

Module-1.

Page No. _____
Date / /

1) Explain Flynn's taxonomy of Parallel computer with suitable example

=>

Flynn's taxonomy classifies a parallel computer according to the number of instruction streams and

The number of the data streams (or data paths) it can simultaneously manage. It categorizes into four.

Four types:

i) SISD (single instruction, single data)

This architecture processes a single instruction stream and a single data stream. It represents traditional non-parallel computers where a single processor executes instruction sequentially on one piece of data at a time.

ex:- A classic uniprocessor computer like an older PC running a single program without parallel processing

Features:

- i) SEMD (single instruction, multiple data) In this architecture, a single instruction stream operates on multiple data streams simultaneously. All processors execute the same instruction

- carry out different data elements at the same time. This particularly efficient tasks involving repetitive operation on large datasets.

ex: - Graphics processing units used in

modern computers for rendering graphics or performing scientific computations. coherent same operation is applied to many pixels concurrently.

iii) MISD (multiple instruction, single data)

This architecture involves multiple instruction stream operation on a single data stream.

This is a less common architecture type. It implies different processors performing different operations on the same data simultaneously.

ex: - while rare in general purpose computing some specialized fault tolerant systems or pipeline architectures might exhibit MISD characteristics.

where different stages of pipeline operate on the same data in sequence but with

different instructions

iv) MIMD (multiple instruction, multiple data) This is the most common type of parallel architecture, a single instruction stream operates where multiple processors can execute different data streams simultaneously. This allows for true parallel processing, where independent tasks can run concurrently.

ex: - modern multi-core processors (eg. Intel, AMD, Ryzen) or large-scale supercomputers and distributed computing systems like server farms, where multiple independent programs or threads can run in parallel on different data sets.

By compare and contrast SIMD and MIMD systems in terms of architecture, performance and applications

Terms	Feature	SIMD	MIMD
Architecture	Instruction Stream	single instruction stream	multiple instruction stream
	Data Stream	multiple data stream	multiple data streams

Terms	Aspect	Page No.	
		Date	Page No.
Control Unit	centralized	independent	
	control for all processing units	control units for each processor	
Synchronization	highly synchronous	Asynchronous	
	3rd	complex coordination	
Performance & efficiency		High for uniform tasks	High for diverse tasks
Scalability	SEMD	MEMD	
	Limited by uniformity of operations	Highly scalable across heterogeneous tasks	
Overhead		Low control overhead	Higher due to independent instructions
Parallelism		Data-level parallelism	Task-level parallelism
Application	Image/Video processing	Excellent for pixel-wise operations	Less efficient unless tasks vary significantly

Terms	Aspect	Page No.	
		Date	Page No.
Control Unit	centralized	independent	
	control for all processing units	control units for each processor	
Synchronization	highly synchronous	Asynchronous	
	3rd	complex coordination	
Performance & efficiency		High for uniform tasks	High for diverse tasks
Scalability	SEMD	MEMD	
	Limited by uniformity of operations	Highly scalable across heterogeneous tasks	
Overhead		Low control overhead	Higher due to independent instructions
Parallelism		Data-level parallelism	Task-level parallelism
Application	Image/Video processing	Excellent for pixel-wise operations	Less efficient unless tasks vary significantly

⑧ Discuss the architecture and working of a vector processor and working of a vector processor. How do vector registers and retrieved memory enhance performance?

Ans:- Typical recent systems have the following characteristics

4 Vector registers:

* These are registers capable of storing a vector of operands & operating simultaneously on their contents.

* The vector length is fixed by the system and can range from 4 to 856 64-bit elements.

=> vectorised and pipelined functional units:

* Note that the same operation is applied to each element in the vector, or in the case of operations like addition, the same operation is applied to each pair of corresponding elements in the two vectors. Thus vector operations are SIMD.

