

COA-Module 4 Notes

Q1) What are the characteristics of Memory?

A:

(LACTPOPP)

1. Location:

CPU: Memory within the central processing unit.

Internal: Main memory or cache, directly accessible by the CPU.

External: Storage outside the main system, such as disks and tapes, used for backup and additional storage.

2. Capacity:

Word Size: The natural unit of data organization (e.g., 8-bit, 16-bit).

Number of Words or Bytes: Total number of words or bytes stored in the memory.

3. Unit of transfer:

Internal: Governed by data bus width, typically a word.

External: Generally a larger block, such as a disk sector.

Addressable Unit: The smallest unit uniquely addressable in the memory system, often a byte.

4. Access method:

Sequential: Data is accessed in a set sequence, starting from the beginning (e.g., tape storage). Access time depends on the location of data and the previous location

Direct: Individual blocks have unique addresses, and data access can be achieved by jumping to a specific block, followed by a sequential search (e.g., disk storage). Access is by jumping to vicinity plus sequential search. Access time depends on location and previous location

Random: Any location can be accessed independently and immediately (e.g., RAM). Access time is independent of location or previous access.

Associative: Data is accessed by comparing contents rather than by a specific address (e.g., cache memory). Access time is independent of location or previous access

5. Performance(SRAM,DRAM):

Access Time: Time between a request and data availability at the required location.

Memory Cycle Time: Minimum time between successive read requests.

Transfer Rate: Speed at which data can be transferred to/from memory.

6. Physical type:

Semiconductor: RAM types, such as SRAM and DRAM.

Magnetic: Disk drives and magnetic tapes.

Optical: Media such as CDs and DVDs.

7. Physical characteristics:

Decay: Rate at which stored charge decays, affecting data stability.

Volatility: Determines whether memory retains data without power (e.g., RAM is volatile, ROM is non-volatile).

Erasable: Some types can be erased and rewritten (e.g., Flash memory).

Power Consumption: Amount of power needed to maintain data, particularly in volatile memory.

8. Organisation- Direct Mapping, Associative Mapping:

Direct Mapping: A specific block in main memory maps directly to a specific line in the cache.

Associative Mapping: A memory block can be loaded into any line of cache; often used in fully associative or set-associative caches for flexibility.

Q2)Explain Memory Hierarchy

A:

The memory hierarchy is designed to optimize the trade-offs between speed, cost, and storage capacity by layering various types of memory. The hierarchy includes the following main levels:

1. Registers (In CPU)

- **Registers** are the smallest and fastest memory units located directly within the CPU.
- They hold data temporarily while the CPU performs operations, like arithmetic or data manipulation.
- Since they are directly accessible by the CPU, registers provide very high-speed access compared to other memory types.
- Common examples of registers include the **Accumulator** (for storing results of operations), **Program Counter** (for keeping track of instruction locations), and **General Purpose Registers**.

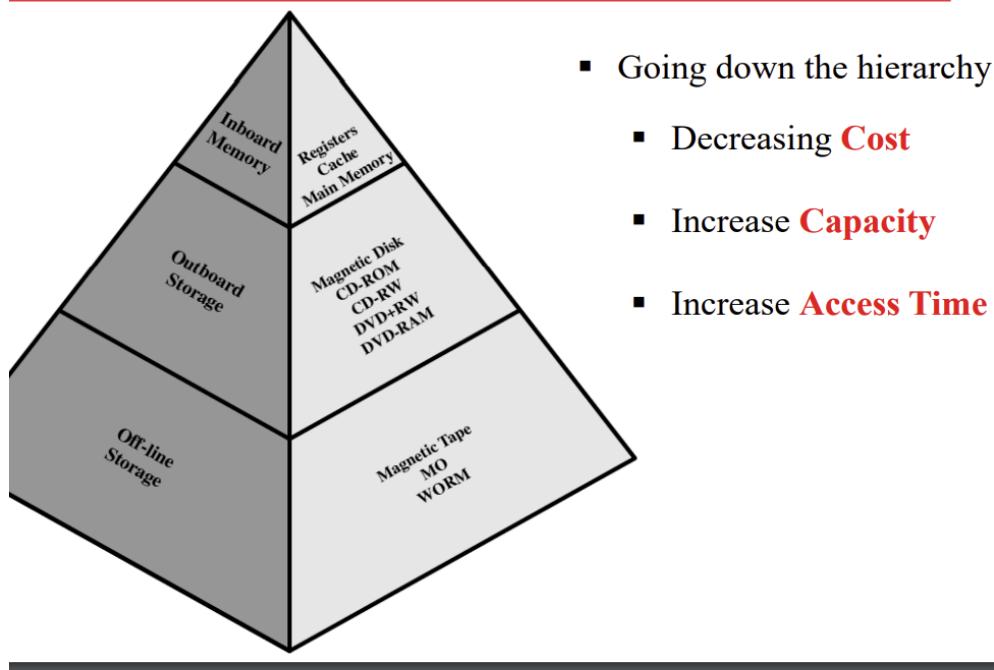
2. Internal or Main Memory

- This is the primary storage area that the CPU uses to store data and instructions during program execution.
- **Main Memory** typically refers to **RAM (Random Access Memory)**, which is fast but volatile (data is lost when the system powers off).
- Internal memory may include **one or more levels of cache**:
 - **Cache** is a small, very fast type of memory that sits between the CPU and main memory. It's used to store frequently accessed data, reducing the time needed for the CPU to retrieve information.
 - Cache levels are usually labeled as **L1, L2, and L3**—with L1 being the smallest and fastest, and L3 being larger but slower.
- **RAM** provides the CPU with quick access to data and instructions necessary for current processes, and it is larger than cache but slower than registers.

3. External Memory (Backing Store)

- **External memory**, also called **secondary storage** or **backing store**, is used for long-term storage of data and programs.
- This type of memory is **non-volatile**, meaning data is retained even when the system is powered off.
- Examples include **hard disk drives (HDDs)**, **solid-state drives (SSDs)**, **USB drives**, and **optical discs** (e.g., CDs and DVDs).
- External memory is much slower than internal memory (RAM or cache) but offers much larger storage capacity. It's used to store data that isn't immediately needed by the CPU, like files, applications, and the operating system.

Memory Hierarchy - Diagram



Q3)What is RAM?

A:

RAM (Random Access Memory) is a type of volatile memory that allows data to be read from or written to any location directly (randomly), without going through other memory locations sequentially. It's the primary memory in a

computer system, playing a critical role in the performance of applications and processes.

- **Random Access:** Allows data to be accessed directly from any location without sequential order.
- **Read/Write Memory:** Supports both reading and writing, enabling temporary storage of data and instructions for active processes.
- **Volatile:** Requires power to retain data; all information is lost when the system shuts down.
- **Temporary Storage:** Acts as a workspace for the CPU, holding data and applications in use for quick access.
- **Types:**
 - **Static RAM (SRAM):** Faster, used in cache memory.
 - **Dynamic RAM (DRAM):** Slower, forms the main system memory.

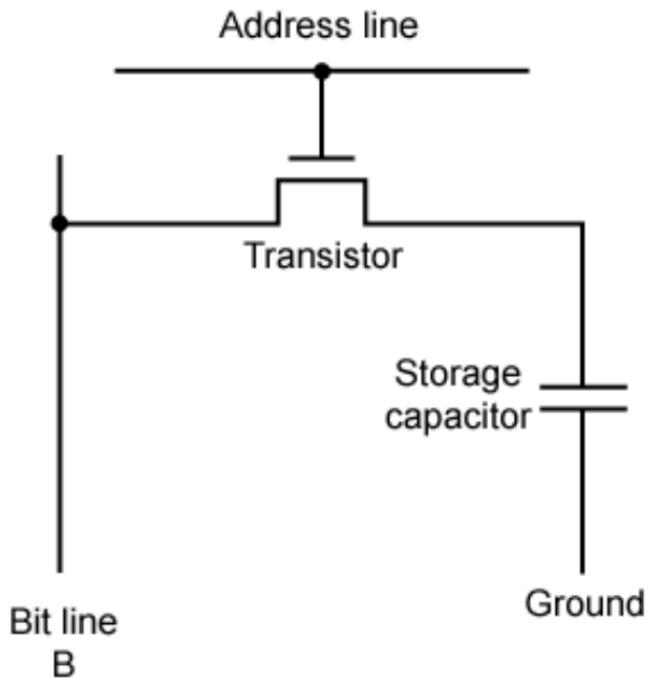
Q4)What is Dynamic RAM?

A:

- **Bits Stored as Charge in Capacitors:** DRAM stores each bit as a small charge in a capacitor.
- **Charge Leakage:** The charge in capacitors leaks over time, which means data cannot be stored permanently.
- **Need for Refreshing:** Even when powered, DRAM requires **periodic refreshing** to maintain data integrity due to charge leakage.
- **Simpler Construction:** DRAM cells are simpler than SRAM cells, making it **less expensive** to manufacture.
- **Refresh Circuits Required:** Additional **refresh circuits** are necessary to recharge the capacitors regularly, contributing to slower performance.
- **Speed:** DRAM is generally slower than SRAM because of the time required for refreshing and recharging.
- **Main Memory:** DRAM is widely used as **main memory** in computers due to its high density and cost-effectiveness.

- **Analog Nature:** DRAM is essentially analog; the **level of charge** in each capacitor determines the value (0 or 1) of the stored bit.

Dynamic RAM Structure



Q5) What is Static RAM?

