(f)

**Figure:** Processors: (a) Trillium TPU v6 (b) Nvidia GeForce 50 series (c) Cerebras' third-generation wafer-scale engine (WSE-3) (d), (e) SambaNova DataScale SN40L (f) Intel Core Ultra 9 285K



# Outline

1 Objectives

2 Course admin

3 **Parallel computers**

- Parallel Computing — What is it?
- Why parallelism?
- **Supercomputers**

4 Classification of Parallel Computers

- Control Structure Based Classification

5 A Quantitative Look at Parallel Computation

6 Limitations of memory system performance

# Top performing supercomputers

No 1 supercomputer in the top500 list from Nov 2024 – El Capitan:

<https://www.top500.org/system/180307/>;

---

Site	DOE/NNSA/LLNL
Manufacturer	HPE
Cores	11,039,616
Linpack Performance (Rmax)	1,742.00 PFlop/s
Theoretical Peak (Rpeak)	2,746.38 PFlop/s
Nmax	25,446,528
Power	29,580.98 kW
Processor	AMD 4th Gen EPYC 24C 1.8GHz
Interconnect	Slingshot-11
Operating System	TOSS
Compiler	g++ 12.2.1 and hipcc 6.2.0
Math library	AMD rocBLAS 6.0.2 and Intel MKL 2016
MPI	HPE Cray MPI

---

# No 1 Supercomputer in 11/2022 & 11/2023 — Frontier

---

Site	DOE/SC/Oak Ridge National Laboratory
Manufacturer	HPE ( <a href="#">Frontier</a> ; <a href="#">Frontier URL</a> )
Cores	8,730,112
Processor	AMD Optimized 3rd Generation EPYC 64C 2GHz
Interconnect	Slingshot-11
Linpack Performance (Rmax)	1,102.00 PFlop/s
Theoretical Peak (Rpeak)	1,685.65 PFlop/s
Nmax	24,440,832
Power	21,100.00 kW (Submitted)
Operating System	HPE Cray OS

---

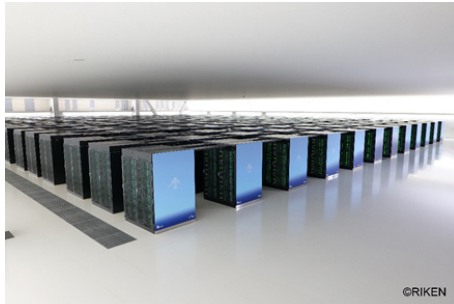
# No 1 Supercomputer in 11/2020 & 11/2021 — Fugaku

---

Site	RIKEN Center for Computational Science
Manufacturer	Fujitsu ( <a href="#">Fugaku</a> ; <a href="#">Fugaku virtual tour</a> )
Cores	7,630,848
Linpack Performance (Rmax)	442,010 TFlop/s
Theoretical Peak (Rpeak)	537,212 TFlop/s
HPCG	16,004.5 TFlop/s
Nmax	21,288,960
Power	29,899.23 kW (Submitted)
Memory	5,087,232 GB
Processor	A64FX 48C 2.2GHz
Interconnect	Tofu interconnect D
Operating System	RHEL
Compiler	Fujitsu software technical computing suite v4.0
Math library	Fujitsu software technical computing suite v4.0
MPI	Fujitsu software technical computing suite v4.0

---

# No 1 Supercomputer in 11/2020 & 11/2021 — Fugaku cont.



**Figure:** Fugaku, a supercomputer from RIKEN and Fujitsu Limited.

# The Clusters at Mathematical Sciences Lab

As of 2024, the clusters (or partitions) at Wits Mathematical Sciences Lab (MSL)

- stampede: For general purpose use or jobs that can leverage InfiniBand (if enabled). It has 32 nodes, each with two Xeon E5-2680 CPUs, two GTX1060 GPUs (6GB per GPU, 12GB per node), and 32GB of system RAM. MaxTime=4320 and MaxNodes=16.
- bigbatch: For bigger jobs that can leverage a bigger GPU and additional system RAM. It has 48 nodes each with a single Intel Core i9-10940X CPU (14 cores), NVIDIA RTX3090 GPU (24GB), and 128GB of system RAM. MaxTime=4320 and MaxNodes=16. Partition bigbatch has the potential for 10Gb networking (Not sure if this has been implemented).

