**NAME: - NIKHITA TAKSANDE**

**ROLL NO: - 3**

**BRANCH: - CSE**

**SEMESTER: - 2nd Mtech**

**SUBJECT: - HPCA**

**Q.1 The main memory of a computer is organized as 64 blocks, with a block size of 8 words. When cache has 8 blocks frames. Show the mapping from main memory to cache memory according to the following techniques with calculation of number of bits for address field.**

**i) Direct mapping.**

**ii) 2 way set associative**

**iii) Fully associative**

**Ans: i) Direct mapping**

Since there are 8 cache blocks

We need 3 bits for the cache index

Because $2^3=8$

Cache index=3 bits

Tag field =remaining bits of the address

Tag field=5

## ii) 2 way set associative

In total we need 4 bits for the cache index

Cache index=4 bits

Tag field=4

## iii) Fully associative

Since there are 64 blocks in main memeiry
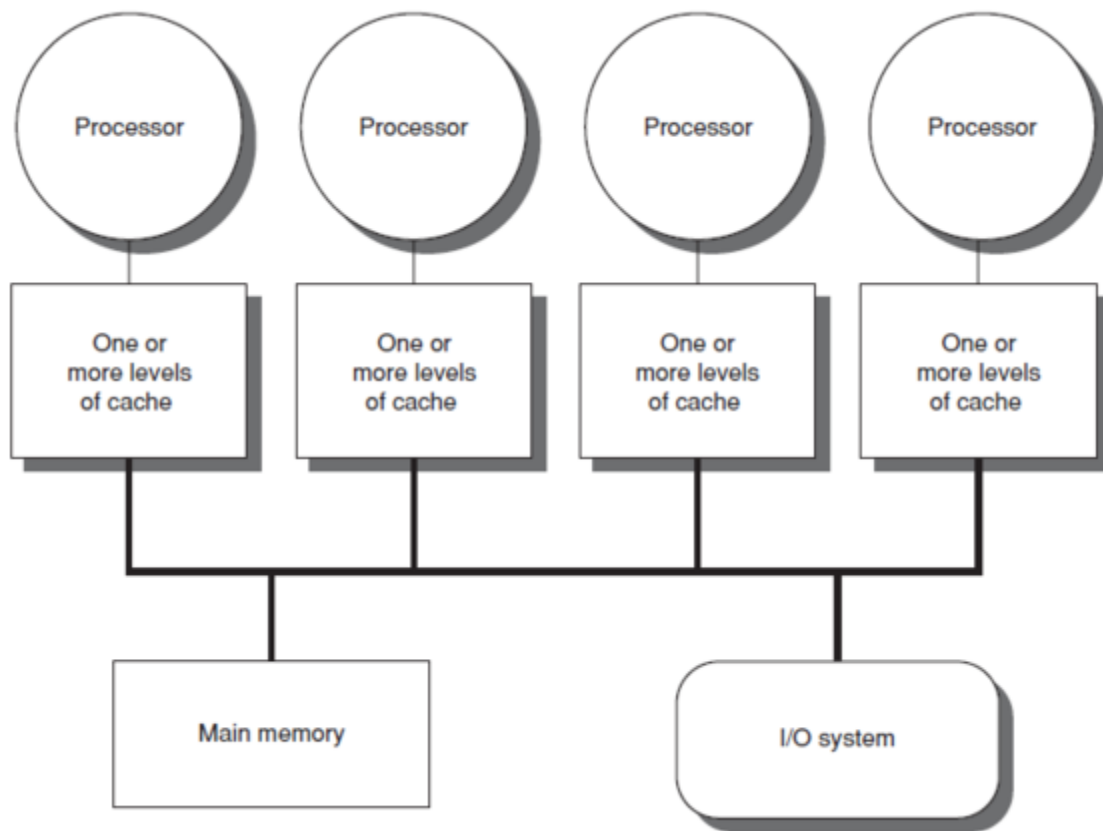
We need 6 bits for the tag field

Since 2^6=64

Tag field=6

## Q.2 Draw and explain Hierarchical cache architecture?

**Ans:** Cache hierarchy, or multi-level cache, is a memory architecture that uses a hierarchy of memory stores based on varying access speeds to cache data. Highly requested data is cached in high-speed access memory stores, allowing swifter access by central processing unit (CPU) cores.