=> Vector instructions: These are instructions that operate on vectors rather than scalars. If the vector length is vector length these instructions have the great virtual that a sample loop such as
`for(i=0; i<n; i++)
X[i] += y[i];`

requires only a single load, add and store for each block of vector-length elements while a conventional system requires a load add

=> Interleaved memory: The memory system consists of multiple "banks" that can be accessed in parallel or less independently. After accessing one bank, the next will be ready before it can be reaccessed, but as if next bank can be accessed much sooner.

=> Strided memory access and hardware scatter gather:
In strided memory, access the program accesses elements of a vector located at fixed intervals. Scatter/gather (in this context) is writing (scatter) or reading (gather) elements of a vector located at irregular intervals.

4) With the help of a diagram, explain the working of shared-memory and distributed-memory system. Highlight their advantages and limitations.

Shared-memory system :- In a shared memory system, multiple processors share a common memory space.

each processors can access any memory location, and changes made by one processor are visible to all others.

Working:

1) processors Access shared memory: each processor can access any memory location in the shared memory space.

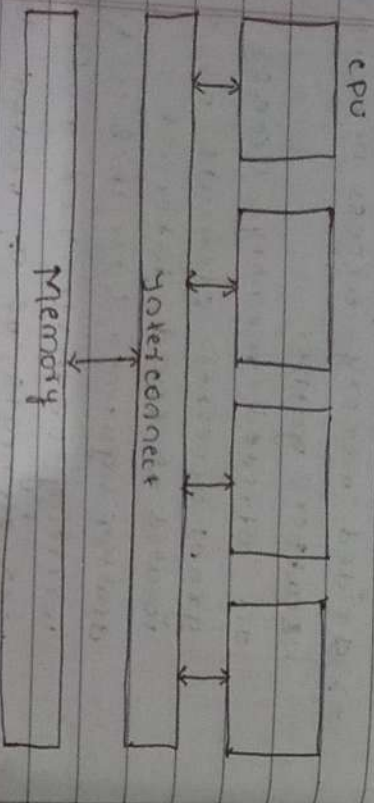


Fig:- A shared memory system.

2) cache coherence: To maintain data consistency, cache coherence protocols are used to ensure that change made by one processor are reflected in the cache of other processors.

3) Synchronization: processors use synchronization mechanisms such as locks or semaphores to coordinate

access to shared resources.

Advantages:

- Shared memory system are relatively easy to program.
- Since processors can access shared memory, directly communication overhead is relatively low.

Limitations:

- Shared memory systems can be limited in scalability.
- The shared memory can become a bottleneck, limiting the performance of the system.

Distributed-memory system:- Each processor has its own private memory space. Processors communicate with each other by exchanging message over an interconnect network.

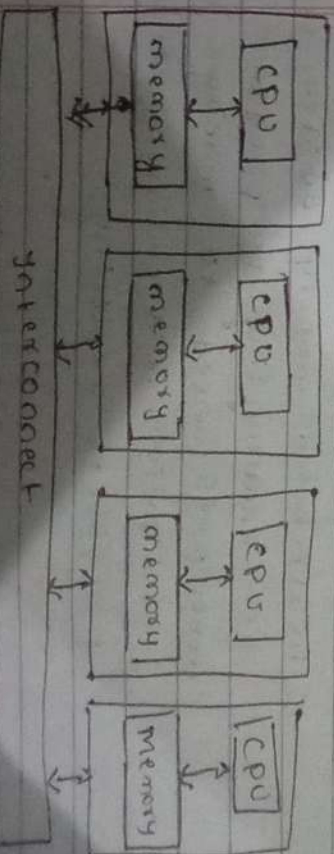


Fig:- A distributed memory system.

Working:-

Its processors communicate via message. Processors exchange message over the interconnect n/w to communicate with each other.

2) Data distribution: Data is distributed across processors and each processor operates on its local data.

3) Message passing: processors use message-passing protocols to exchange data and coordinate their actions.

Advantages

→ Distributed memory system extendable to a large no of processors.

