

CE425 HIGH PERFORMANCE COMPUTING

- ASSESSMENT
 - CA – 40%
 - TESTS – 30%
 - PRESENTATIONS – 10%
 - EXAM – 60%
- SYLLUBUS & RECOMMENDED BOOKS
 - CLASS REPRESENTATIVE
 - Lecture Notes:
www.Lechaamwe.weebly.com

1/13/2022

1

HPC

- Overview
- Clarification of parallel machines
- Some General Parallel Terminology
- Shared memory and message passing

1/13/2022

2

Overview

- “supercomputing” vs “high performance computing” vs “cloud computing”
 - Supercomputers are computers capable of supercomputing
 - Supercomputing uses computers to research, design products and support complex decisions.
 - it also includes software systems and testing tools, and algorithms to solve complex computing

1/13/2022

3

Overview

- High-performance computers are parallel computers composed of traditional ultra high-speed computers and multiple CPU.
- HPC is almost equivalent to supercomputing, mainly for scientific computing, engineering simulation, animation rendering and other fields, belonging to the computing-intensive applications.

1/13/2022

4

Overview

- The design and manufacture of private cars, the orbit simulation of space shuttles, the design of Nike shoes, the development of drugs and so on all belong to the category of high performance computing. At first, high performance computing was used in cryptography,

1/13/2022

5

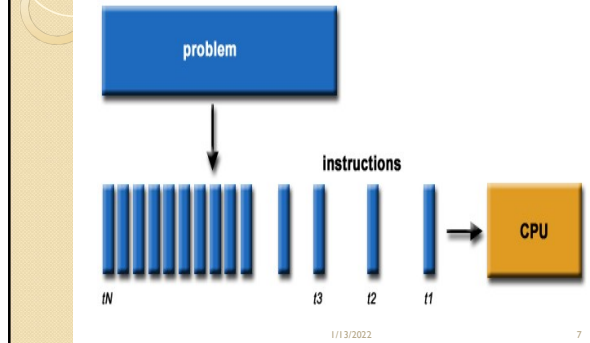
What is Parallel Computing?

- Traditionally, software has been written for **serial** computation:
 - To be run on a single computer having a single Central Processing Unit (CPU);
 - A problem is broken into a discrete series of instructions.
 - Instructions are executed one after another.
 - Only one instruction may execute at any moment in time.

1/13/2022

6

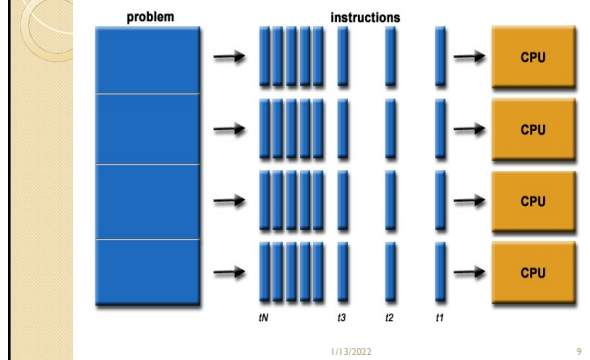
What is Parallel Computing?



What is Parallel Computing?

- **Parallel computing** is the simultaneous use of multiple computer resources to solve a computational problem:
 - To be run using multiple CPUs
 - 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 CPUs

What is Parallel Computing?



What is Parallel Computing?

- The compute resources might be:
 - A single computer with multiple processors;
 - An arbitrary number of computers connected by a network;
 - A combination of both.

What is Parallel Computing?

- The computational problem should be able to:
 - Be broken apart into discrete pieces of work that can be solved simultaneously;
 - Execute multiple program instructions at any moment in time;
 - Be solved in less time with multiple computer resources than with a single computer resource.

Why Use Parallel Computing?

- **Save time and/or money:**
 - In theory, throwing more resources at a task will shorten its time to completion, with potential cost savings.
 - Parallel computers can be built from cheap, commodity components.

Why Use Parallel Computing?

- **Solve larger problems:**
 - Many problems are so large and/or complex that it is impractical or impossible to solve them on a single computer,
 - especially given limited computer memory.
 - Web search engines/databases processing millions of transactions per second

1/13/2022

13

Why Use Parallel Computing?

- **Provide concurrency:**
 - A single computer resource can only do one thing at a time.
 - Multiple computing resources can be doing many things simultaneously.
 - For example, the Access Grid (accessgrid.org) provides a global collaboration network where people from around the world can meet and conduct work "virtually".

1/13/2022

14

Why Use Parallel Computing?