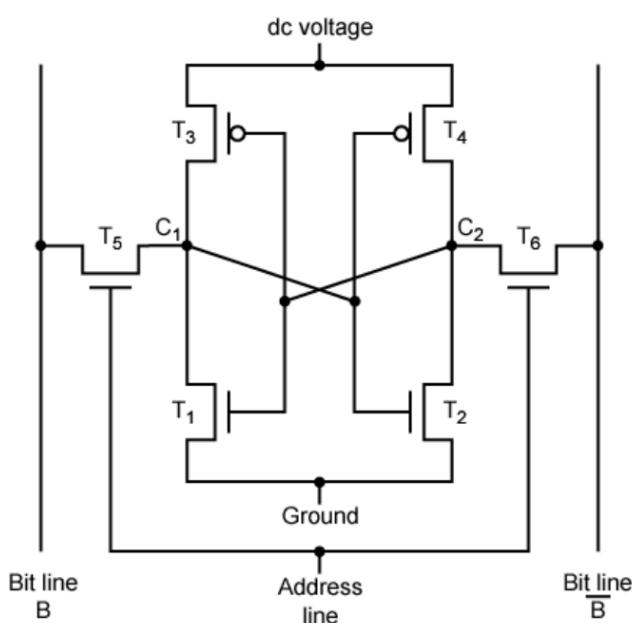
A:

- **Bits Stored as On/Off Switches:** SRAM stores each bit using a set of transistors configured as a flip-flop, which holds the data in a stable state without needing refreshing.
- **No Charge Leakage:** Unlike DRAM, SRAM does not rely on capacitors, so there is no charge to leak, and **data remains stable** as long as power is supplied.
- **No Refreshing Required:** Since SRAM does not lose charge, it does not need refreshing, which allows for **faster access times** compared to DRAM.

- **More Complex Construction:** SRAM cells use more transistors per bit (typically 4-6), making it more **complex and physically larger** than DRAM.
- **Higher Cost:** Due to its complex construction, SRAM is **more expensive** than DRAM per bit.
- **Faster Performance:** SRAM's lack of refresh cycles makes it significantly faster, which is why it is often used for **cache memory** close to the CPU.
- **Digital Nature:** Unlike the analog nature of DRAM, SRAM operates digitally, with each bit represented by a stable, discrete on/off state (flip-flop).

In summary, **SRAM** is a fast, stable type of memory often used for cache due to its **non-volatile behavior during power, higher speed, and reliability**, though it comes at a higher cost and with less density than DRAM.

Static RAM Structure



Q7) Differentiate between SRAM and DRAM.

A:

Feature	SRAM (Static RAM)	DRAM (Dynamic RAM)
Volatile	Yes	Yes
Power needed to preserve data	Requires constant power to retain data	Requires constant power to retain data
Dynamic cell	No, uses flip-flops for storing data	Yes, uses capacitors to store data
Simpler to build	More complex, requires more transistors	Simpler, uses fewer components
Smaller	Larger due to more complex structure	Smaller, more compact
Density	Less dense (fewer bits per unit area)	More dense (more bits per unit area)
Cost	More expensive	Less expensive
Needs refresh	No refresh needed	Requires periodic refreshing of data
Larger memory units	Typically used for smaller, faster caches	Used for larger memory systems (e.g., RAM)
Speed	Faster (due to lack of refresh cycles)	Slower (due to refresh cycles and simpler design)
Used in	Primarily used in cache memory	Used in main memory (system RAM)

Q8) What is ROM?

A:

ROM (Read-Only Memory) is a type of non-volatile memory that is primarily used for storing firmware or software that doesn't change frequently. Unlike RAM, which is volatile (loses data when power is turned off), ROM retains its contents even when power is removed. ROM typically stores the boot-up instructions and firmware for hardware devices, such as the computer's BIOS, embedded systems, or other permanent software that doesn't need to be modified often.

Characteristics of ROM:

- **Non-volatile:** Data is retained even when the device is powered off.
- **Read-only:** The data stored in ROM can be read, but typically cannot be modified (depending on the type).
- **Used for firmware:** Typically stores low-level system instructions, like the BIOS in a computer or firmware in devices like printers and routers.
- **Durable and reliable:** Since the data is etched permanently (in the case of traditional ROM), it is highly resistant to data corruption.

Q9)What is PROM?

A:

PROM (Programmable Read-Only Memory) is a type of non-volatile memory that is programmed during the manufacturing process, allowing it to retain data even when the power is turned off. PROM is often used to store firmware or other software that doesn't need to change once it's been written.

- **Written during manufacture:** PROM is programmed during its manufacture and cannot be altered afterward.
- **Programmable (“once”):** Data is written once and cannot be modified after programming.
- **Small amount of data to be written:** Typically used for storing small, critical data like firmware.
- **Less expensive:** PROM is cost-effective because it only needs a one-time programming process.
- **Non-volatile, written only once:** Retains data even without power and cannot be rewritten once programmed.
- **Writing performed electrically at the time of chip fabrication:** Data is written to the chip using electrical signals during manufacturing, with no possibility of modification afterward.

Q10) What is EPROM?

A:

- **Erasable Programmable (EPROM) – Erased by UV:** EPROM can be erased by exposing it to ultraviolet (UV) light, which resets the stored data.
- **Read and written electrically:** Data can be both read and written to an EPROM chip using electrical signals.
- **All storage cells should be erased electrically to initial state by exposure to UV radiation:** To erase EPROM, the chip must be exposed to UV light, which clears all the data, allowing it to be reprogrammed.
- **Can be altered multiple times and holds data virtually indefinitely:** EPROM can be rewritten many times and the data remains intact for a long period if not exposed to UV light.
- **More expensive than PROM:** EPROM is more expensive than PROM due to the additional equipment required for erasure and reprogramming.

Q11) What is EEPROM?

A: **Electrically Erasable Programmable Read-Only Memory (EEPROM):**

- **Can be written anytime without erasing prior contents:** EEPROM allows individual bytes to be rewritten without affecting the other stored data.
- **Write operation takes longer than read:** Writing data to EEPROM takes more time than reading it, due to the more complex process involved.
- **More expensive than EPROM, less dense:** EEPROM is more costly than EPROM because it offers more flexibility but at the expense of lower storage density.

Q12)What are the types of ROM?

A:

- 1. Programmable Read-Only Memory (PROM):**
 - **Empty of data when manufactured:** PROM comes without any data preloaded and can be programmed by the user once.
 - **May be permanently programmed by the user:** After programming, the data is permanently written and cannot be changed.
- 2. Erasable Programmable Read-Only Memory (EPROM):**
 - **Can be programmed, erased, and reprogrammed:** EPROM can be written to, erased with UV light, and reprogrammed multiple times.
 - **The EPROM chip has a small window on top allowing it to be erased by shining ultra-violet light on it:** A transparent window allows UV light to erase the stored data for reprogramming.
 - **After reprogramming, the window is covered to prevent new contents from being erased:** Once reprogrammed, the window is covered to prevent accidental erasure.
 - **Access time is around 45–90 nanoseconds:** EPROM chips have relatively fast access speeds, with typical read times between 45 and 90 nanoseconds.
- 3. Electrically Erasable Programmable Read-Only Memory (EEPROM):**
 - **Reprogrammed electrically without using ultraviolet light:** EEPROM can be erased and rewritten electrically, making it more convenient than EPROM.
 - **Must be removed from the computer and placed in a special machine to do this:** Though EEPROM is electrically erasable, it still needs to be physically removed and placed in a programmer to update data.

- **Access times between 45 and 200 nanoseconds:** EEPROM has slightly slower access times compared to EPROM, with read times ranging from 45 to 200 nanoseconds.

4. Flash ROM:

- **Similar to EEPROM:** Flash ROM shares many similarities with EEPROM but has more advanced features for efficient reprogramming.
- **However, can be reprogrammed while still in the computer:** Unlike EEPROM, Flash ROM can be updated without removing it from the device, making it more convenient for regular updates.
- **Easier to upgrade programs stored in Flash ROM:** Flash ROM is commonly used in devices that require frequent updates, such as firmware upgrades in embedded systems.
- **Used to store programs in devices e.g., modems:** Flash ROM is used in devices like modems, routers, and smartphones to store the system's firmware.
- **Access time is around 45–90 nanoseconds:** Flash ROM has fast access times, similar to EPROM, typically between 45 and 90 nanoseconds.

