

INTRODUCTION TO HIGH PERFORMANCE COMPUTING PART 1

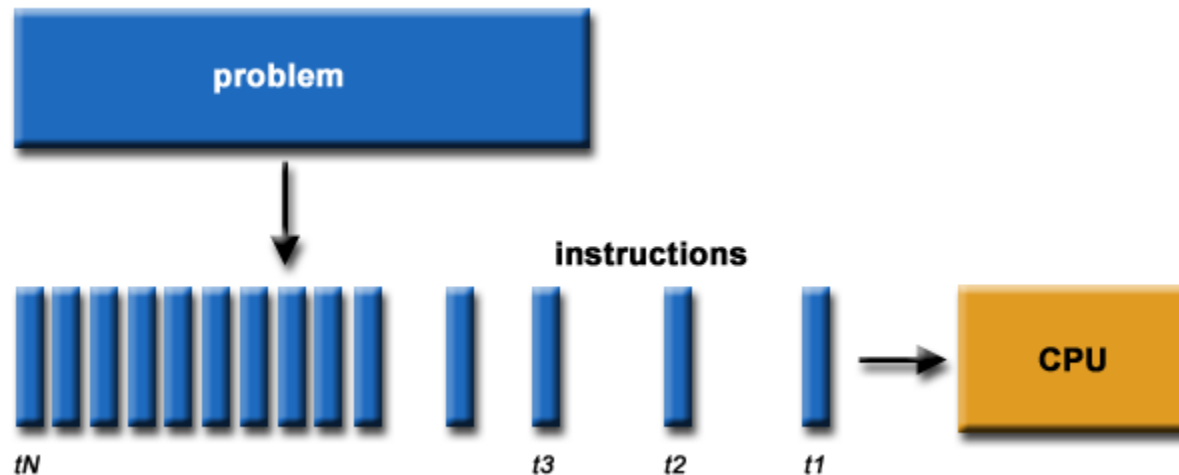
Based on Introduction to Parallel Computing
Lawrence Livermore National Laboratory
https://computing.llnl.gov/tutorials/parallel_comp/

Overview

- This presentation begins a series of presentations – “Introduction to High Performance Computing”. Together over the next few weeks we will journey into the world of parallel and cluster computing.
- We will begin with an overview and some concepts and terminology associated with parallel computing, and then the topics of **parallel memory architectures and programming models** are explored.
- These topics are followed by a discussion on a number of **issues** related to designing parallel programs.
- The last portion of the presentation is spent examining how to parallelize several different types of serial programs.

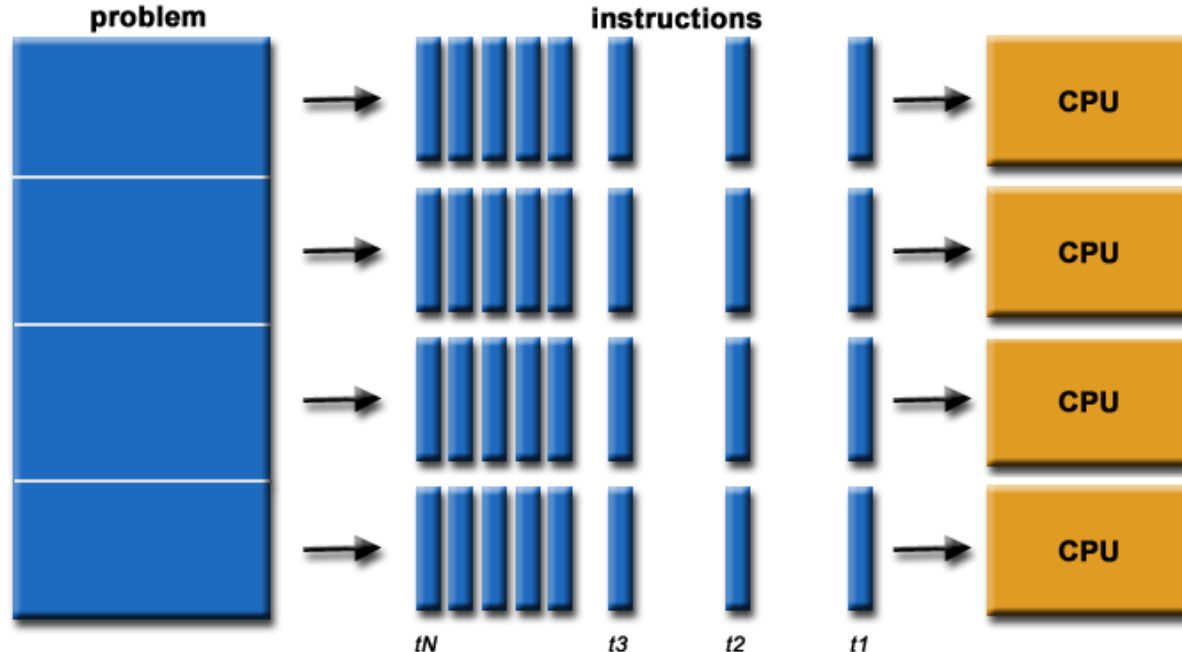
What is Parallel Computing? (1)