- **Use of non-local resources:**
 - Using computer resources on a wide area network, or even the Internet when local computer resources are scarce.
 - For example: SETI@home (setiathome.berkeley.edu) uses 2.9 million computers in 253 countries.
 - Folding@home (folding.stanford.edu) uses over 450,000 cpus globally

1/13/2022

15

Why Use Parallel Computing?

- **Limits to serial computing:**
 - Both physical and practical reasons pose significant constraints to simply building ever faster serial computers:
 - Transmission speeds - the speed of a serial computer is directly dependent upon how fast data can move through hardware.
 - Absolute limits are the speed of light (30 cm/nanosecond) and the transmission limit of copper wire (9 cm/nanosecond).
 - Increasing speeds necessitate increasing proximity of processing elements.

1/13/2022

16

Why Use Parallel Computing?

- Limits to miniaturization - processor technology is allowing an increasing number of transistors to be placed on a chip.
 - However, even with molecular or atomic-level components, a limit will be reached on how small components can be.

1/13/2022

17

Why Use Parallel Computing?

- Economic limitations - it is increasingly expensive to make a single processor faster.
 - Using a larger number of moderately fast commodity processors to achieve the same (or better) performance is less expensive.

1/13/2022

18

Disadvantages of Parallel Computing

- It addresses Parallel architecture that can be difficult to achieve.
- In the case of clusters, better cooling technologies are needed in parallel computing.
- The multi-core architectures consume high power consumption.
- The parallel computing system needs low coupling and high cohesion, which is difficult to create.

1/13/2022 19

Disadvantages of Parallel Computing

- The code for a parallelism-based program can be done by the most technically skilled and expert programmers.
- Due to synchronization, thread creation, data transfers, and more, the extra cost sometimes can be quite large; even it may be exceeding the gains because of parallelization.
- Moreover, for improving performance, the parallel computing system needs different code tweaking for different target architectures

1/13/2022 20

Classification of Parallel Machines

- **Models of Computation (Flynn 1966)**
- Any computer, whether sequential or parallel, operates by executing instructions on data.
- a stream of **instructions** (the algorithm) tells the computer what to do.
- a stream of **data** (the input) is affected by these instructions.

1/13/2022 21

Classification of Parallel Machines

- Depending on whether there is one or several of these streams, we have four classes of computers.
- Single Instruction Stream, Single Data Stream : SISD.
- Multiple Instruction Stream, Single Data Stream : MISD.
- Single Instruction Stream, Multiple Data Stream : SIMD.
- Multiple Instruction Stream, Multiple Data Stream : MIMD.

1/13/2022 22

SISD Computers

- This is the standard sequential computer.
- A single processing unit receives a single stream of instructions that operate on a single stream of data.

```

    graph LR
      CONTROL[CONTROL] -- INSTRUCTION STREAM --> PROCESSOR[PROCESSOR]
      PROCESSOR -- DATA STREAM --> MEMORY[MEMORY]
  
```

1/13/2022 23

MISD Computers

- N processors, each with its own control unit, share a common memory.

```

    graph LR
      MEMORY[MEMORY] -- DS --> P1[PROCESSOR ONE]
      MEMORY -- DS --> P2[PROCESSOR TWO]
      MEMORY -- DS --> Pn[PROCESSOR N]
      C1[CONTROL ONE] -- IS ONE --> P1
      C2[CONTROL TWO] -- IS TWO --> P2
      Cn[CONTROL N] -- IS N --> Pn
  
```

1/13/2022 24

MISD Computers

- There are N streams of instructions (algorithms / programs) and **one** stream of data.
- Parallelism is achieved by letting the processors do different things at the same time on the same datum.
- MISD machines are useful in computations where the same input is to be subjected to several different operations.

1/13/2022

25

Example

- Checking whether a number Z is prime.
- A simple solution is to try all possible divisions of Z .
- Assume the number of processors, N , is given by $N = Z - 2$.
- All processors take Z as input and tries to divide it by its associated divisor.
- So it is possible in **one step** to check if Z is prime.

1/13/2022

26

MISD Computers

- More realistically if $N < Z - 2$ then a subset of divisors would be assigned to each processor.
- For most applications MISD are very awkward to use and no commercial machines exist with this design.

1/13/2022

27

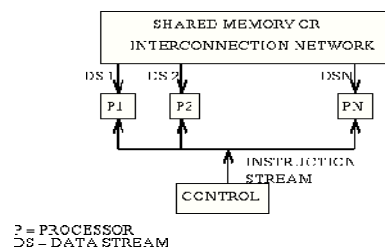
SIMD Computers

- All N identical processors operate under the control of a single instruction stream issued by a central control unit.
- (to ease understanding assume that each processor holds the same identical program.)
- There are N data streams, one per processor so different data can be used in each processor.