# The Clusters at Mathematical Sciences Lab cont.

- **biggpu:** For mature code that can meaningfully leverage very large amounts of GPU and system RAM. It has 4 nodes, each node has two Intel Xeon Platinum 8280L CPUs (28 cores per CPU, 56 cores per node) with two NVIDIA Quadro RTX8000 GPUs (48GB per GPU, 96GB per node), and 1TB of system RAM. MaxTime=4320 and MaxNodes=3. biggpu has the potential for 10Gb networking (Not sure if this has been implemented). Whenever possible, limit jobs on biggpu to using only a single node so others can use the partition simultaneously. Please use biggpu responsibly!

# Outline

- 1 Objectives
- 2 Course admin
- 3 Parallel computers
  - Parallel Computing — What is it?
  - Why parallelism?
  - Supercomputers
- 4 **Classification of Parallel Computers**
  - **Control Structure Based Classification**
- 5 A Quantitative Look at Parallel Computation
- 6 Limitations of memory system performance



# Outline

- 1 Objectives
- 2 Course admin
- 3 Parallel computers
  - Parallel Computing — What is it?
  - Why parallelism?
  - Supercomputers
- 4 **Classification of Parallel Computers**
  - **Control Structure Based Classification**
- 5 A Quantitative Look at Parallel Computation
- 6 Limitations of memory system performance

# Flynn's taxonomy

Flynn's taxonomy is widely used since 1966 for classification of parallel computers. The classification is based on two independent dimensions of *instruction stream* and *data stream* with two possible states: *single* or *multiple*.

- SISD: Single instruction stream single data stream. This is the traditional CPU architecture: at any one time only a single instruction is executed, operating on a single data item.

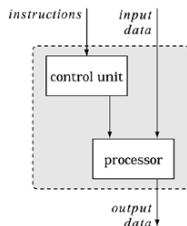


Figure: The SISD architecture

# SIMD

- SIMD: Single instruction stream multiple data stream. In this computer type there can be multiple processors, each operating on its own data item, but they are all executing the same instruction on that data item. SIMD computers excel at operations on arrays, such as

*for*( $i = 0; i < N; i++$ )     $a[i] = b[i] + c[i];$

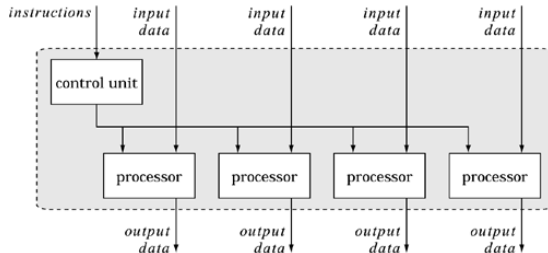


Figure: The SIMD architecture

# MIMD

- MISD: Multiple instruction stream single data stream. Each processing unit executes different instruction streams on a single data stream. Very few computers are in this type.
- MIMD: Multiple instruction stream multiple data stream. Multiple processors operate on multiple data items, each executing independent, possibly different instructions. Most current parallel computers are of this type.

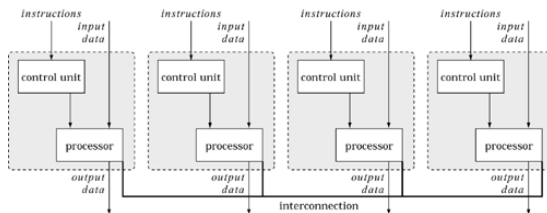


Figure: The MIMD architecture

- Most of MIMD machines operate in *single program multiple data* (SPMD) mode, where the programmers starts up the same executable on the parallel processors.

# A Further Decomposition of MIMD

The MIMD category is typically further decomposed according to memory organization: shared memory and distributed memory.

- **Shared memory:** In a shared memory system, all processes share a single address space and communicate with each other by writing and reading shared variables.
- One class of shared-memory systems is called SMPs (symmetric multiprocessors).

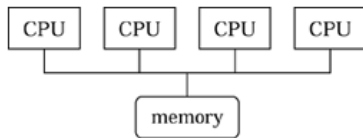


Figure: The SMP architecture

# NUMA

- The other main class of shared-memory systems is called non-uniform memory access (NUMA). The memory is shared, it is uniformly addressable from all processors, but some blocks of memory may be physically more closely associated with some processors than others.
- To mitigate the effects of non-uniform access, each processor has a cache, along with a protocol to keep cache entries coherent — cache-coherent non-uniform memory access systems (ccNUMA).