5. ROM Cartridges:

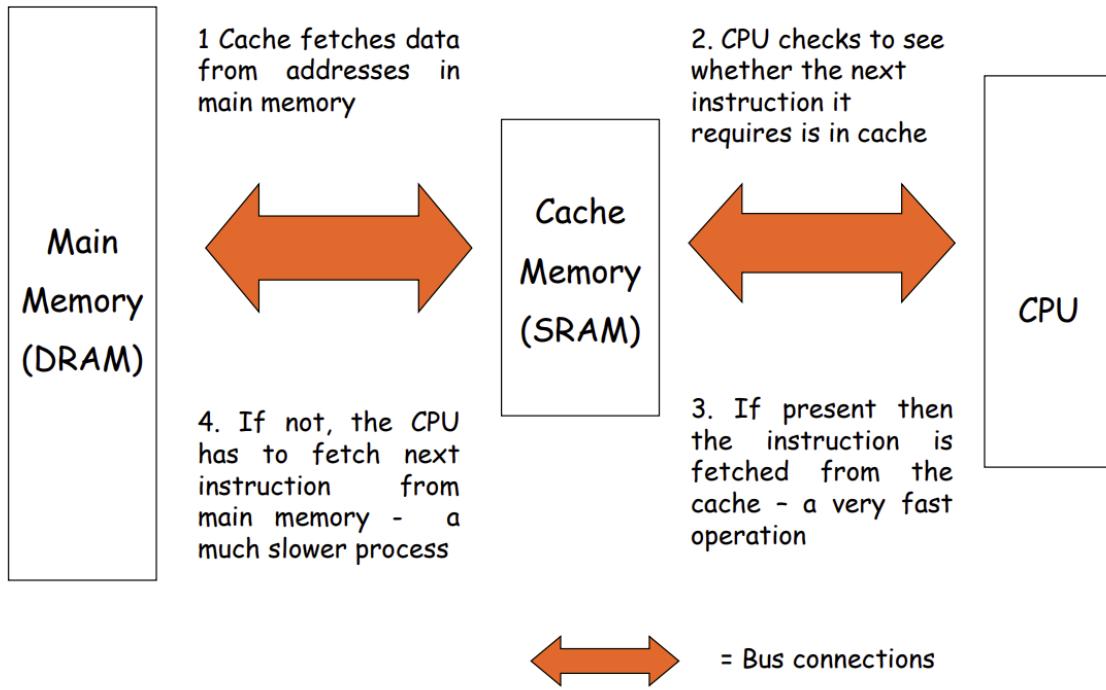
- **Commonly used in game machines:** ROM cartridges are used to store video games and other software in consoles, providing a portable and easily replaceable storage medium.
- **Prevents software from being easily copied:** ROM cartridges offer some level of protection against software piracy because the content is usually read-only and not easily copied or modified.

Q13)What is Cache? Show the operation of cache memory.

A:

Cache is a small amount of fast memory- which sits between the main memory and CPU. It may be located on a CPU chip or module. It is used to store data required in regular cases.

The operation of cache memory

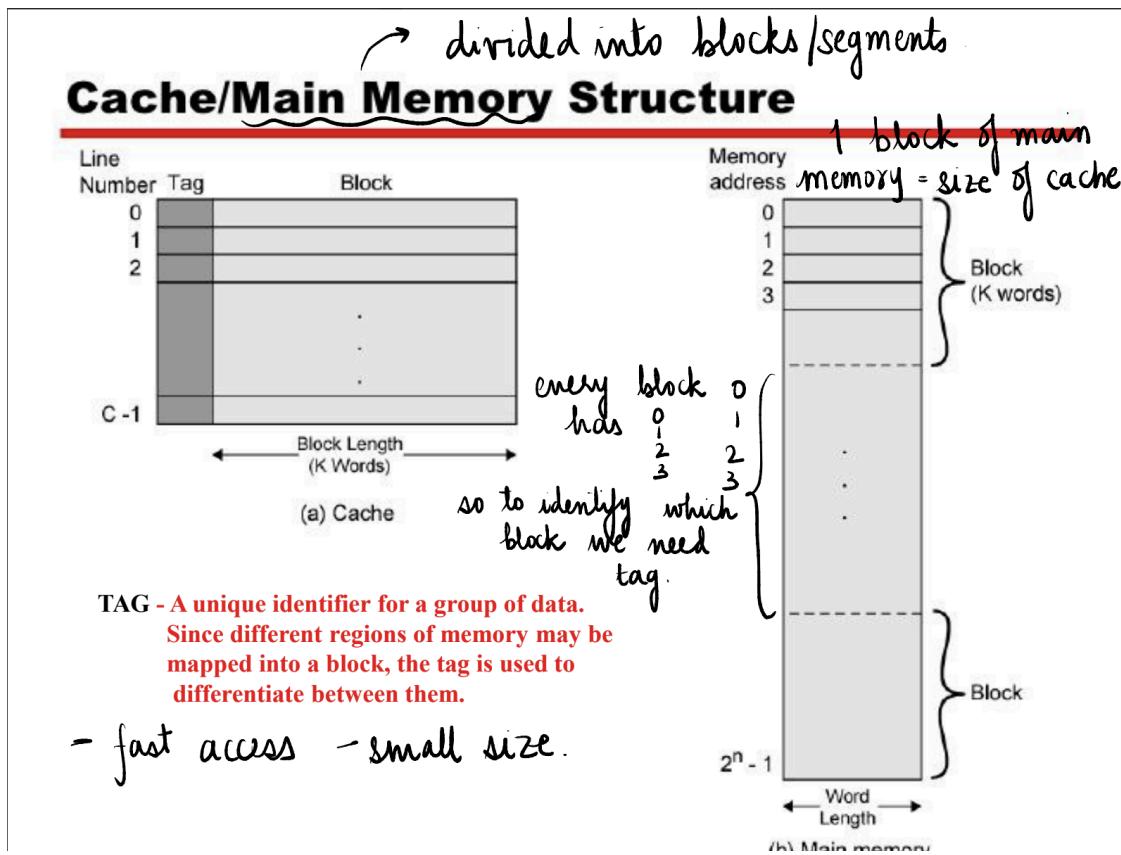


Q14)Give an overview of cache operation.

A:

- CPU requests contents of memory location.
- Check cache for this data.
- If present, get from cache (fast).
- If not present, read required block from main memory to cache.

- Then deliver from cache to CPU.
- Cache includes tags to identify which block of main memory is in each cache slot.



Q15) What is Tag?

A:

TAG - A unique identifier for a group of data. Since different regions of memory may be mapped into a block, the tag is used to differentiate between them.

Q16) Explain the differences between L1, L2, and L3 cache in terms of speed, size, and location.

A:

- **L1 Cache:** This is the smallest and fastest cache, located directly on the CPU chip. It has a very small capacity (usually 16KB to 128KB) and is used for quick access to frequently used data and instructions.
- **L2 Cache:** Larger than L1, with more capacity (typically 128KB to several MB), and located either on the CPU chip or on a separate chip near the CPU. It is slower than L1 but still faster than RAM and helps improve CPU performance by storing data that L1 might not have.
- **L3 Cache:** This is the largest and slowest among the three, often shared by multiple CPU cores. It can be several MB to tens of MB in size. It helps improve the performance of both L1 and L2 by storing a larger pool of data, with access speeds roughly twice that of RAM.

Q17) Explain the key factors involved in cache design.

A:

- **Size:** The size of the cache determines how much data can be stored, affecting performance and cost; larger caches generally improve speed but increase cost.
- **Mapping Function:** The mapping function defines how data from main memory is assigned to cache locations, impacting cache efficiency and access time.
- **Replacement Algorithm:** This algorithm decides which data to evict when the cache is full, with common strategies like LRU (Least Recently Used) or FIFO (First In First Out).
- **Write Policy:** Determines how data is written to cache and main memory, with strategies like write-through (immediate write to memory) and write-back (write only when data is evicted).
- **Block Size:** Refers to the amount of data fetched from memory in a single cache line, affecting the trade-off between data locality and cache efficiency.

- **Number of Caches:** This refers to how many levels of cache exist (L1, L2, L3) and their respective sizes, impacting overall system performance.
- **Cost:** More cache increases the overall cost of the system due to the need for more physical memory and more complex cache management.
- **Speed:** Larger caches can improve system speed by reducing the need to access slower main memory, but there's a point where adding more cache yields diminishing returns.
- **Checking Cache for Data Takes Time:** While caches improve speed, checking multiple levels of cache for data introduces overhead, particularly in systems with deep cache hierarchies.

Q18)What does the Mapping Technique mean?

A: Mapping Technique means how data from main memory is stored in cache.

Q19)What is: i)Tag ii)Word iii)Line

A:

1. **Tag Field:** The tag field stores the higher-order bits of the memory address and is used to check if the requested data is in the cache.
2. **Line Field:** The line field selects which cache line the data might be stored in by using the middle portion of the memory address.
3. **Word Field:** The word field selects the specific word within the cache line when multiple words are stored in each line.

Q20)What is direct mapping?

A:

Direct Mapping is a cache mapping technique where each block of main memory is mapped to exactly one cache line. In other words, for a given memory address, there is only one specific location in the cache where the data can be stored, determined by a simple mapping function.

Key Points:

- Each block of memory is assigned to a specific cache line based on the address.
- The memory address is divided into three parts: **Tag**, **Line (or Index)**, and **Block Offset**.
- The **Line field** determines which cache line to look in, and the **Tag field** is compared with the stored tag in that cache line to check for a hit or miss.

Advantages of Direct-mapping

- It requires very less hardware complexity.
- It is a very simple method of implementation.
- Access time is very low.