1/13/2022

28

SIMD Computers



1/13/2022

29

SIMD Computers

- The processors operate **synchronously** and a **global clock** is used to ensure **lockstep** operation.
- i.e. at each step (global clock tick) all processors execute the same instruction, each on a different datum

1/13/2022

30

SIMD Computers

- Array processors such as the ICL DAP (Distributed Array Processor)
- and pipelined vector computers such as the CRAY 1 & 2 and CYBER 205 fit into the SIMD category.
- SIMD machines are particularly useful to solve problems which have a regular structure. i.e. the same instruction can be applied to subsets of the data.

1/13/2022 31

Example

- Adding two matrices $A + B = C$.
- Say we have two matrices A and B of order 2 and we have 4 processors.
- $A_{11} + B_{11} = C_{11} \dots A_{12} + B_{12} = C_{12}$
- $A_{21} + B_{21} = C_{21} \dots A_{22} + B_{22} = C_{22}$
- The same instruction is issued to all 4 processors (add the two numbers) and all processors execute the instructions simultaneously.
- It takes one step as opposed to four steps on a sequential machine.

1/13/2022 32

SIMD Computers

- An instruction could be a simple one (eg adding two numbers) or a complex one (eg merging two lists of numbers).
- Similarly the datum may be simple (one number) or complex (several numbers).
- Sometimes it may be necessary to have only a subset of the processors execute an instruction i.e. only some data needs to be operated on for that instruction.

1/13/2022 33

MIMD Computers (multiprocessors /multicomputers)

- This is the most general and most powerful of our classification.
- We have N processors, N streams of instructions and
- N streams of data.

1/13/2022 34

MIMD Computers

The diagram illustrates the architecture of MIMD computers. At the top is a box labeled "SHARED MEMORY OR INTERCONNECTION NETWORK". Below it, three data streams labeled "DS 1", "DS 2", and "DS N" point downwards to three processor boxes: "PROCESSOR ONE", "PROCESSOR TWO", and "PROCESSOR N". Below each processor box is a control unit box labeled "CONTROL ONE", "CONTROL TWO", and "CONTROL N" respectively. Upward-pointing arrows labeled "IS 1", "IS 2", and "IS N" connect each control unit to its corresponding processor. At the bottom, a legend states "DS = DATA STREAM IS = INSTRUCTION STREAM".

1/13/2022 35

MIMD Computers

- Each processor operates under the control of an instruction stream issued by its own control unit.(i.e. each processor is capable of executing its own program on a different data.
- This means that the processors operate **asynchronously** (typically) i.e. can be doing different things on different data at the same time.
- As with SIMD computers communication of data or results between processors can be via a shared memory or interconnection network.

1/13/2022 36

MIMD Computers

- MIMD computers with shared memory are known as **multiprocessors** or **tightly coupled machines**.
- Examples are ENCORE, MULTIMAX, SEQUENT & BALANCE.
- MIMD computers with an interconnection network are known as **multicomputers** or **loosely coupled machines**.
- Examples are INTEL iPSC, NCUBE/7 and transputer networks.

1/13/2022 37

MIMD Computers

- Note:** Multicomputers are sometimes referred to as distributed systems.
- This is **INCORRECT**.
- Distributed systems should, for example, refer to a network of personal workstations (such as SUN's) and
- even though the number of processing units can be quite large the communication in such systems is currently too slow to allow close operation on **one** job.

1/13/2022 38

Potential of the 4 classes

POTENTIAL OF THE 4 CLASSES

$\begin{matrix} I \\ A \ B \end{matrix}$	$\begin{matrix} \text{S(SD)} \\ \text{S(MT)} \end{matrix}$	$\begin{matrix} \rightarrow A + B \\ \rightarrow A + B \\ \rightarrow C + D \end{matrix}$
$\begin{matrix} + \\ A \ B \\ C \ D \end{matrix}$	M(SD)	$\begin{matrix} \rightarrow A + B \\ \rightarrow A * B \\ \rightarrow A * B \end{matrix}$
$\begin{matrix} I \\ + \\ A \ B \end{matrix}$	M(MT)	$\begin{matrix} \rightarrow A + B \\ \rightarrow C + D \\ \rightarrow C * D \end{matrix}$
$\begin{matrix} + \\ + \\ A \ B \\ C \ D \end{matrix}$	M(MT)	$\begin{matrix} \rightarrow A + B \\ \rightarrow C + D \\ \rightarrow C * D \end{matrix}$

1/13/2022 39

Some General Parallel Terminology

- Supercomputing / High Performance Computing (HPC)**
 - Using the world's fastest and largest computers to solve large problems
- Node**
 - A standalone "computer in a box". Usually comprised of multiple CPUs/processors/cores.
 - Nodes are networked together to comprise a supercomputer.