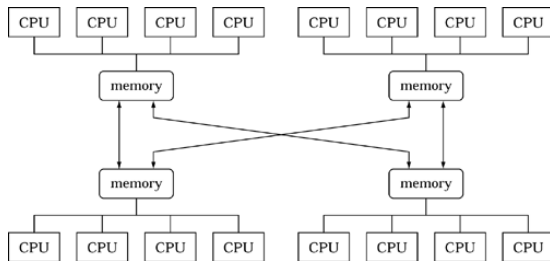


Figure: The NUMA architecture

# Distributed memory systems

- Each process has its own address space and communicates with other processes by message passing (sending and receiving messages).

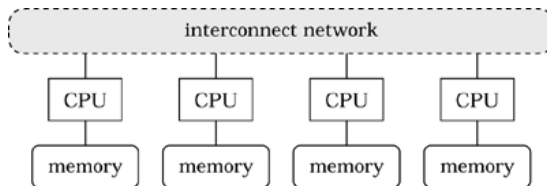


Figure: The distributed memory architecture



# Clusters

- Clusters are distributed memory systems composed of off-the-shelf computers connected by an off-the-shelf network.

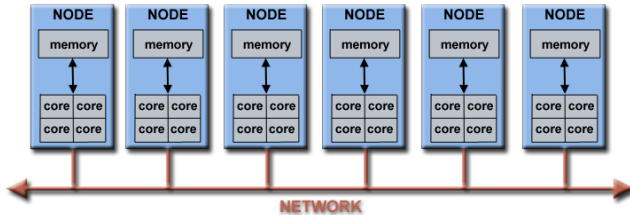


Figure: A cluster

# Outline

- 1 Objectives
- 2 Course admin
- 3 Parallel computers
  - Parallel Computing — What is it?
  - Why parallelism?
  - Supercomputers
- 4 Classification of Parallel Computers
  - Control Structure Based Classification
- 5 A Quantitative Look at Parallel Computation**
- 6 Limitations of memory system performance

# A simple performance modelling

- Consider a computation consisting of three parts: a setup section, a computation section, and a finalization section.
- The total running time of this program on one processing element (PE) is given as:

$$T_{total}(1) = T_{setup} + T_{compute} + T_{finalization} \quad (1)$$

- What happens when we run this computation on a parallel computer with multiple PEs?

$$T_{total}(P) = T_{setup} + \frac{T_{compute}(1)}{P} + T_{finalization} \quad (2)$$

- An important measure of how much additional PEs help is the relative speedup  $S$ , which describes how much faster a problem runs:

$$S(P) = \frac{T_{total}(1)}{T_{total}(P)} \quad (3)$$

- A related measure is the efficiency  $E$ , which is the speedup normalized by the number of PEs.

$$E(P) = \frac{S(P)}{P} = \frac{T_{total}(1)}{PT_{total}(P)} \quad (4)$$

- Ideally, we would want the speedup to be equal to  $P$ , the number of PEs. This is sometimes called perfect linear speedup.

# Amdahl's Law

- The terms that cannot be run concurrently are called the serial terms. Their running times represent some fraction of the total, called the serial fraction, denoted  $\gamma$ .

$$\gamma = \frac{T_{\text{setup}} + T_{\text{finalization}}}{T_{\text{total}}(1)} \quad (5)$$

- The fraction of time for parallelizable part is then  $1 - \gamma$ . The total computation time with  $P$  PEs becomes

$$T_{\text{total}}(P) = \gamma T_{\text{total}}(1) + \frac{(1 - \gamma) T_{\text{total}}(1)}{P} \quad (6)$$

# Amdahl's Law cont.

- Then we obtain the well-known Amdahl's law:

$$S(P) = \frac{T_{total}(1)}{T_{total}(P)} = \frac{T_{total}(1)}{(\gamma + \frac{1-\gamma}{P}) T_{total}(1)} = \frac{1}{\gamma + \frac{1-\gamma}{P}} \quad (7)$$

Taking the limit as P goes to infinity in Eq. 7,

$$S(P) = \frac{1}{\gamma} \quad (8)$$

Eq. 8 gives an upper bound on the speedup.

- Amdahl's Law says the speedup you can achieve is limited by the fraction for serial computation in your problem (i.e., the number of PEs does not determine the upper bound of speedup you can achieve).

