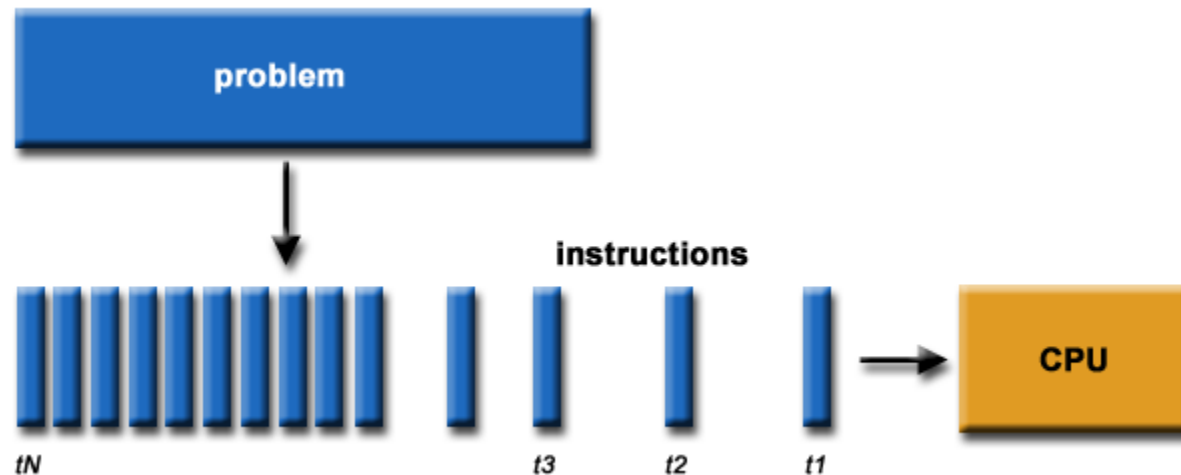# Introduction to Parallel Computing
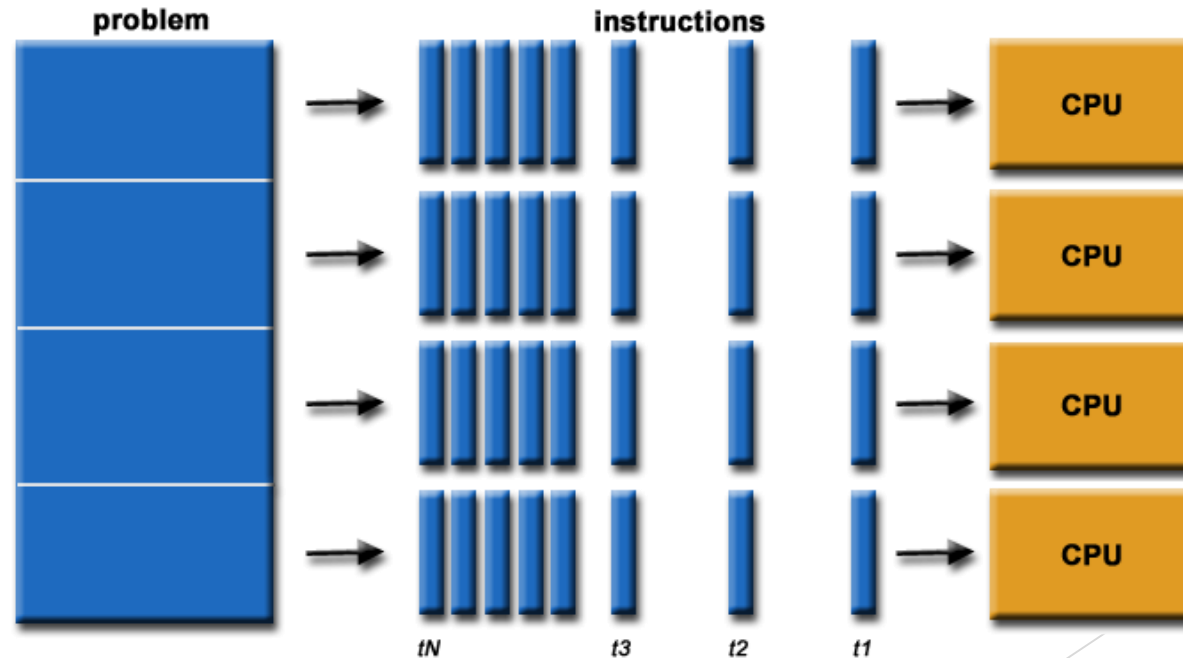
# What is Parallel Computing? (1)

- 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.
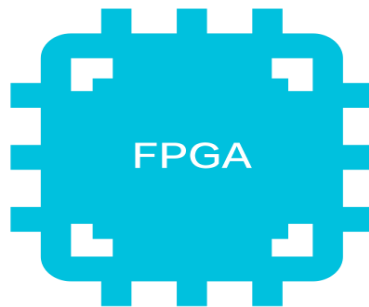  - Only one instruction may execute at any moment in time.