→ Distributed memory system can be designed to accommodate different types of processors of interconnect n/w.

Limitations

→ Distributed memory system can be more challenging to program as developers need to manage data distribution and communication b/w processors explicitly.

→

The communication overhead can be significant especially for application with complex communication patterns.

5)

Describe various interconnection networks used in parallel systems. compare bus, crossbar, ring, mesh and hypercube topologies with respect to scalability & performance.

Ans:-

Interconnection n/w in parallel systems enable communication b/w processors memory and other components.

The various interconnection n/w's

→ crossbar network: Allows any-to-any connections b/w processors and memory high performance but costly with p^2 switches for p processors.

→ Omega network: A multistage n/w using 2×2 switches. less expensive than crossbar with approx $\log p$ switches.

→ Mesh network: processors connected in a grid good for local communication patterns.

→ Hypercube network: connects processors with $\log p$ neighbors.

for P processors use full for loop parallel algorithms.

→ Bus network: A shared channel for communication. Simple but can be a bottleneck.

→ Ring network: Processors connected in a loop. Moderate scalability and performance.

Comparison of bus crossbar, ring, mesh and hypercube topologies regarding scalability and performance in parallel systems

i) Bus

Scalability: Low, limited by bandwidth and contention

Performance: Low for many processors due to shared channel.

ii) Crossbar

Scalability: Limited due to point-to-point connections for processors.

Performance: High any to any connections possible without contention.

iii) Ring:

Scalability: Moderate, easy to add nodes but performance degrades with size.

Performance: Moderate, data must traverse node to reach destination.

Mesh Scalability: Good for 2D/3D structures, local connections scale well.

Performance: Good for local communication patterns, otherwise, latency increases with distance.

Hypercube:-

Scalability: High for certain algorithms, reversing its connectivity.

Performance: High for algorithms matching hypercube structure, log steps for P processors.

Q6)

Explain the cache coherence problem in shared memory systems? Discuss

snooping and directory-based cache coherence protocols with example. Cache coherence problem:-

In shared memory systems, multiple processors access a common memory space. Each processor type

-caching has its own cache to improve performance by reducing memory access latency. However, when multiple processors cache the same memory location, inconsistencies can arise if one processor modifies its cached copy. This leads to the cache coherence problem.

Key issues:-

- 1) Data inconsistency :- when one processor updates a cache copy of memory location, other processors might still have outdated copies in their caches.
- 2) stale data :- processors might read stale data from their cache instead of the updated value.

Snooping cache coherence protocol
In snooping protocols, each cache monitors the shared bus for memory transactions. When a cache sees a write to a location it has cached, it updates or invalidates its copy accordingly.

Example :-

- 1) Initially, both processor X (P_X) and processor A (P_A) cache the value of X as 10.
- 2) P_X updates X to 20 in its cache and broadcasts this change on the bus.
- 3) P_A's cache snoops the bus, sees the updates, and invalidates its copy of X.

Directory based cache coherence protocol
In directory-based protocols, a centralized directory keeps track of which caches have copies of each memory location and their states.

Ex:-

- 1) The directory tracks that both P₁ and P₂ have cached X.
- 2) When P₁ wants to write to X, it sends a request to the directory.
- 3) The directory sends invalidate messages to P₂, which then invalidates its copy of X.

Module - 2

1) Explain the concept of GPU programming

-> how does it differ from traditional CPU programming?

Discuss challenges such as branching and scheduling in GPU threads

Ans:-

GPU (Graphics processing unit) programming involves using GPUs to perform general purpose computing tasks beyond graphics rendering. GPUs are designed for parallel processing with thousands of cores that can handle many threads simultaneously.

-> Differences from CPU programming

How does it differ from traditional CPU programming?

-> parallelism: GPUs are optimized for massive parallelism, whereas CPUs are designed for serial execution.

1) Core Architecture: GPUs consist of numerous smaller CPU cores.