1/13/2022 40

Parallel Terminologies

- CPU / Socket / Processor / Core**
 - In the past, a CPU (Central Processing Unit) was a singular execution component for a computer.
 - Then, multiple CPUs were incorporated into a node.
 - Then, individual CPUs were subdivided into multiple "cores", each being a unique execution unit.

1/13/2022 41

Parallel Terminologies

- CPUs with multiple cores are sometimes called "sockets" - vendor dependent.
- The result is a node with multiple CPUs, each containing multiple cores.

1/13/2022 42

Parallel Terminologies

- **Task**
 - A logically discrete section of computational work.
 - A task is typically a program or program-like set of instructions that is executed by a processor.
 - A parallel program consists of multiple tasks running on multiple processors.

1/13/2022

43

Parallel Terminologies

- **Pipelining**
 - Breaking a task into steps performed by different processor units, with inputs streaming through, much like an assembly line; a type of parallel computing.

1/13/2022

44

Parallel Terminologies

- **Shared Memory**
 - From a strictly hardware point of view, describes a computer architecture where all processors have direct (usually bus based) access to common physical memory.
 - In a programming sense, it describes a model where parallel tasks all have the same "picture" of memory and can directly address and access the same logical memory locations regardless of where the physical memory actually exists.

1/13/2022

45

Parallel Terminologies

- **Symmetric Multi-Processor (SMP)**
 - Hardware architecture where multiple processors share a single address space and access to all resources; shared memory computing.

1/13/2022

46

Parallel Terminologies

- **Distributed Memory**
 - In hardware, refers to network based memory access for physical memory that is not common.
 - As a programming model, tasks can only logically "see" local machine memory and must use communications to access memory on other machines where other tasks are executing.

1/13/2022

47

Parallel Terminologies

- **Communications**
 - Parallel tasks typically need to exchange data. There are several ways this can be accomplished, such as through a shared memory bus or over a network, however the actual event of data exchange is commonly referred to as communications regardless of the method employed.

1/13/2022

48

Parallel Terminologies

- **Synchronization**
 - The coordination of parallel tasks in real time, very often associated with communications.
 - Often implemented by establishing a synchronization point within an application where a task may not proceed further until another task(s) reaches the same or logically equivalent point.

1/13/2022

49

Parallel Terminologies

- Synchronization usually involves waiting by at least one task, and can therefore cause a parallel application's wall clock execution time to increase.
- **Granularity**
 - In parallel computing, granularity is a qualitative measure of the ratio of computation to communication.

1/13/2022

50

Parallel Terminologies

- **Coarse:** relatively large amounts of computational work are done between communication events
- **Fine:** relatively small amounts of computational work are done between communication events

1/13/2022

51

Parallel Terminologies

- **Observed Speedup**
 - Observed speedup of a code which has been parallelized, defined as:

$$\frac{\text{wall-clock time of serial execution}}{\text{wall-clock time of parallel execution}}$$
 - One of the simplest and most widely used indicators for a parallel program's performance.

1/13/2022

52

Parallel Terminologies

- **Parallel Overhead**
 - The amount of time required to coordinate parallel tasks, as opposed to doing useful work.
 - Parallel overhead can include factors such as:
 - Task start-up time
 - Synchronizations
 - Data communications
 - Software overhead imposed by parallel compilers, libraries, tools, operating system, etc.
 - Task termination time

1/13/2022

53

Parallel Terminologies

- **Massively Parallel**
 - Refers to the hardware that comprises a given parallel system - having many processors.
 - The meaning of "many" keeps increasing, but currently, the largest parallel computers can be comprised of processors numbering in the hundreds of thousands.

1/13/2022

54

Parallel Terminologies

- **Embarrassingly Parallel**
 - Solving many similar, but independent tasks simultaneously; little to no need for coordination between the tasks.
- **Scalability**
 - Refers to a parallel system's (hardware and/or software) ability to demonstrate a proportionate increase in parallel speedup with the addition of more processors.

1/13/2022

55

• END

1/13/2022

56