# What is Parallel Computing? (2)

▶ In the simplest sense, *parallel computing* is the simultaneous use of multiple compute 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

# Parallel Computing: Resources

▶ The compute resources can include:

    ▶ A single computer with multiple processors;

    ▶ A single computer with (multiple) processor(s) and some specialized computer resources (GPU, FPGA ...)

    ▶ An arbitrary number of computers connected by a network;

    ▶ A combination of both.

# 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.

# Parallel Computing: what for? (3)

- Example applications include:

    - parallel databases, data mining

    - web search engines, web based business services

    - computer-aided diagnosis in medicine

    - advanced graphics and virtual reality, particularly in the entertainment industry
    - networked video and multi-media technologies
- Ultimately, parallel computing is an attempt to maximize the infinite but seemingly scarce commodity called time.

# Why Parallel Computing?

- This is a legitime question! Parallel computing is complex on any aspect!

- The primary reasons for using parallel computing:
  - Save time - wall clock time
  - Solve larger problems
  - Provide concurrency (do multiple things at the same time)

# 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.

▶ 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.

# Flynn Taxanomy

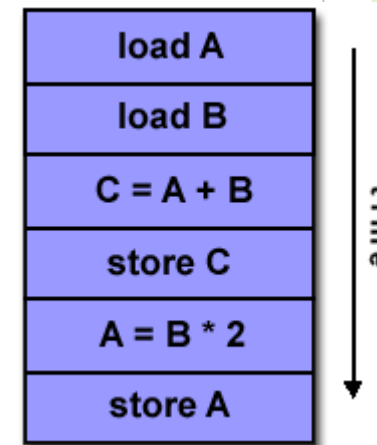▶ The matrix below defines the 4 possible classifications according to Flynn

| SISD | SIMD |
|------|------|
| Single Instruction, Single Data | Single Instruction, Multiple Data |
| MISD | MIMD |
| Multiple Instruction, Single Data | Multiple Instruction, Multiple Data |

# Flynn Taxanomy

# 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 is the oldest and until recently, the most prevalent form of computer

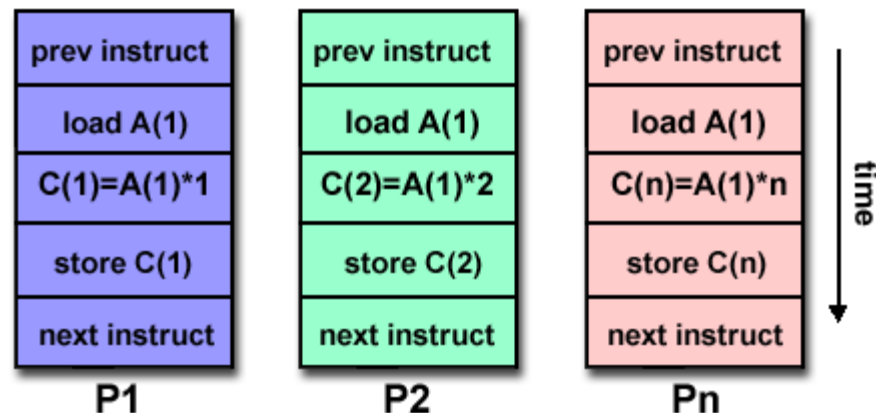▶ Examples: most PCs, single CPU workstations and mainframes

# 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 ̶
- This type of machine typically has an instruction dispa̶ bandwidth internal network, and a very large array of̶ instruction units.
- Best suited for specialized problems characterized by̶ 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-̶
  - Vector Pipelines: IBM 9000, Cray C90, Fujitsu VP, NEC SX̶