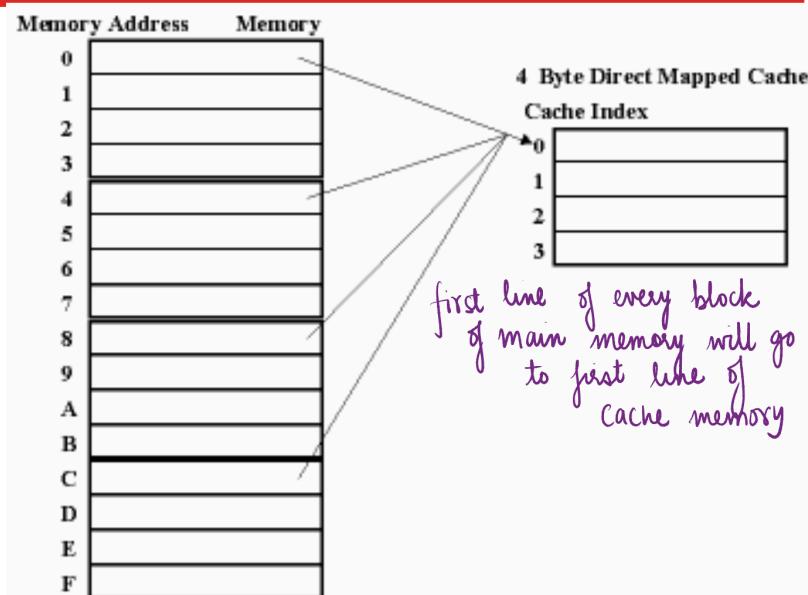
Disadvantages of Direct-mapping

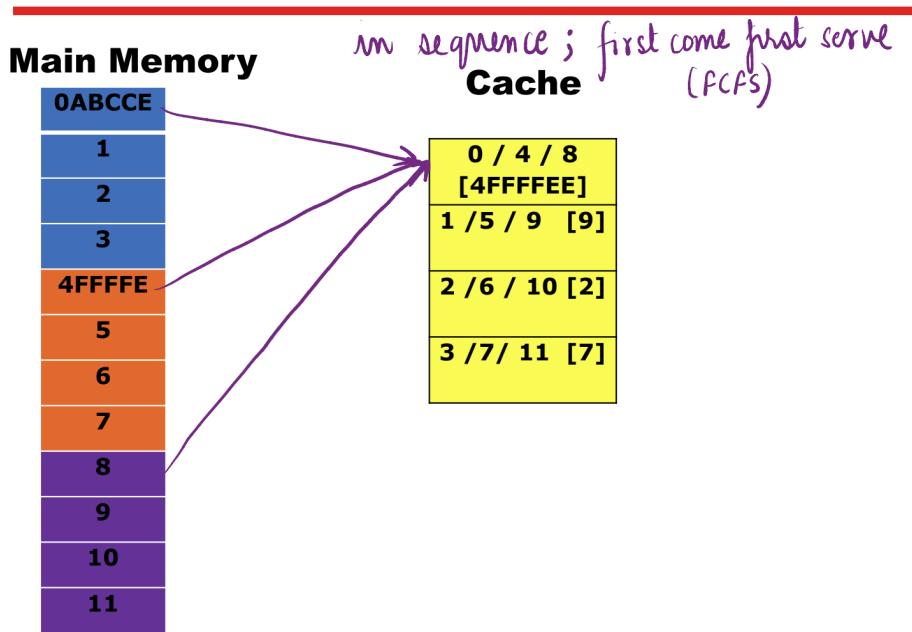
- Cache performance is unpredictable in direct mapping.
- Handling of spatial locality is poor.

- Use of cache space is inefficient.
- Conflict misses are high.

Simple • Inexpensive • Fixed location for given block —If a program accesses 2 blocks that map to the same line repeatedly, cache misses are very high

DIRECT MAPPING CONCEPT





Q21) What is Fully associative mapping?

A:

Associative Mapping (also known as **Fully Associative Mapping**) is a cache mapping technique where any block of memory can be stored in any cache line, allowing the cache to store data more flexibly.

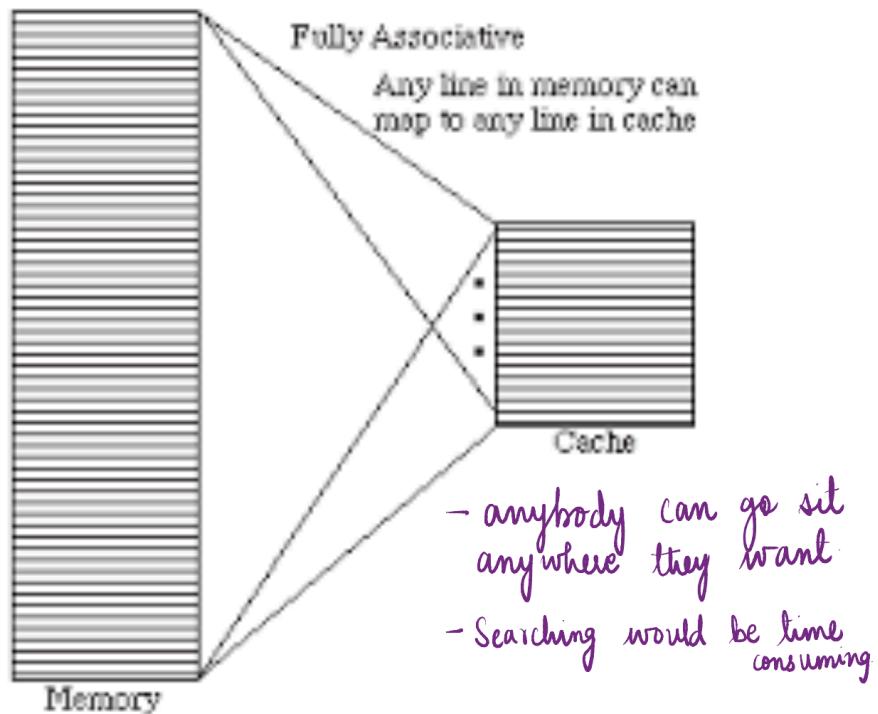
Key Points:

- In associative mapping, the **Tag** field holds the entire address (without any predefined index or line field), and every cache line can store any memory block.

- There is no fixed relationship between a memory block and a cache line, so the cache controller must search through all the cache lines to check if the required data is present (cache lookup).
- **Tag field** is compared with all cache lines to determine a hit or miss.

A main memory block can load into any line of cache • Memory address is interpreted as tag and word • Tag uniquely identifies block of memory • Every line's tag is examined for a match • Cache searching gets expensive.

FULLY ASSOCIATIVE MAPPING



Q22)What is Set associative mapping?

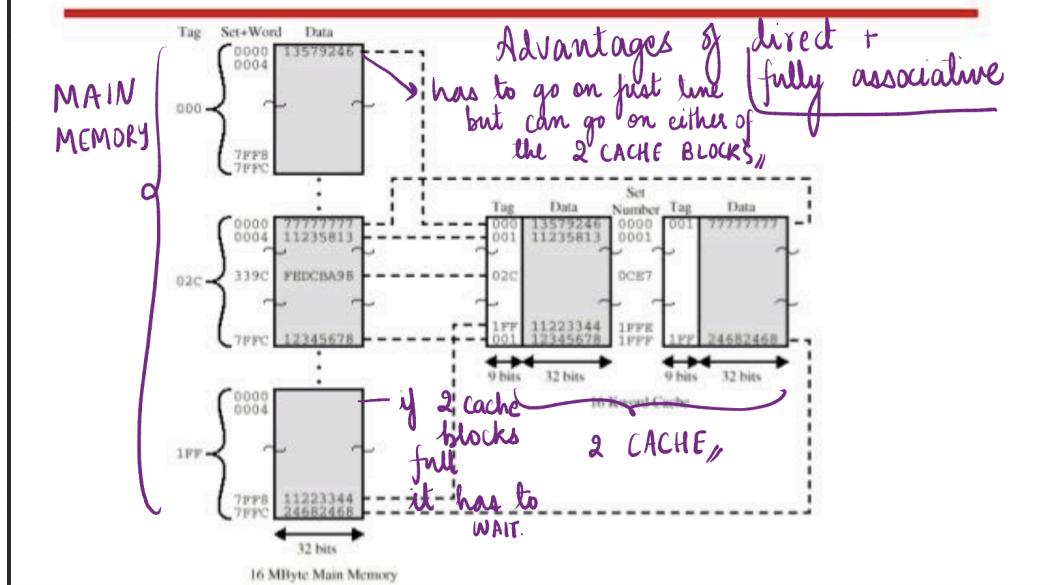
A:

Set-Associative Mapping is a cache mapping technique that combines aspects of both **direct mapping** and **associative mapping**. In this method, the cache is divided into multiple sets, and each set contains a fixed number of lines (or slots). A given block of memory can be stored in any of the lines within a specific set.

Key Points:

1. **Cache divided into sets:** The cache is organized into multiple sets, each containing a small number of cache lines.
2. **Each set contains a number of lines:** Each set can store a predefined number of cache lines (e.g., 2, 4, etc.).
3. **Block maps to any line in a set:** A memory block is mapped to one of the lines in a specific set, but not to a particular line in the cache. The block can go into any available line within the set.
4. **Example (2-way set associative):** If there are 2 lines per set, it's called **2-way associative mapping**, meaning that a given memory block can be placed in either of the 2 lines of one set.

Two Way Set Associative Mapping Example



Brushing up on Associative & Set Associative Mapping:

Associative Mapping:

- A main memory block can load into any line of cache
- Memory address is interpreted as tag and word
- Tag uniquely identifies block of memory
- Every line's tag is examined for a match
- Cache searching gets expensive

Set Associative Mapping:

- Cache is divided into a number of sets
- Each set contains a number of lines
- A given block maps to any line in a given set

- —e.g. Block B can be in any line of set i
- e.g. 2 lines per set
- —2 way associative mapping
- —A given block can be in one of 2 lines in only one set

Q23)What is Cache Coherence?

A:

Cache Coherence is the concept of maintaining a consistent view of data across multiple caches in a multiprocessor system. When multiple processors have their own caches, each cache may hold copies of the same memory data, which can lead to inconsistencies if one processor updates its copy while others have outdated versions.

Key Mechanisms in Cache Coherence:

1. Write Policies:

- **Write-Through:** Updates to a cache are immediately written to main memory, ensuring all caches see the latest data but increasing memory traffic.
- **Write-Back:** Updates occur only in the cache and are written to memory later, reducing traffic but requiring additional coherence mechanisms.

2. Coherence Protocols:

- **MESI Protocol:** Defines states (Modified, Exclusive, Shared, Invalid) for cache lines to track and manage shared data consistency.
- **Snoopy Protocols:** Each cache controller monitors a shared bus for changes, updating or invalidating its data as needed to maintain consistency.