# Amdahl's Law cont.

## Example 1

Suppose we are able to parallelize 90% of a serial program. Further suppose the speedup of this part is  $P$ , the number of processes we used (which is a perfect linear speedup). If the serial time is  $T_{serial} = 20s$ , then the runtime of the parallelized part is  $(0.9 \times T_{serial})/P = 18/P$ . The runtime of unparallelized part is  $0.1 \times 20 = 2s$ . The overall parallel runtime will be

$$T_{parallel} = 18/P + 2,$$

and the speedup will be

$$S = \frac{T_{serial}}{T_{parallel}} = \frac{20}{18/p + 2}.$$

As  $P$  gets larger,  $18/P$  gets close to 0, and the denominator gets close to 2, in turn,  $S$  gets close to 10. That means  $S \leq 10$ , no matter how many number of processes you use in your program.

# Gustafson's Law

- In contrast to Eq. 2, we obtain  $T_{total}(1)$  from the serial and parallel parts when executed on  $P$  PEs.

$$T_{total}(1) = T_{setup} + PT_{compute}(P) + T_{finalization} \quad (9)$$

- Now, we define the so-called scaled serial fraction, denoted  $\gamma_{scaled}$ , as

$$\gamma_{scaled} = \frac{T_{setup} + T_{finalization}}{T_{total}(P)}, \quad (10)$$

and then

$$T_{total}(1) = \gamma_{scaled} T_{total}(P) + P(1 - \gamma_{scaled}) T_{total}(P). \quad (11)$$



# Gustafson's Law cont.

- Using Eq. 11, we obtain the scaled speedup, sometimes known as Gustafson's law.

$$\begin{aligned} S(P) &= \frac{T_{total}(1)}{T_{total}(P)} = \frac{\gamma_{scaled} T_{total}(P) + P(1 - \gamma_{scaled}) T_{total}(P)}{T_{total}(P)} \\ &= P(1 - \gamma_{scaled}) + \gamma_{scaled} = P + (1 - P)\gamma_{scaled}. \end{aligned} \quad (12)$$

- Suppose we take the limit in  $P$  while holding  $T_{compute}$  and  $\gamma_{scaled}$  constant. That is, we are increasing the size of the problem so that the total running time remains constant when more processors are added. In this case, the speedup is linear in  $P$ .

## Example 2

Now, for the previous example, let's assume the scaled serial fraction of the same problem is 0.1, that is  $\gamma_{scaled} = 0.1$ , and  $p = 16$ . Then the scaled speedup is  $S = 16(1 - 0.1) + 0.1 = 14.5$ , which is greater than the upper bound determined by Amdahl's Law (10).

- There are another two performance metrics we frequently use in the process of the course: efficiency and scalability.
- Efficiency:

$$E = \frac{S}{P}, \quad (13)$$

where  $S$  is the speedup, and  $P$  is the number of processes used to achieve the speedup. Note  $0 < E \leq 1$ . Efficiency is better when  $E$  is closer to 1, which simply means you utilized the  $P$  processors efficiently.

- Scalability: In general, a technology is scalable if it can handle ever-increasing problem size.
- In parallel program performance, scalability refers to the following measure. Suppose we run a parallel program using certain number of processors and with a certain problem size, and obtained an efficiency  $E$ . Further suppose that now we want to increase the number of processors. In this case, if we can find a rate at which the problem size increases so that we can still maintain the efficiency  $E$ , we say the parallel program scalable.

## Example 3

Suppose  $T_{\text{serial}} = n$ , where  $n$  is also the problem size. Also suppose  $T_{\text{parallel}} = n/p + 1$ . Then

$$E = \frac{n}{p(n/p + 1)} = \frac{n}{n + p}.$$

To see if the problem is scalable, we increase  $p$  by a factor of  $k$ , then we get

$$E = \frac{n}{n + p} = \frac{xn}{xn + kp}.$$

If  $x = k$ , then we have the same efficiency. That is, if we increase the problem size at the same rate that we increase the number of processors, then the efficiency remain constant, hence, the program is scalable.

- **Strong scalability:** When we increase the number of processes, we can keep the efficiency fixed without increasing the problem size, the program is strongly scalable.
- **Weak scalability:** If we can keep the efficiency fixed by increasing the problem size at the same rate as we increase the number of processes, then the program is said to be weakly scalable.