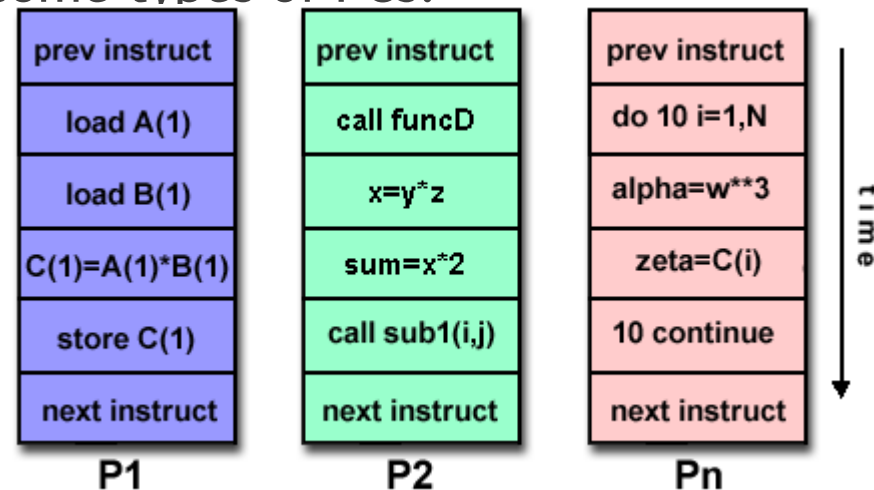| P1 | P2 | Pn |
|---|---|---|
| prev instruct | prev instruct | prev instruct |
| load A(1) | load A(2) | load A(n) |
| load B(1) | load B(2) | load B(n) |
| C(1)=A(1)*B(1) | C(2)=A(2)*B(2) | C(n)=A(n)*B(n) |
| store C(1) | store C(2) | store C(n) |
| next instruct | next instruct | next instruct |

time

# 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).

▶ 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

- 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.

| prev instruct | prev instruct | prev instruct |
|---|---|---|
| load A(1) | call funcD | do 10 i=1,N |
| load B(1) | x=y*z | alpha=w**3 |
| C(1)=A(1)*B(1) | sum=x*2 | zeta=C(i) |
| store C(1) | call sub1(i,j) | 10 continue |
| next instruct | next instruct | next instruct |
| P1 | P2 | Pn |

# Some General Parallel Terminology

- **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 (yields 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.

# 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 statement 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 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. 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.

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

# Granularity

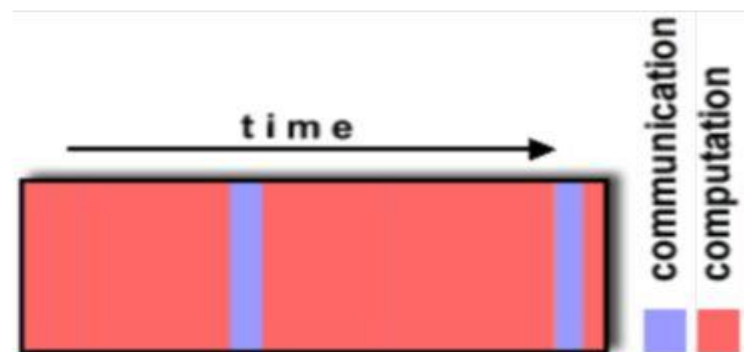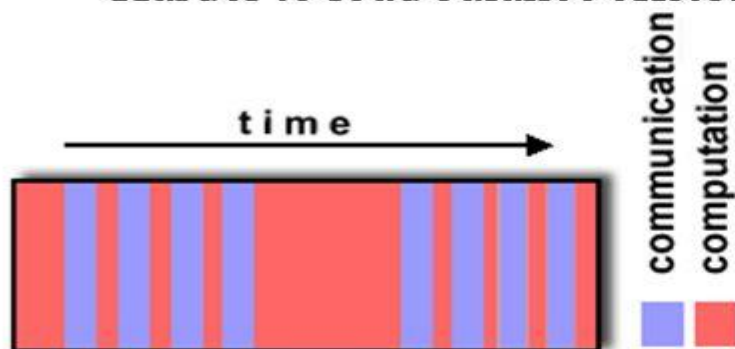Granularity is the ratio of computation to communication.

Periods of computation are typically separated from periods of communication by synchronization events.