Solutions

- **Write-Through:** Each write to a cache is also written to main memory, allowing other caches to monitor memory for updates and keep data consistent. However, it increases memory traffic.
- **Write-Back:** Updates are made only in the cache and written to memory later. This reduces memory traffic but requires extra mechanisms to handle consistency across caches.
- **Write Invalidate:** When one processor writes to its cache, it invalidates that data in other caches. This ensures only one processor has a valid copy during the update, which reduces conflicts.
- **Write Update:** When data is updated in one cache, the change is sent to other caches holding that data, ensuring they all receive the latest version immediately. This can increase bus traffic but ensures consistency.

Q24)What are the Software Solutions?

A:

In some systems, **cache coherence** can be managed at the **software level** by the **compiler and operating system (OS)**, which reduces the hardware burden for maintaining consistency across caches. Here's how this approach works:

- **Compile-Time Overhead:** The compiler is responsible for identifying which data is likely to change during program execution. This **overhead is handled at compile time**, reducing the need for hardware-based coherence mechanisms.
- **Data Marking:** The compiler marks data that may be frequently modified, flagging it for the OS. By doing so, the system knows which data should not be cached or should be handled carefully to avoid coherence issues.
- **Operating System Control:** The OS manages access to marked data by **preventing it from being cached** or by handling it in a way that reduces coherence conflicts. For example, the OS may keep frequently updated data in main memory rather than cache.
- **Design Complexity Shift:** By handling coherence in software, **design complexity is shifted from hardware to software**. This approach can simplify hardware design and reduce power consumption, though it

may introduce additional overhead in the software development process.

Q25)What are the Hardware Solutions?

A:

In hardware-based cache coherence solutions, **cache coherence protocols** and mechanisms are implemented directly in the hardware, allowing dynamic management of data consistency across multiple caches without programmer intervention. Here's how it works:

- **Cache Coherence Protocols:** Protocols like **MESI (Modified, Exclusive, Shared, Invalid)** ensure that each cache knows the state of data it holds, allowing caches to track changes and manage consistency automatically.
- **Dynamic Problem Recognition:** The hardware can **recognize potential coherence issues at runtime**, meaning it detects and resolves inconsistencies as they arise, rather than relying on pre-compiled instructions or OS management.
- **Efficient Cache Use:** Hardware coherence protocols allow caches to **operate more efficiently** by coordinating data updates across caches, minimizing redundant data fetching and reducing cache misses.
- **Transparency to Programmer:** The coherence mechanisms work in the background, making **coherence management invisible to the programmer**.

programmer. This transparency reduces the complexity of writing code for multiprocessor systems.

- **Snoopy Protocols:** Snoopy protocols distribute cache coherence responsibility among cache controllers. Each cache "snoops" on the shared bus, monitoring it for transactions related to shared data. If a cache detects a change to data it holds, it updates or invalidates its copy as needed. This approach is especially effective in bus-based multiprocessor systems, where each cache can monitor bus traffic to maintain consistency.

Q26)What is Snoopy Protocol?

A:

Snoopy Protocols are hardware-based cache coherence protocols that maintain data consistency across caches in multiprocessor systems. They do this by **distributing coherence management** to each cache controller, allowing caches to monitor ("snoop on") the shared system bus for changes in shared data. Here's an elaboration:

- **Distributed Responsibility:** Instead of relying on a central controller, **each cache controller is responsible for maintaining coherence**. When one cache updates a shared data line, it broadcasts the update across the bus.

- **Shared Line Recognition:** Caches are aware of which data lines are shared across multiple caches. If a cache holds data that may be shared, it monitors the bus for updates to ensure it has the most recent data.
- **Announcement of Updates:** When one processor modifies a data line, its cache controller **announces the update on the bus**. Other caches that hold the same data line listen to this announcement and either update or invalidate their copies as required.
- **Ideal for Bus-Based Multiprocessors:** Snoopy protocols work well in **bus-based multiprocessor systems**, where a single bus connects multiple processors and caches. This shared bus structure allows all caches to "snoop" on each other's actions.
- **Increased Bus Traffic:** A trade-off of snoopy protocols is **higher bus traffic**, as each update or modification requires a broadcast across the bus. This can lead to congestion, especially in systems with many processors, potentially affecting performance. However, the protocol ensures that all caches remain consistent.

Q27) Explain Write Through.

A:

1. **All writes go to main memory as well as cache:**

- In multi-core systems, when a CPU writes data, it updates both its local cache and main memory. This approach, known as a *write-through* policy, ensures that main memory always has the latest data and that other CPUs can read accurate values if they bypass the cache. While this guarantees consistency, it can create additional overhead as each write operation involves both the cache and main memory.

2. Multiple CPUs can monitor main memory traffic to keep local (to CPU) cache up to date:

- To ensure consistency across multiple CPUs, each CPU monitors changes in main memory and in other CPUs' caches. This technique, known as *cache coherence*, allows each CPU's cache to stay updated when another CPU modifies a shared piece of data. The CPUs use a protocol (like MESI) to manage this coherence, making sure they don't work with outdated data.

3. Lots of traffic:

- The need to update main memory and synchronize caches across CPUs generates substantial data traffic on the memory bus. Each cache update, write operation, or coherence message adds to this traffic, which can slow down the system as all CPUs compete to access and update memory simultaneously.

4. Slows down writes:

- The combination of extra steps—writing to cache, updating main memory, and maintaining coherence—adds latency to each write operation. This slows down overall performance, as write operations are no longer instant and require coordination across the system. The more CPUs involved, the greater the potential slowdown due to increased coherence traffic.

Q28)Explain Write Back.

A:

Write Back- where updates are initially made only to the cache and not directly to main memory, along with some of the challenges this creates.

Here's a detailed breakdown of each point:

1. Updates initially made in cache only:

- When a CPU modifies data, the update is initially applied only to the data in the cache, rather than immediately writing it to main memory. This approach (known as "write-back caching") reduces the frequency of main memory accesses, which can improve speed since the CPU only accesses the faster cache for these updates.

2. Update bit for cache slot is set when update occurs:

- Each cache slot has an *update bit* (often called a "dirty bit") that is set whenever data in that cache slot is modified. This bit marks the cache line as "dirty," indicating that it has been modified and now contains data different from what's stored in main memory.

3. If block is to be replaced, write to main memory only if update bit is set:

- When the cache needs to replace an old data block with a new one, it checks the update bit for the old block. If the bit is set, it means the block contains data not yet written to main memory, so it writes this data back to main memory before replacing it. If the update bit is not set, the cache can simply discard the old data without writing to main memory, saving time.

4. Other caches get out of sync:

- Since the update is made only in one CPU's cache, other CPUs' caches may have outdated data for the same memory address. Without coordination, these caches will be "out of sync," meaning they don't have the latest data. Cache coherence protocols (like MESI) are used to handle this issue, ensuring that all CPUs eventually receive updates.

5. I/O must access main memory through cache:

- When I/O devices need to access data that the CPU has modified in the cache, they must go through the cache (or be updated

with cache data) to get the latest information. This coordination ensures that I/O devices work with the correct data, but it adds complexity, as I/O must be aware of cache contents to avoid outdated information.

Q29)Explain Write Update

A:

Write Update:

- When one processor updates a memory location, this change is sent to all other processors' caches. This ensures that all caches have the latest data. It helps maintain consistency but increases memory traffic.

Multiple readers and writers:

- Multiple processors can read from or write to shared memory simultaneously, requiring coordination to prevent conflicts. Protocols ensure only one processor writes to a data point at a time. This prevents data corruption and maintains consistency.

Updated word is distributed to all other processors:

- When a processor updates a data value, the new value is sent to all other processor caches. This keeps all processors working with the latest data. However, it generates additional memory traffic.

Some systems use an adaptive mixture of both solutions:

- Some systems adaptively combine "write-through" and "write-back" caching. This helps balance the need for data consistency with performance. The system adjusts its approach based on current traffic and update frequency.

Q30) Explain Write Invalidate.

A:

Write Invalidate:

- In a "write-invalidate" strategy, when one processor writes to a memory location, it first invalidates that data in all other caches. This ensures that no other processor will use stale data. Only the writing processor will hold the valid version until it's needed by others.

Multiple readers, one writer:

- In this approach, multiple processors can read a data line, but only one processor can write to it at a time. This prevents conflicts, as only one processor has the authority to modify data. Other processors will see the latest data when they request access again.

When a write is required, all other caches of the line are invalidated:

- When a processor needs to write, it sends an invalidate signal to other processors' caches. This makes all other copies of that data line "invalid." This way, only the writing processor has a valid, exclusive copy of the data line.

Writing processor then has exclusive access until line required by another processor:

- After invalidating other caches, the writing processor has exclusive access to the data. This continues until another processor requests that data, at which point the cache line might be shared again. This avoids conflicts by giving temporary ownership to the writer.

State of every line is marked as modified, exclusive, shared, or invalid:

- Each cache line is labeled with a state: Modified (only one processor has changed it), Exclusive (one processor holds it without modifying), Shared (multiple processors have read-only access), or Invalid (data is

outdated). These states help the system maintain coherence efficiently.

MESI protocol:

- The MESI (Modified, Exclusive, Shared, Invalid) protocol is a cache-coherence protocol used to implement the above states. It helps manage cache states in a multi-processor system, ensuring that data consistency is maintained while optimizing memory performance.

Q31)Explain MESI Protocol.