Cache hierarchy is a form and part of memory hierarchy and can be considered a form of tiered storage.[1] This design was intended to allow CPU cores to process faster despite the memory latency of main memory access. Accessing main memory can act as a bottleneck for CPU core performance as the CPU waits for data, while making all of main memory high-speed may be prohibitively expensive. High-speed caches are a compromise allowing high-speed access to the data most-used by the

CPU,



Accessing main memory for each instruction execution may result in slow processing, with the clock speed depending on the time required to find and fetch the data. In order to hide this memory latency from the processor, data caching is used. Whenever the data is required by the processor, it is fetched from the main memory and stored in the smaller memory structure called a cache. If there is any further need of that data, the cache is searched first before going to the main memory. This structure resides closer to the processor in terms of the time taken to search and fetch data with respect to the main memory. The advantages of using cache can be proven by calculating the average access time (AAT) for the memory hierarchy with and without the cache

**Disadvantages**

- Cache memory comes at an increased marginal cost than main memory and thus can increase the cost of the overall system.[19]

- Cached data is stored only so long as power is provided to the cache.

- Increased on-chip area required for memory system.[20]

- Benefits may be minimized or eliminated in the case of a large programs with poor temporal locality, which frequently access the main memory.

# Q.3 Compare and Contrast between the cache coherence and synchronization mechanism.

# Ans:

Purpose:

- Cache coherence is about ensuring data consistency across different cache levels in a multiprocessor system.

- Synchronization mechanisms are about ensuring correct execution order and preventing concurrent access issues at the software level.

Level of Operation:

- Cache coherence is managed by hardware and deals with the cache and memory subsystem.

- Synchronization mechanisms are managed by software (though supported by hardware primitives) and deal with thread and process management.

Complexity:

- Cache coherence protocols can be complex and are implemented at the hardware level to minimize the performance impact.

- Synchronization mechanisms can be simpler in terms of implementation but can add significant overhead if not used correctly.

Interaction:

- These two mechanisms interact in the sense that cache coherence ensures that when one thread updates a variable, other threads see the update (given proper synchronization).

- Proper synchronization ensures that updates made by one thread are correctly propagated and seen by others, leveraging the underlying cache coherence protocols to maintain data consistency.

Scenarios:

- Cache coherence is critical in systems with multiple processors or cores, where each has its own cache.

- Synchronization mechanisms are critical in any concurrent programming scenario, whether on single-core or multi-core systems.

## Q.4 Consider following pipeline reservation table.

| X | shaded | shaded | X |
|---|--------|--------|---|
|   | X      |        |   |
|   |        | X      |   |

**a) What are the forbidden latency?**

**b) Draw state transition diagram?**

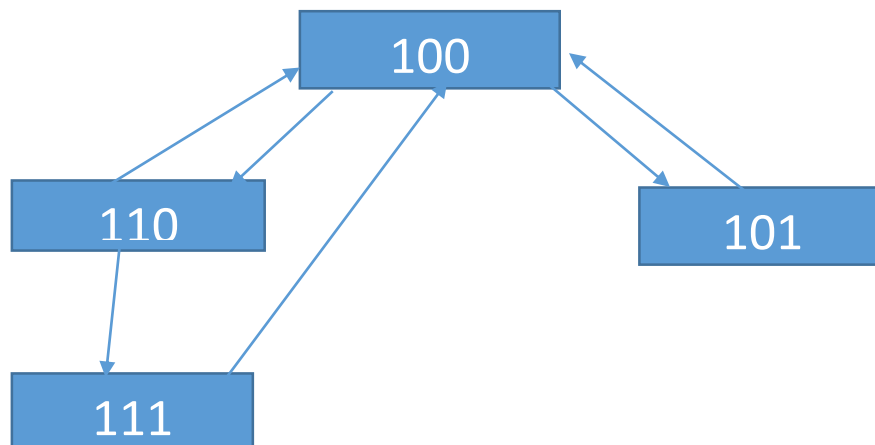**c) List all simple cycles & greedy cycles?**

**d) Determine optimal constant latency cycle and minimal average latency (MAL)**

**Ans:**

a) What are the forbidden latency?