**Fine-grain Parallelism:** Relatively small amounts of computational work are done between communication events.

Facilitates load balancing and Implies high communication overhead and less opportunity for performance enhancement

**Coarse-grain Parallelism:** Relatively large amounts of computational work are done between communication/synchronization events.
Harder to load balance efficiently

# 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 BG/L pushes this number to 6 digits.

# Some General Parallel Terminology
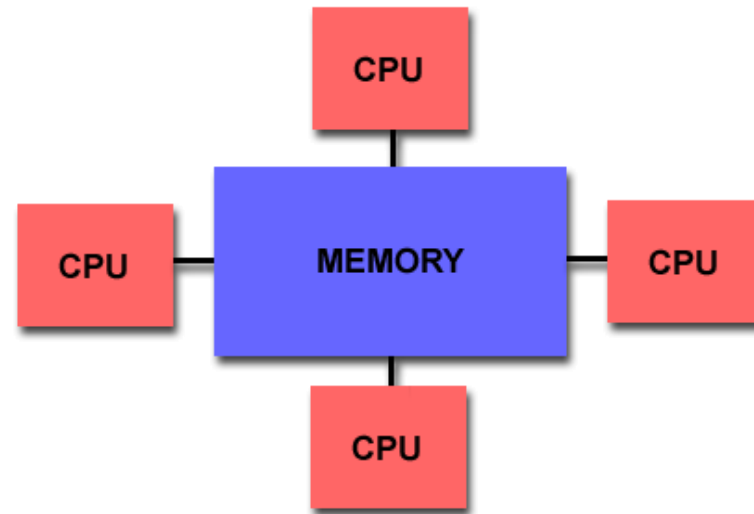
- **Scalability**
  - Refers to a parallel system's (hardware and/or software) ability to demonstrate a proportionate increase in parallel speedup with the addition of more processors. Factors that contribute to scalability include:
    - Hardware - particularly memory-cpu bandwidths and network communications
    - Application algorithm
    - Parallel overhead related
    - Characteristics of your specific application and coding

# Shared Memory

▶ Shared memory parallel computers vary widely, but generally have in common the ability for all processors to access all memory as global address space.
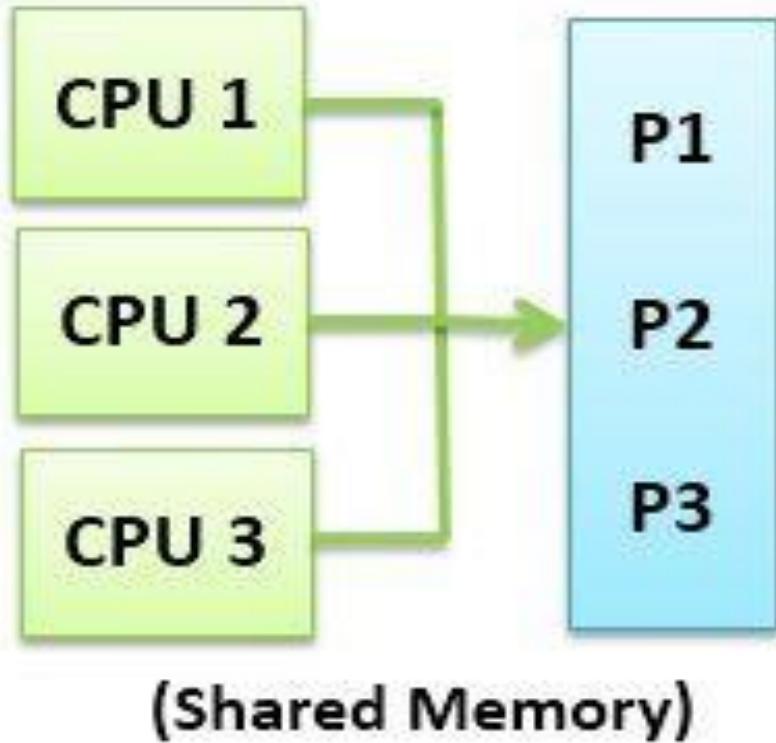


▶ Multiple processors can operate independently but share the same memory resources.

▶ Changes in a memory location effected by one processor are visible to all other processors.

▶ Shared memory machines can be divided into two main classes based upon memory access times: *UMA* and *NUMA*.