# Outline

- 1 Objectives
- 2 Course admin
- 3 Parallel computers
  - Parallel Computing — What is it?
  - Why parallelism?
  - Supercomputers
- 4 Classification of Parallel Computers
  - Control Structure Based Classification
- 5 A Quantitative Look at Parallel Computation
- 6 Limitations of memory system performance

# Memory performance - latency

- Latency: A memory system, possibly with multiple levels of cache, takes in a request for a memory word and returns a block of data of size  $b$  containing the requested word after  $l$  ns. Here,  $l$  is referred to as the latency of the memory.

## Example 4 (Effect of latency on memory performance)

Consider a processor operating at 1GHz (1ns clock) connected to a DRAM with a latency of 100ns (no caches). Assume that the processor has two multiply-add units and is capable of executing four instructions in each cycle of 1ns. The peak performance rating is therefore 4GFLOPS. However, since the memory latency is 100ns, and if the block size is only one word (4 bytes), then every time a memory request is made, the processor must wait 100ns (latency) before it can process a word. In the case of each floating point operation requires one data fetch, the peak speed of this computation is limited to one floating point operation every 100ns, or a speed of 10MFLOPS, which is only a small fraction of the peak processor performance.



# Effect of cache on memory performance

- Improving effective memory latency using caches: Caches are used to address the speed mismatch between the processor and memory.
- **Cache hit ratio:** The fraction of data references satisfied by the cache is called the cache hit ratio of the computation on the system.
- **Memory bound computation:** The effective computation rate of many applications is bounded not by the processing rate of the CPU, but by the rate at which data can be pumped into the CPU. Such computations are referred to as being memory bound. The performance of memory bound program is critically impacted by the cache hit ratio.
- **Temporal locality:** The notion of repeated reference to a data item in a small time window is called temporal locality of reference.

## Effect of cache on memory performance cont.

### Example 5

We still consider the 1GHz processor with a 100ns latency DRAM as in Example 4. In this case, let's consider adding a cache of size 32KB ( $1K = 2^{10} \approx 10^3$ ) with a latency of 1ns. We use this setup to multiply two matrices  $A$  and  $B$  of dimensions  $32 \times 32$  (for this, the input matrices and output matrix can all fit in the cache). Fetching the two input matrices (2K words) into the cache takes approximately  $200\mu s$ . The total floating point operations in multiplying the two matrix is 64K (multiplying two  $n \times n$  matrices takes  $2n^3$  floating point operations), which takes approximately 16K cycles, i.e., 16Kns (or  $16\mu s$ ) at 4 instructions per cycle. So, the total time for the computation is  $200\mu s + 16\mu s = 216\mu s$ . This corresponds to a peak computation rate of  $64K/216\mu s \approx 303\text{MFLOPS}$ . Compared to Example 4, this is a much better rate (due to cache), but is still a small fraction of the peak processor performance (4GFLOPS).

# Effect of bandwidth on memory performance

- **Bandwidth:** The rate at which data can be moved between the processor and the memory. It is determined by the bandwidth of the memory bus as well as the memory units.

## Example 6

We continue with Examples 4 and 5. Now consider increasing the block size to 4 words (memory bandwidth in this case; also, if a memory request returns a contiguous block of 4 words, the 4-word unit is also considered as a cache line). Now fetching the two input matrices into the cache takes  $200/4 = 50\mu s$ . Consequently the total time for the computation is  $50\mu s + 16\mu s = 66\mu s$ , which corresponds to a peak computation rate of  $64K/66\mu s \approx 993\text{MFLOPS}$ .

## Effect of bandwidth on memory performance cont.

- Increasing the memory bandwidth, or building wide data bus connected to multiple memory banks, is expensive to construct. In a more practical system, consecutive words are sent on the memory bus on subsequent bus cycles after the first word is retrieved. For example, with a 32-bit data bus, the first word is put on the bus after 100ns (the latency) and one word is put on each subsequent bus cycle. Then the 4 words will become available after  $100 + 3 \times (\text{memory bus cycle})\text{ns}$ . Assuming a data bus operating at 200MHz, this adds 15ns to the cache line access time.
- **Spatial locality:** Successive computations require contiguous data. If the computation (or data access pattern) in a program does not have spatial locality, then effective bandwidth can be much smaller than the peak bandwidth.

We discussed

- Definitions of parallel computers and parallel computing
- Why parallel computing
- Supercomputers
- Classification of parallel computers
- Simple modelling of parallel computing performance
- Limitation of memory system performance

# Exercise

Complete all questions in Sec1.6, PC\_notes\_2024.