Forbidden latency is 3

b) Draw state transition diagram?



c) List all simple cycles & greedy cycles

| Simple cycle | Average latency |
|---|---|
| 2 | 2 |
| 1,4 | 2.5 |
| 1,1,4 | 2 |
| 4 | 4 |
| 2,4 | 3 |

d) Determine optimal constant latency cycle and minimal average latency (MAL)

Ans:

The optimal latency having average latency is lower (2),(1,4)(1,1,4)

The minimal value in simple cycle will be 2

## Q.5 Consider the following reservation table.

| X |   |   |   |   | X |   |
|---|---|---|---|---|---|---|
|   |   | X |   |   |   | X |
|   | X |   | X |   |   |   |
|   |   | X |   | X |   |   |

**a) Determine latencies in the forbidden list F and collision vector, C.**

**b) Draw the state transition diagram.**

**c) List all simple cycles and greedy cycles.**

**d) Determine MAL.**

**e) If clock period T = 20ns, determine the maximum throughput of the pipeline**

**Ans:**

a) Determine latencies in the forbidden
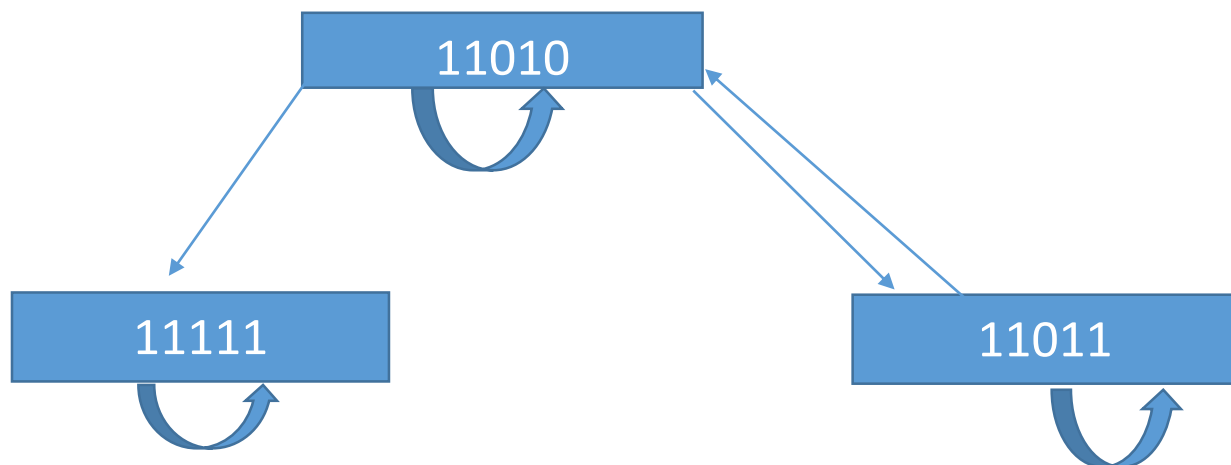
 List F and collision vector

Forbidden latency is 2,4,5

Collision vector

5 4 3 2 1

1 1 0 1 0


b) Draw the state transition diagram.



c) List all simple cycles and greedy cycles.

| Simple cycle | Average latency |
| --- | --- |
| 1,6 | 3.5 |
| 3,6 | 4.5 |
| 3 | 3 |
| 6 | 6 |

d) Determine MAL.

Minimum average latency is 3

e) If clock period T = 20ns, determine the maximum throughput of the pipeline

maximum average throughput for T=20 ns

1/MAL=1/MAL*T=1/20ns

=1/60*10^-9

=10^9/60sec

=10^3/60  *10^6/sec

=10000/60

166.66 MPS