A:

Here's a brief explanation of each aspect of the MESI protocol:

1. Commonly implemented for Cache Coherence:

- The MESI protocol is widely used to maintain cache coherence in multi-core systems. It ensures that processors don't work with stale or inconsistent data. By defining cache states, it coordinates how and when data is shared or invalidated across caches.

2. Four states of a cache line:

- MESI defines four possible states for each cache line: Modified, Exclusive, Shared, and Invalid. Each state represents the current status of data consistency between caches and main memory. This helps optimize data access while maintaining consistency.

3. Invalid:

- A cache line marked as "Invalid" is outdated or unused, meaning the data in this line is no longer valid for that cache. This happens when the data has been changed in another cache or is no longer relevant to the current processor.

4. Exclusive:

- When a cache line is in the "Exclusive" state, this cache is the only one holding a copy of the data, and it matches the data in main memory. No other caches hold this line, so the processor can modify it without notifying others initially.

5. Shared:

- In the "Shared" state, multiple caches hold copies of this line, and the data matches main memory. Any of these caches can read the data, but if a write is needed, the line must be invalidated or changed to a different state to keep consistency.

6. Modified:

- A cache line in the "Modified" state has been changed by a processor, so it differs from the main memory copy. This cache has the only up-to-date version, and before any other cache can use it, the line must be written back to main memory.

Q32)What is Interleaved Memory?

A:

Here's an explanation of interleaved memory in three lines for each point:

1. Interleaved memory compensates for the slow speed of DRAM:

- DRAM access can be slower than the CPU's processing speed, creating delays. Interleaved memory divides memory into banks that can be accessed independently. This allows simultaneous access, reducing wait times for memory operations.

2. Spreads memory addresses evenly across banks:

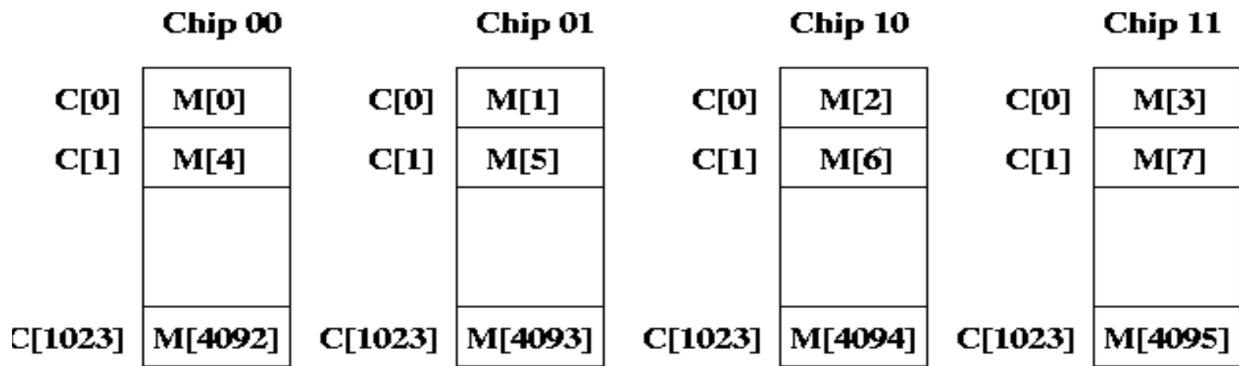
- Memory addresses are divided across multiple banks in an alternating fashion. For example, address 0 might go to Bank 1, address 1 to Bank 2, and so on. This arrangement allows multiple addresses to be accessed at once, boosting efficiency.

3. Enables contiguous memory reads and writes:

- With interleaving, contiguous addresses are spread across banks, so consecutive data can be read or written in parallel. This parallel access is faster than accessing a single memory bank sequentially.

4. Results in higher memory throughput due to reduced waiting:

- By accessing multiple banks at once, interleaved memory minimizes idle CPU time waiting for data. This results in higher overall memory throughput and faster system performance for memory-intensive tasks.



Q33) What is Associative Memory?

A:

Here's a breakdown of associative memory with explanations in 2-3 lines for each point:

1. Content-addressed or associative memory:

- Associative memory, also known as content-addressable memory (CAM), retrieves data based on its content rather than its specific memory address. This allows data to be accessed by matching content, making it efficient for search-based operations.

2. Reference clues are "associated" with actual memory contents:

- In associative memory, reference clues are used to search for the desired data. The memory searches for matches based on the

content provided, and once a match (or set of matches) is found, the relevant data is returned, bypassing the need for specific addresses.

3. Humans retrieve information by linking to related information:

- Associative memory resembles how humans remember information by linking it to other related pieces. This makes retrieval faster, as related data is associated and retrieved based on content rather than an exact “location,” similar to the way human memory functions.

Q33)What is Interleaved Memory?

A:

1. Allows more programs to be opened simultaneously:

- Virtual memory allows the system to run more programs than can fit in physical RAM by using the hard disk for temporary storage of memory pages. This makes it seem like there's more RAM available than physically exists.

2. 32 or 64MB of RAM available for CPU usage:

- The system uses a fixed amount of RAM (e.g., 32MB or 64MB) for active processes, but through virtual memory, it can simulate much more by swapping data to and from the hard disk.

3. Users expect all their programs to run at once:

- Virtual memory enables users to open multiple applications simultaneously, like email programs, web browsers, and word processors, even when there's not enough physical RAM to support them all at once.

4. Find RAM for areas that have not been used recently and copy them onto the hard disk:

- The operating system identifies sections of memory that have not been accessed recently and moves them to the hard disk. This frees up RAM for more active processes, maintaining multitasking without crashes.

5. Frees up space in RAM to load new applications:

- When RAM space is needed for a new program, virtual memory ensures that unused data is swapped out to the hard disk. This keeps the system running smoothly while allowing new applications to load.

6. Copying happens automatically (feels like unlimited RAM space):

- The process of swapping data between RAM and the hard disk is handled automatically by the operating system. This creates the illusion of unlimited RAM, even though physical RAM is limited.

7. Hard disk space is much cheaper than RAM chips, providing economic benefits:

- Since hard disk storage is cheaper than RAM, using it as temporary storage for memory pages is a cost-effective way to expand available memory without purchasing more expensive RAM.

8. Read/write speed of a hard drive is not geared toward small pieces of data:

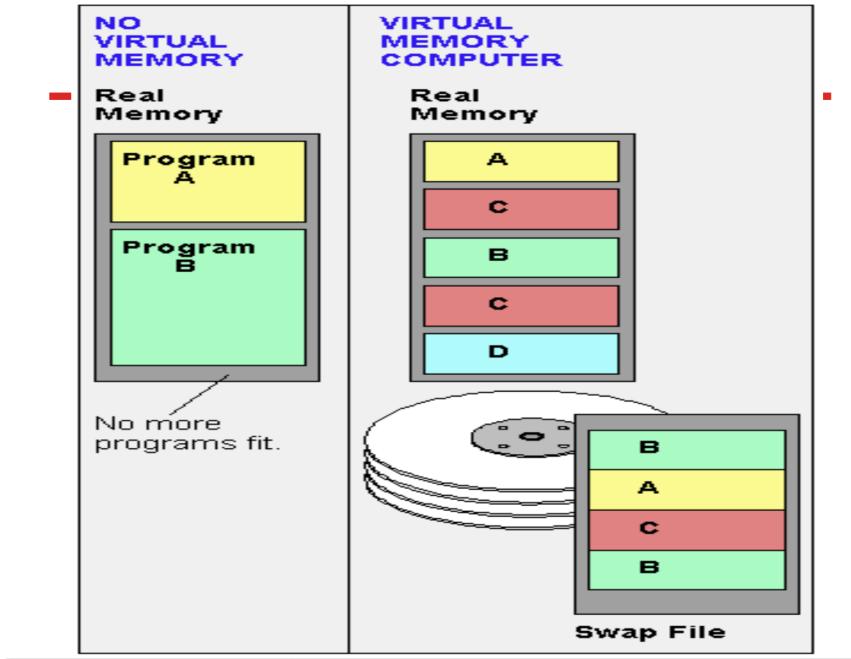
- Hard drives are slower than RAM and not optimized for accessing small chunks of data quickly, which can slow down performance when the system relies heavily on swapping data.

9. Operating system swaps information between RAM and the hard disk:

- The operating system is responsible for constantly swapping data between the faster RAM and the slower hard disk. This allows more programs to run, but requires careful management to avoid excessive swapping.

10. Thrashing:

- When the system spends too much time swapping data between RAM and the hard disk, it leads to "thrashing." This makes the computer feel incredibly slow because the CPU is mostly occupied with swapping rather than processing tasks.



Q34) What is Paging?

A:

Paging is a memory management scheme that eliminates the need for contiguous memory allocation by dividing both the physical memory and the logical memory into equal-sized blocks. Here's an explanation of the points related to paging:

1. Unequal fixed size / Variable Size partitions (Inefficient):

- Traditional memory allocation methods used unequal or variable-sized partitions, leading to fragmentation. This fragmentation could result in wasted memory space. Paging overcomes this by using fixed-sized partitions, making memory management more efficient.

2. Primary memory is divided into small equal fixed-sized partitions

(256, 512, 1K) called page frames:

- In paging, the physical memory (RAM) is divided into fixed-sized blocks, known as page frames. These sizes could range from 256 bytes to 1 KB or larger, depending on the system's architecture.

This division helps manage memory more predictably and reduces fragmentation.

3. Processes are divided into the same-sized blocks (pages) called

paging:

- The process's memory is divided into equal-sized units called pages, corresponding to the page frames in memory. Each page is a logical unit of memory that can be mapped to a specific frame in physical memory. This makes it easier to manage and access memory locations efficiently.