2) Memory hierarchy: CPUs have a different memory hierarchy with a focus on high bandwidth memory.

Challenges in GPU threads

-> Branching

Branch divergence: In a GPU, threads are executed in groups called warps. When threads within a warp take different branches during execution, this reduces performance as the GPU must execute each branch path separately, leading to idle threads.

Impact: Branch divergence can significantly reduce GPU performance, especially if the branches are complex or have many divergent paths.

Mitigation: Minimize branching within warps or use techniques like predication where each thread conditionally executes instruction based on a predicate.

Scheduling

Thread scheduling: GPUs use a GWT (single instruction, multiple thread) execution model, efficient scheduling of the threads is crucial to maximize performance.

Page No. _____
Date / /

Occupancy: refers to the no of active warps per multiprocessor. high occupancy is essential for hiding memory latency and maximizing throughput.

Mitigation: optimize thread block sizes, minimize thread divergence and ensure efficient memory access patterns to improve occupancy and overall performance.

② what are hybrid programming model in parallel computing? explain with an example how shared-memory and distributed-memory APIs can be combined.

Ans:- Hybrid programming models in parallel computing combine multiple parallel programming paradigms to achieve better performance scalability and efficiency. These models integrate different programming models, such as shared-memory and distributed memory models, to leverage their strengths types of hybrid models-

1) Shared, memory and distributed

Page No. _____
Date / /

memory: combined models like openmp (shared-memory) with mpiccs (distributed-memory) to parallelize computations within and across nodes.

As GPU and CPU: use models like CUDA (for NVIDIA GPUs), or open (for multiple platforms) in conjunction with CPU-based parallel programming models like openmp or mpi.

Combining shared memory and distributed memory APIs.

→ Let's consider an example using openmp (shared memory) and mpi (distributed memory) to parallelize a computation.

Example: parallelizing a loop

Suppose we have a large loop that needs to be parallelized across multiple nodes in a cluster with each node having multiple cores.

1) MPI (distributed memory): use MPI to divide the loop iterations among multiple nodes. each node will process a portion of the iterations.

2) openmp (shared memory) with in each node, use openmp to parallelize the loop iterations across multiple cores.

```
#include <mpi.h>
#include <omp.h>

int main (int argc, char **argv)
{
    int start (1, argc, &argv);
    int rank, size;
    MPI - comm - rank (MPI - comm -
    world (0) rank);
    MPI - comm - size (MPI - comm -
    world (0) size);
    int num - iterations = 100;
    int iterations - per - node = num -
    iterations / size;
    int start = rank * iterations - per -
    node;
    int end = (rank + 1) * iterations - per -
    node;
    #pragma omp parallel for
    for (int i = start; i < end; i++)
        // perform computation
    }
    MPI - finalize ();
    return 0;
}
```

Ans

3) Define and compare speedup and efficiency in MIMD systems. Give state with suitable numerical examples.

→ Speedup and efficiency in MIMD

System

Speedup
Definition: Speedup is the ratio of the execution time of a sequential algorithm to the execution time of parallel algorithm.
Formula: $S = \frac{T_{\text{sequential}}}{T_{\text{parallel}}}$

Efficiency
Definition: Efficiency measure the effectiveness of parallelization by comparing the speedup to the no. of processors used.

Formula: $E = \frac{S}{P}$ where P is the no. of processors

Speedup vs efficiency: Speedup focuses on the performance gain while efficiency focuses on the effective utilization of resources.

Numerical Example:

→ Suppose a sequential algorithm

Ans

tasks 100 seconds to execute. A parallel version of the algorithm using 4 processors takes 25 seconds.

1. Speed up calculation.

$$S = T_{\text{sequential}} / T_{\text{parallel}} = 100 / 25 = 4.$$

2. Efficiency calculation.

$$E = S / P = 4 / 4 = 1 \text{ (or } 100\%)$$

In this example, the parallel algorithm achieves a speed up of 4 an efficiency of 100% with 4 processors.