- Traditionally, software has been written for ***serial*** (***sometimes called sequential***) computation:
 - To be executed on a single computer having a single Central Processing Unit (CPU);
 - A single 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.



What is Parallel Computing? (2)

- In the simplest sense, **parallel computing** is the simultaneous use of multiple compute resources to solve a **single** computational problem.
 - To be executed using multiple CPUs (or computing resources: GPU, etc.)
 - A single problem is broken into discrete parts that can be solved in parallel
 - Each part is further broken down to a series of instructions
- Instructions from each part execute **simultaneously** on different CPUs (or other hardware that processes instructions)



Parallel Computing: Resources

- The compute resources can include:
 - A single computer with multiple processors (or cores);
 - A single computer with (perhaps multiple) processors and some specialized compute resources (GPU, FPGA, Xeon Phi, etc.);
 - An arbitrary number of computers connected by a network;
 - A combination of any of these.

Parallel Computing: The computational problem

- The computational problem usually demonstrates characteristics such as the ability to be:
 - Broken apart into discrete pieces of work that can be solved simultaneously;
 - Execute multiple program instructions at any moment in time;
 - Solved in **less time** with multiple compute resources than with a single compute resource.
- Not every problem can be solved in parallel
- Many problems that can may not see a benefit
- However, (and fortunately) many natural and scientific problems are naturally pre-disposed to parallelism and can benefit from it.

Parallel Computing: what for? (1)

- Parallel computing is an evolution of serial computing that attempts to emulate what has always been the state of affairs in the natural world: many complex, interrelated events happening at the same time, yet within a sequence.
- Some examples:
 - Planetary and galactic orbits
 - Weather and ocean patterns
 - Tectonic plate drift
 - Rush hour traffic in a city
 - Automobile assembly line
 - Daily operations within a business
 - Building a shopping mall
 - Ordering a hamburger at the drive through!

Parallel Computing: what for? (2)

- Traditionally, parallel computing has been considered to be "the high end of computing" and has been motivated by numerical simulations of complex systems and "Grand Challenge Problems" such as:
 - weather and climate
 - chemical and nuclear reactions
 - biological, human genome
 - geological, seismic activity
 - mechanical devices - from prosthetics to spacecraft
 - electronic circuits
 - manufacturing processes

Parallel Computing: what for? (3)

- Today, commercial applications are providing an equal or greater driving force in the development of faster computers. These applications require the processing of large amounts of data in sophisticated ways. Example applications include:
 - parallel databases, data mining
 - oil exploration
 - web search engines, web based business services
 - computer-aided diagnosis in medicine
 - management of national and multi-national corporations
 - advanced graphics and virtual reality, particularly in the entertainment industry
 - networked video and multi-media technologies
 - collaborative work environments
- *Ultimately, parallel computing aims to maximize the infinite but seemingly scarce commodity called time.*
- *What we are ultimately concerned with is “wall-clock” time.*

Why Parallel Computing? (1)

- This is a legitime question! Parallel computing is complex in many aspects!
- The primary reasons for using parallel computing:
 - Save time - wall clock time
 - Solve larger problems
 - Provide parallelism (do multiple things at the same time)

Why Parallel Computing? (2)

- Other reasons might include:
 - Taking advantage of non-local resources - using available compute resources on a wide area network, or even the Internet when local compute resources are scarce.
 - Cost savings - using multiple "cheap" computing resources instead of paying for time on a supercomputer.
 - Overcoming memory constraints - single computers have very finite memory resources. For large problems, using the memories of multiple computers may overcome this obstacle.