4. Recently referenced pages in memory:

- To optimize performance, frequently accessed pages are kept in physical memory, ensuring faster access. The operating system tries to maintain a working set of pages that are most relevant to the current process. This reduces the need to constantly fetch pages from slower secondary storage.

5. Need a page table for this management:

- The page table is a data structure that stores the mapping of virtual memory addresses (pages) to physical memory addresses (page frames). The operating system uses this table to translate virtual addresses to physical addresses during memory access. It is essential for managing the paging mechanism and ensuring correct memory allocation.

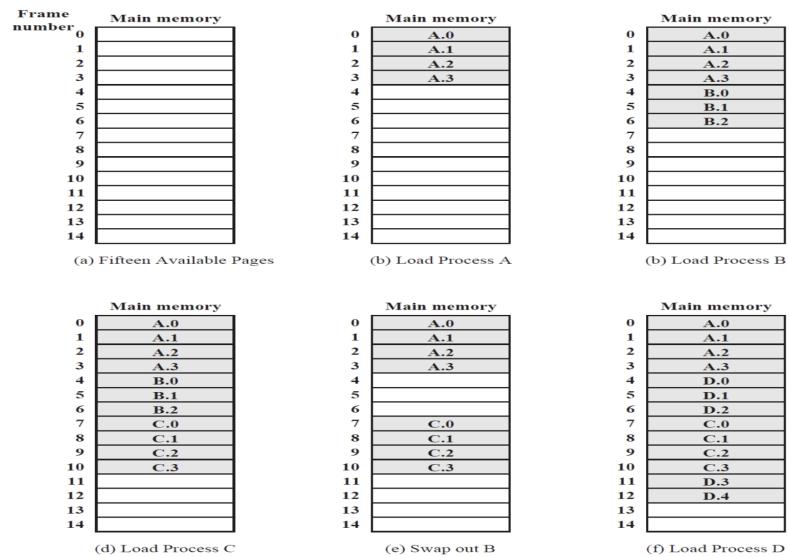


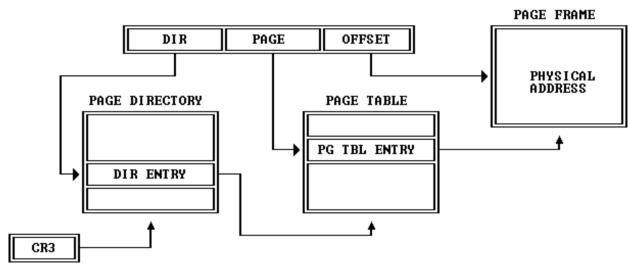
Figure 7.9 Assignment of Process Pages to Free Frames

Page Table Sample

Figure 5-8. Format of a Linear Address



Figure 5-9. Page Translation



Q35)What is Segmentation?

A:

1. Paging → Internal fragmentation:

- Paging can lead to internal fragmentation because memory is divided into fixed-size pages, and if a process doesn't fully use the allocated page, the unused space within the page is wasted. This inefficiency occurs when processes do not perfectly match the page size, leaving small gaps. Segmentation, on the other hand, avoids this by allowing variable-sized memory partitions.

2. Segmentation maps segments representing data structures, modules, etc., into variable partitions:

- In segmentation, a process is divided into segments that represent logical components like data structures, code modules, or stack frames. Unlike paging, segments can vary in size depending on the program's requirements. This allows for a more flexible and efficient use of memory.

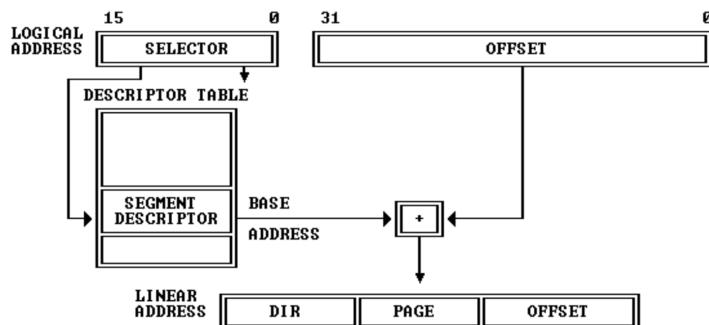
3. Not all segments of a process are loaded at a time:

- Unlike paging, where the entire process is divided into pages that are loaded into memory, in segmentation, only the necessary segments of a process are loaded into memory. This allows for more efficient memory utilization as large or unused segments do not occupy space in RAM.

4. We need a segment table very much like a page table:

- A segment table is used in segmentation to keep track of the mapping between logical addresses and physical memory addresses. Similar to a page table in paging, the segment table helps translate the segment address into the corresponding location in physical memory. This ensures that memory is accessed correctly and efficiently.

Figure 5-2. Segment Translation



Q36)What is Main Memory Allocation?

A:

1. Memory is divided into a set of contiguous locations called regions/segments/pages:

- Memory is organized into smaller, logically distinct units such as regions, segments, or pages. These units can vary in size depending on the memory management method used, like

segmentation or paging. Contiguous memory allocation helps manage and access data more efficiently.

2. Store blocks of data:

- Memory units (regions, segments, or pages) store blocks of data, which can include program instructions, variables, or other information. These blocks are the fundamental units that the CPU accesses during execution. Efficient allocation of these blocks ensures that memory is utilized effectively.

3. Placement of blocks of information in memory is called Memory Allocation:

- Memory allocation refers to how memory blocks (regions, segments, or pages) are assigned to store data. This can involve dynamic allocation, where memory is assigned as needed, or static allocation, where memory is predefined. Effective memory allocation ensures optimal use of available space.

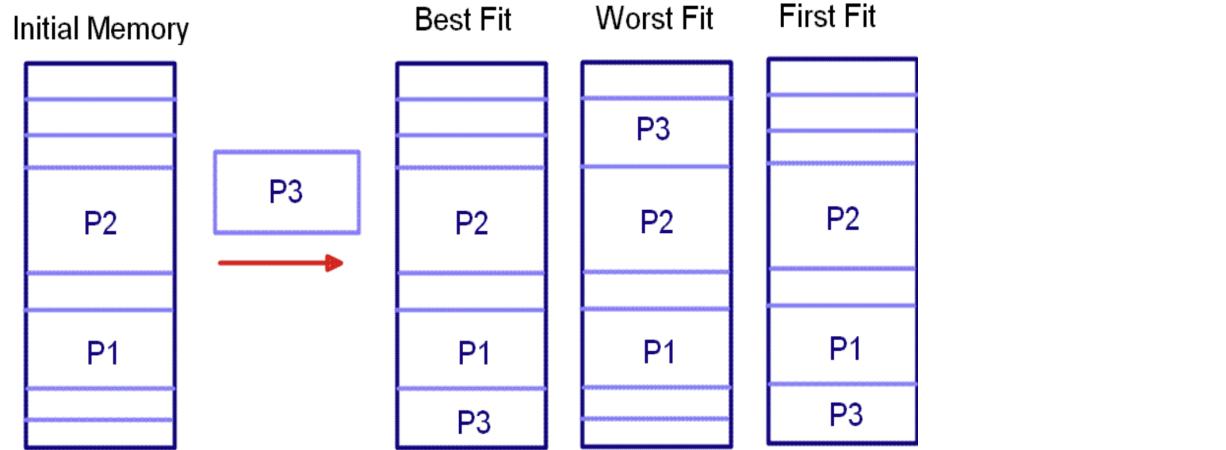
4. Memory Management Systems keep information in a table containing available and free slots:

- A memory management system tracks the status of memory slots using a table that identifies which areas are allocated and which are free. This table helps efficiently manage memory by allocating free memory slots to programs or processes, avoiding fragmentation, and ensuring smooth operation.

Allocation is done only as per needs

- First Fit

- Best Fit



Q37) What is Replacement Algorithm?

A:

Here's a detailed explanation of each memory replacement algorithm in 2-3 lines:

1. **Hardware implemented algorithm (speed):**

- Hardware-implemented algorithms are designed to execute quickly and efficiently by utilizing specialized hardware for managing memory. These algorithms help speed up operations like cache management by offloading tasks from the CPU to dedicated circuits or processors.

2. **Least Recently Used (LRU):**

- LRU replaces the cache slot that hasn't been used for the longest period of time. It assumes that data that has not been accessed recently is less likely to be needed soon. This helps optimize cache usage by keeping frequently accessed data in memory.

3. First In First Out (FIFO):

- FIFO replaces the oldest cache entry, the one that has been in the cache the longest. It is a simple approach, where data is removed in the order it was brought into the cache. While easy to implement, FIFO doesn't always yield the best performance because it doesn't consider the frequency of data access.

4. Random:

- The Random algorithm randomly selects a cache slot to replace when needed. While not optimal, it can be effective in scenarios where other algorithms are too complex or where cache access patterns are unpredictable. Its simplicity makes it easy to implement, though it doesn't necessarily minimize cache misses.

5. OPT-Optimal (Future):

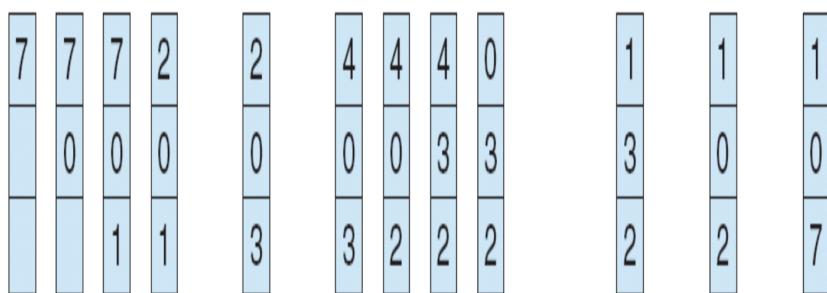
- The Optimal (OPT) algorithm replaces the page that will not be used for the longest time in the future. This theoretical algorithm provides the best possible performance but is not practical because it requires knowledge of future memory accesses, which is typically unavailable in real systems.