# Shared Memory : UMA vs. NUMA

- Uniform Memory Access (UMA):
  - Most commonly represented today by Symmetric Multiprocessor (SMP) machines
  - Identical processors
  - Equal access and access times to memory
  - Sometimes called CC-UMA - Cache Coherent UMA. Cache coherent means if one processor updates a location in shared memory, all the other processors know about the update. Cache coherency is accomplished at the hardware level.
- Non-Uniform Memory Access (NUMA):
  - Often made by physically linking two or more SMPs
  - One SMP can directly access memory of another SMP
  - Not all processors have equal access time to all memories
  - Memory access across link is slower
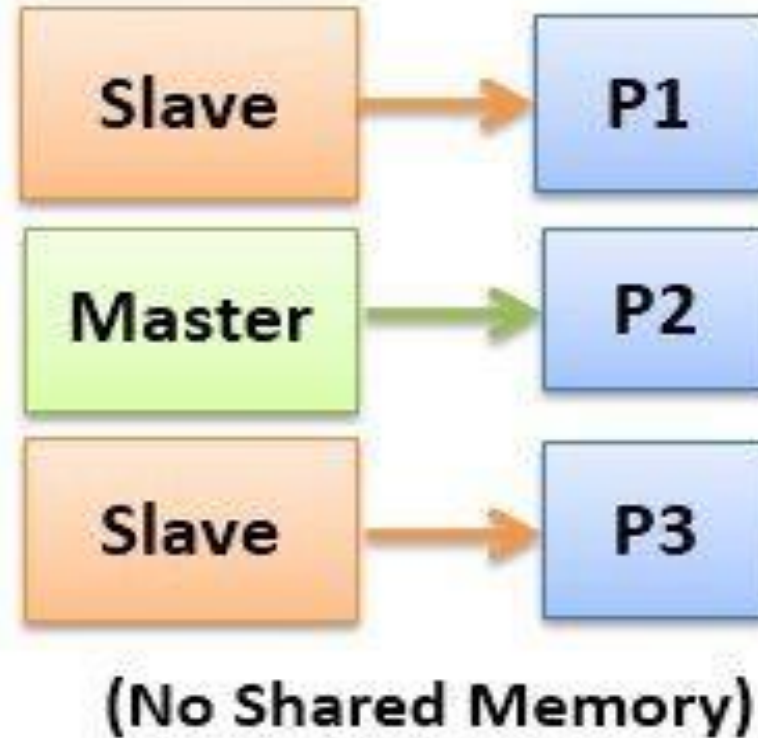  - If cache coherency is maintained, then may also be called CC-NUMA - Cache Coherent NUMA

# Shared Memory: Pro and Con
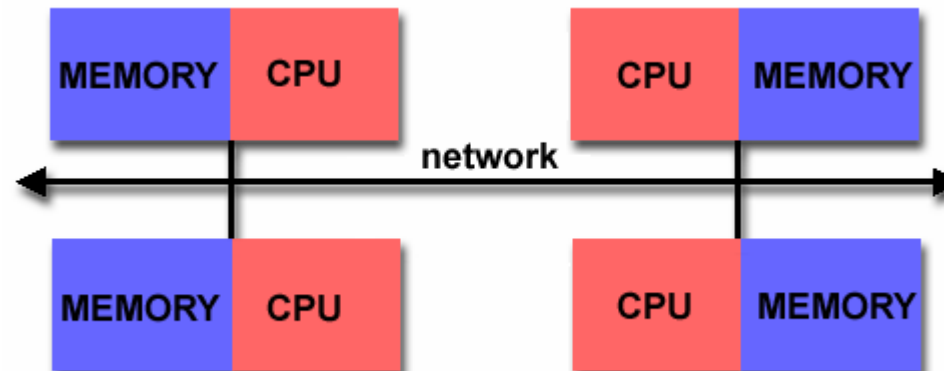
- ▶ Advantages
  - ▶ Global address space provides a user-friendly programming perspective to memory
  - ▶ Data sharing between tasks is both fast and uniform due to the proximity of memory to CPUs
- ▶ Disadvantages:
  - ▶ Primary disadvantage is the lack of scalability between memory and CPUs. Adding more CPUs can geometrically increases traffic on the shared memory-CPU path, and for cache coherent systems, geometrically increase traffic associated with cache/memory management.
  - ▶ Programmer responsibility for synchronization constructs that insure "correct" access of global memory.
  - ▶ Expense: it becomes increasingly difficult and expensive to design and produce shared memory machines with ever increasing numbers of processors.