Limitations of Serial 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.
- **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.
- **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.
- **Space, Cost and Heat** - Placing faster components within close proximity to each other (for speed) poses significant engineering challenges. It gets difficult, expensive and hot to do this!

The future

- during the past 10 years, the trends indicated by ever faster networks, distributed systems, and multi-processor computer architectures (even at the desktop level) clearly show that ***parallelism is the future of computing***.
 - In fact, parallel computing is so commonplace that it seems to get lost in the crowd. You've heard all the headlines on *big data*, *cloud computing*, *machine learning*, etc. etc.
 - Most of the time these technologies explicitly depend upon parallel computing – so much so, that nobody even bothers saying it. It's "obvious".
- The future will take many forms including mixing general purpose solutions (your PC...) and very speciliazed solutions as GPUs, Xeon Phi, FPGAs, etc.

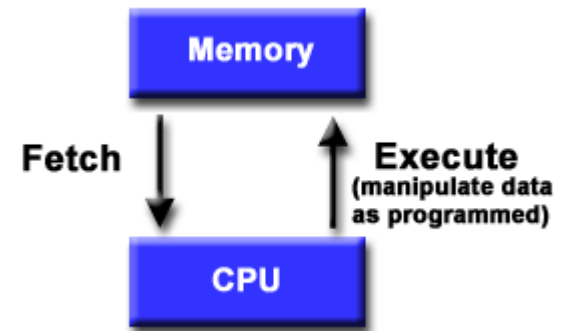
CONCEPTS AND TERMINOLOGY

Von Neumann Architecture

- For over 40 years, virtually all computers have followed a common machine model known as the von Neumann computer. Named after the Hungarian mathematician John von Neumann.
- A von Neumann computer uses the stored-program concept. The CPU executes a stored program that specifies a sequence of read and write operations on the memory.

Basic Design

- Basic design
 - Memory is used to store both program and data instructions
 - Program instructions are coded data which tell the computer to do something
 - Data is simply information to be used by the program
- A central processing unit (CPU) gets instructions and/or data from memory, decodes the instructions and then ***sequentially*** performs them.



Flynn's Classical Taxonomy

- There are different ways to classify parallel computers. One of the more widely used classifications, in use since 1966, is called **Flynn's Taxonomy**.
- Flynn's taxonomy distinguishes multi-processor computer architectures according to how they can be classified along the two independent dimensions of ***Instruction*** and ***Data***. Each of these dimensions can have only one of two possible states: ***Single*** or ***Multiple***.

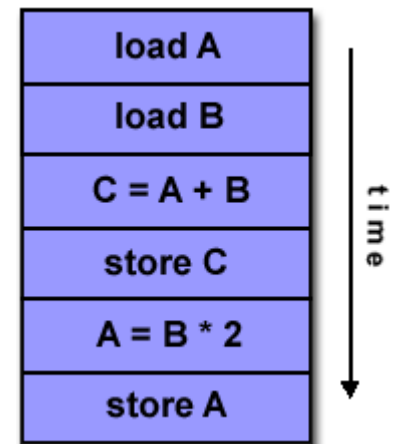
Flynn Matrix

- The matrix below defines the 4 possible classifications according to Flynn's Taxonomy
- We will discuss each of these in turn
- It is important to remember that this is a high-level way of looking at parallel computing – there are 'gray areas'
 - However it is a very useful categorization that is still used today

S I S D Single Instruction, Single Data	S I M D Single Instruction, Multiple Data
M I S D Multiple Instruction, Single Data	M I M D Multiple Instruction, Multiple Data

Single Instruction, Single Data (SISD)

- **A serial (non-parallel) computer**
- Single instruction: only one instruction stream is being acted on by the CPU during any one clock cycle
- Single data: only one data stream is being used as input during any one clock cycle
- **Deterministic execution**
 - This may sound obvious, but wait until you see the others!
- This is the oldest and until recently, the most prevalent form of computer architecture
- Examples: most PCs, single CPU workstations and mainframes (yes, these still exist, although the term has become less popular).

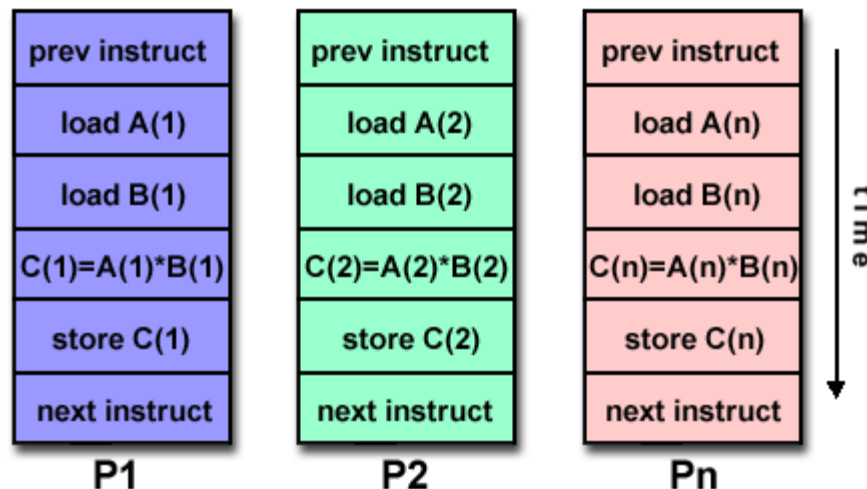


Single Instruction, Multiple Data (SIMD)

- **A type of parallel computer**
- Single instruction: All processing units execute the same instruction at any given clock cycle
- Multiple data: Each processing unit can operate on a different data element
- This type of machine typically has an instruction dispatcher, a very high-bandwidth internal network, and a very large array of very small-capacity instruction units

Single Instruction, Multiple Data (SIMD)

- Best suited for specialized problems characterized by a high degree of regularity, such as image processing.
- Synchronous (lockstep) and deterministic execution
- Two varieties: Processor Arrays and Vector Pipelines
- Examples:
 - Processor Arrays: Connection Machine CM-2, Maspar MP-1, MP-2
 - Vector Pipelines: IBM 9000, Cray C90, Fujitsu VP, NEC SX-2, Hitachi S820

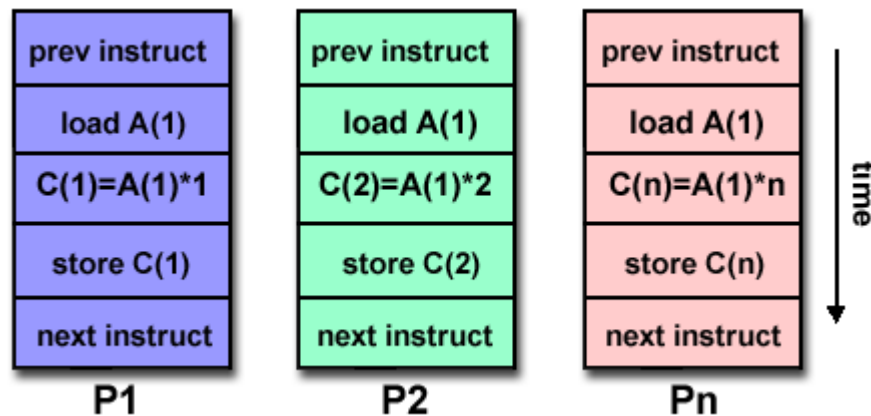


Multiple Instruction, Single Data (MISD)

- A single data stream is fed into multiple processing units.
- Each processing unit operates on the data independently via independent instruction streams.
- Few actual examples of this class of parallel computer have ever existed. One is the experimental Carnegie-Mellon C.mmp computer (1971).
- This does not mean that MISD is dead though. All we are waiting on is a *paradigm shift* that allows MISD to become a viable reality
 - Think, for example, quantum computing?

Multiple Instruction, Single Data (MISD)

- Some conceivable uses might be:
 - multiple frequency filters operating on a single signal stream
- multiple cryptography algorithms attempting to crack a single coded message.