# Distributed Memory

▶ Like shared memory systems, distributed memory systems vary widely but share a common characteristic. Distributed memory systems require a communication network to connect inter-processor memory.

▶ Processors have their own local memory. Memory addresses in one processor do not map to another processor, so there is no concept of global address space across all processors.

▶ Because each processor has its own local memory, it operates independently. Changes it makes to its local memory have no effect on the memory of other processors. Hence, the concept of cache coherency does not apply.

▶ When a processor needs access to data in another processor, it is usually the task of the programmer to explicitly define how and when data is communicated. Synchronization between tasks is likewise the programmer's responsibility.

▶ The network "fabric" used for data transfer varies widely, though it can can be as simple as Ethernet.

# Distributed Memory: Pro and Con

▶ Advantages
  ▶ Memory is scalable with number of processors. Increase the number of processors and the size of memory increases proportionately.
  ▶ Each processor can rapidly access its own memory without interference and without the overhead incurred with trying to maintain cache coherency.
  ▶ Cost effectiveness: can use commodity, off-the-shelf processors and networking.

▶ Disadvantages
  ▶ The programmer is responsible for many of the details associated with data communication between processors.
  ▶ It may be difficult to map existing data structures, based on global memory, to this memory organization.
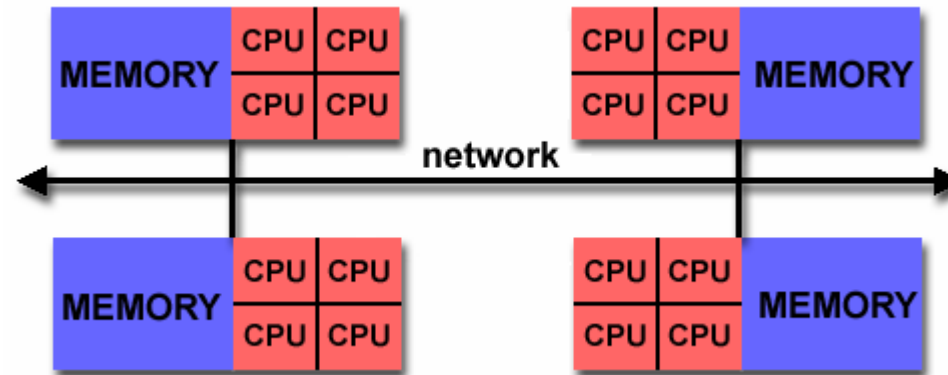  ▶ Non-uniform memory access (NUMA) times

# Hybrid Distributed-Shared Memory

Summarizing a few of the key characteristics of shared and distributed memory machines

| Comparison of Shared and Distributed Memory Architectures | | |
|---|---|---|
| **Architecture** | CC-UMA | CC-NUMA | Distributed |
| **Examples** | SMPs<br>Sun Vexx<br>DEC/Compaq<br>SGI Challenge<br>IBM POWER3 | Bull NovaScale<br>SGI Origin<br>Sequent<br>HP Exemplar<br>DEC/Compaq<br>IBM POWER4 (MCM) | Cray T3E<br>Maspar<br>IBM SP2<br>IBM BlueGene |
| **Communications** | MPI<br>Threads<br>OpenMP<br>shmem | MPI<br>Threads<br>OpenMP<br>shmem | MPI |
| **Scalability** | to 10s of processors | to 100s of processors | to 1000s of processors |
| **Draw Backs** | Memory-CPU bandwidth | Memory-CPU bandwidth<br>Non-uniform access times | System administration<br>Programming is hard to develop and maintain |
| **Software Availability** | many 1000s ISVs | many 1000s ISVs | 100s ISVs |

Note: the header "Architecture" row has four columns (Architecture, CC-UMA, CC-NUMA, Distributed) while the title row spans all columns.

# Hybrid Distributed-Shared Memory

▶ The largest and fastest computers in the world today employ both shared and distributed memory architectures.



▶ The shared memory component is usually a cache coherent SMP machine. Processors on a given SMP can address that machine's memory as global.

▶ The distributed memory component is the networking of multiple SMPs. SMPs know only about their own memory - not the memory on another SMP. Therefore, network communications are required to move data from one SMP to another.

▶ Current trends seem to indicate that this type of memory architecture will continue to prevail and increase at the high end of computing for the foreseeable future.

▶ Advantages and Disadvantages: whatever is common to both shared and distributed memory architectures.