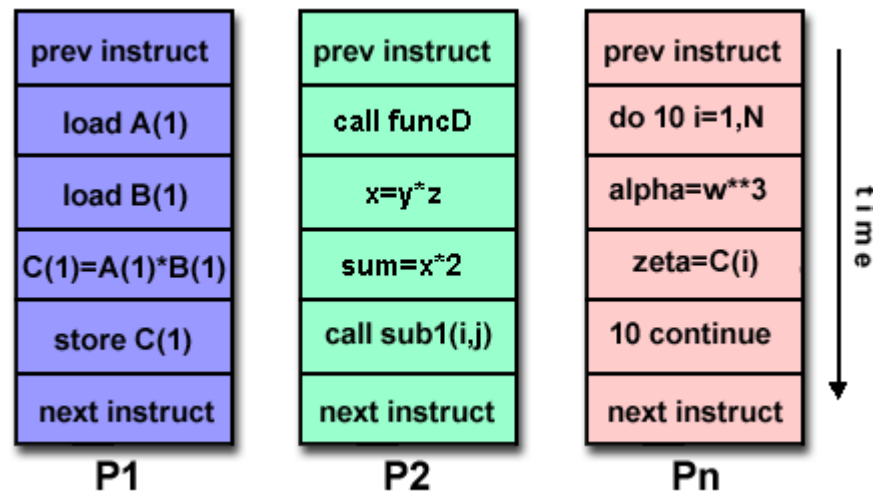


Multiple Instruction, Multiple Data (MIMD)

- Currently, the most common type of parallel computer. Most modern computers fall into this category.
- Multiple Instruction: every processor may be executing a different instruction stream
- Multiple Data: every processor may be working with a different data stream
- Remember however, that to meet the traditional definition of parallel computing, these multiple data streams must all be part of the same (big) problem!

Multiple Instruction, Multiple Data (MIMD)

- **Execution can be synchronous or asynchronous, deterministic or non-deterministic**
- Examples: most current supercomputers, networked parallel computer "grids" and multi-processor SMP computers - including some types of PCs.



Some General Parallel Terminology

Like everything else, parallel computing has its own "jargon". Some of the more commonly used terms associated with parallel computing are listed below. Most of these will be discussed in more detail later.

- **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.

- **Parallel Task**

- A task that can be executed by multiple processors safely (yielding correct results)

- **Serial Execution**

- Execution of a program sequentially, one statement at a time. In the simplest sense, this is what happens on a one processor machine. **However, virtually all parallel tasks will have sections of a parallel program that must be executed serially. We will discuss this more later.**

Some General Parallel Terminology

- **Parallel Execution**

- Execution of a program by more than one task, with each task being able to execute the same or different statements at the same moment in time.

- **Shared Memory**

- From a strictly hardware point of view, describes a computer architecture where all processors have **direct** (usually bus based) access to **common (or shared) 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.

- **Distributed Memory**

- In hardware, refers to network based memory access for physical **memory that is not common (or shared) directly**. 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.

Some General Parallel Terminology

- **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.
 - ***Tasks that do not need to exchange data (or much data) are sometimes called ‘embarrassingly parallel’. More on this later.**

- **Synchronization**

- The coordination of parallel tasks in real time, very often associated with communications. Commonly 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. **We will see how this is done in code later.**
- Synchronization usually involves waiting by at least one task, and can therefore cause a parallel application's wall clock execution time to increase.

Some General Parallel Terminology

- **Granularity**

- In parallel computing, granularity is a qualitative measure of the ratio of computation to communication.
- **Coarse:** relatively large amounts of computational work are done between communication events. Generally, less comms than fine.
- **Fine:** relatively small amounts of computational work are done between communication events. Generally, more comms than coarse.

- **Observed Speedup**

- Observed speedup of a code which has been parallelized can be described (roughly) as comparing:
 - wall-clock time of serial execution
 - wall-clock time of parallel execution
- One of the simplest and most widely used indicators for a parallel program's performance. **We will revisit this in detail later.**

Some General Parallel Terminology

- **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

- **Massively Parallel**

- Refers to the hardware that comprises a given parallel system - having “many” processors. The meaning of “many” keeps increasing, but currently, “many” means tens of thousands of processors or more.

Some General Parallel Terminology

- **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. Factors that contribute to scalability include:
 - Hardware - particularly memory-cpu bandwidths and network communications
 - Application algorithm
 - Parallel overhead related issues
 - Characteristics of your specific application and coding
- The important thing to remember is that scalability is not a given
 - If it works on two processors, or two hundred, or two thousand, it might not work on 50,000
 - This can be due to issues with hardware, software, the problem itself, etc.