LRU Replacement

LRU Page Replacement

reference string

7 0 1 2 0 3 0 4 2 3 0 3 2 1 2 0 1 7 0 1



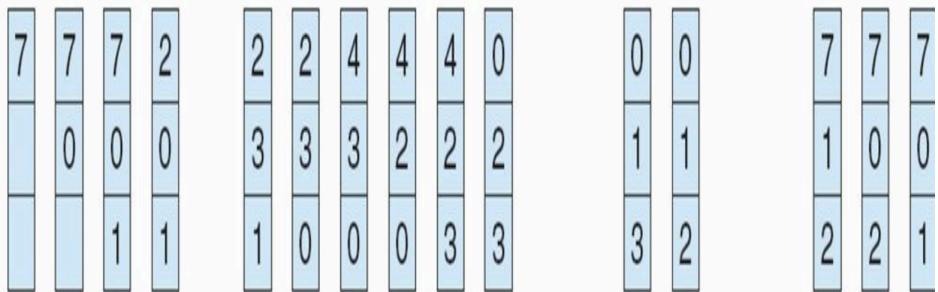
page frames

FIFO Replacement

FIFO Page Replacement

reference string

7 0 1 2 0 3 0 4 2 3 0 3 2 1 2 0 1 7 0 1



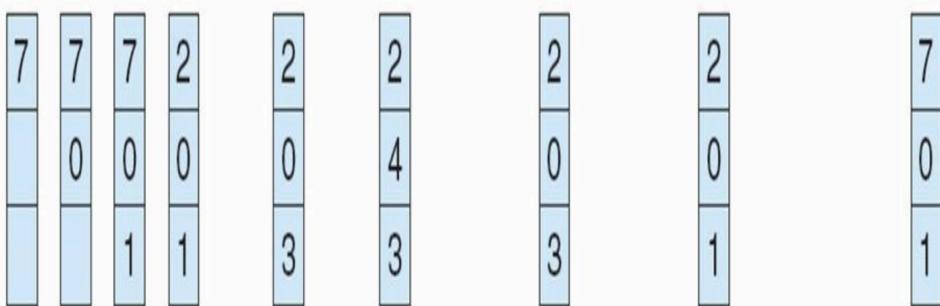
page frames

OPT Replacement

Optimal Page Replacement

reference string

✓ ✓ ✓ ✓ ✓ ~ max no. of hits
7 0 1 2 0 3 0 4 2 3 0 3 2 1 2 0 1 7 0 1
occur least no. of times to replace.



page frames

2 -
0 -
1 -

Q38) Name some Secondary Storage devices.

A:

Secondary Storage

- Magnetic disks
- Floppy disks
- Magnetic Tape
- RAID
- Optical Memory
- CD-ROM
- DVD

Q39) What is Raid Level?

A:

RAID Levels 0 - 6

REDUNDANT ARRAY OF INDEPENDENT DISKS

1. Storage is an important consideration when setting up a server:

- When setting up a server, selecting the appropriate storage solution is critical to ensure data reliability, speed, and availability. Storage should be chosen based on performance,

redundancy, and the specific requirements of the applications or services running on the server.

2. Almost all of the important information that you and your users care about will at one point be written to a storage device to save for later retrieval:

- Data storage is essential for long-term data retention and retrieval. Whether it's user data, system logs, or application files, reliable storage is necessary to ensure that critical information can be accessed when needed, without risk of loss or corruption.

3. Single disks can serve you well if your needs are straightforward:

- For simple applications with minimal performance and redundancy demands, a single disk may suffice. However, it lacks fault tolerance, meaning if the disk fails, all data on it is lost.
Single disks are suitable for small-scale or less critical environments.

4. However, if you have more complex redundancy or performance requirements, solutions like RAID can be helpful:

- RAID (Redundant Array of Independent Disks) is a solution for systems that require better redundancy, fault tolerance, and/or performance. By combining multiple disks in various configurations, RAID can protect against disk failures and improve read/write speeds, depending on the RAID level chosen.

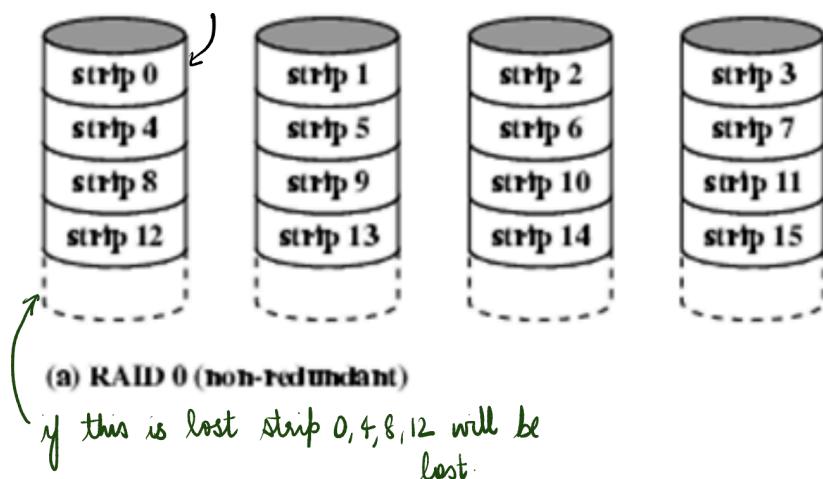
RAID Levels Overview:

- RAID 0 (Striping):
 - Data is split across multiple disks (striped) for improved performance but no redundancy. If one disk fails, all data is lost.

→ data will be stored in disks.

RAID Level 0- Non Redundant

each disk carries data

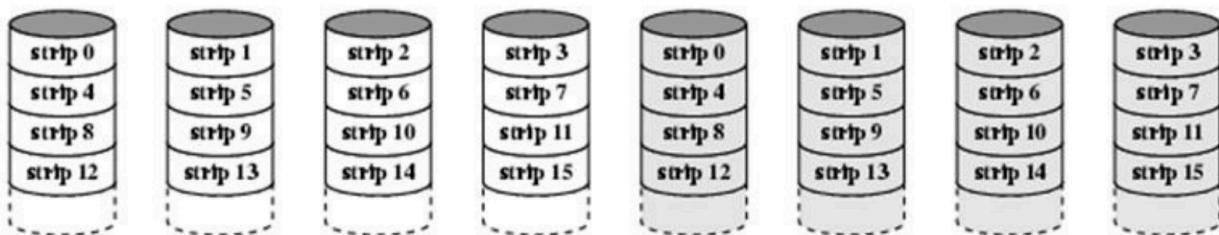


- RAID 1 (Mirroring):

- Data is duplicated on two or more disks. Provides redundancy, meaning if one disk fails, the data is still available on the other disk, but with no performance improvement.

RAID Level -1 Mirrored

duplicate copy is available

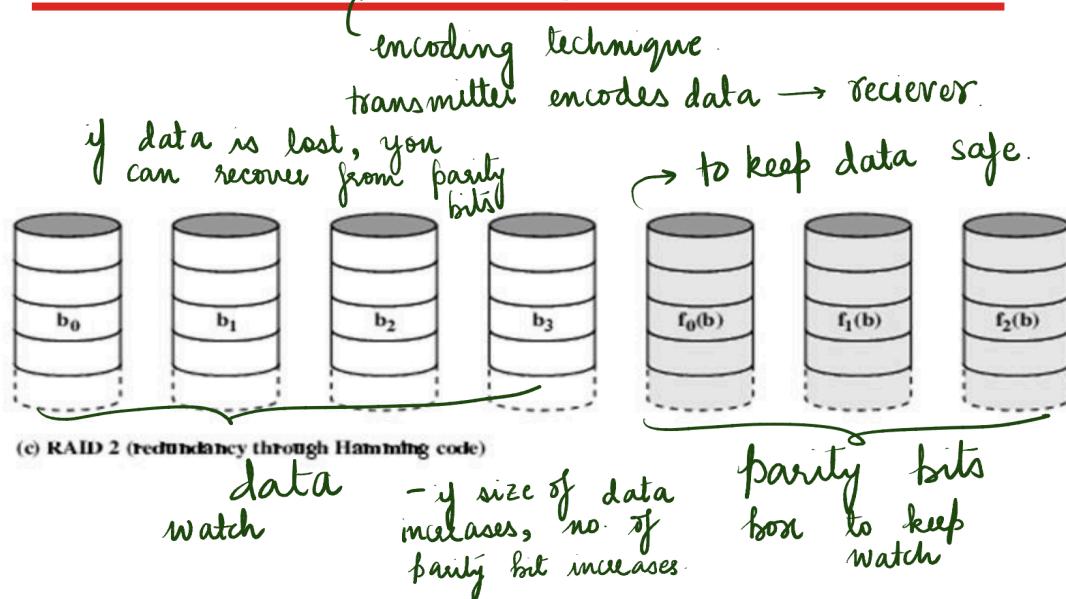


(b) RAID 1 (mirrored)

RAID Level 2 - Hamming Code:

RAID 2 uses bit-level striping with **Hamming code** for error correction. Each bit of data is stored across different disks, and error correction information is stored on additional disks using the Hamming code. This setup provides fault tolerance, but it is rarely used today due to inefficiency and the existence of better alternatives like RAID 3 and RAID 5

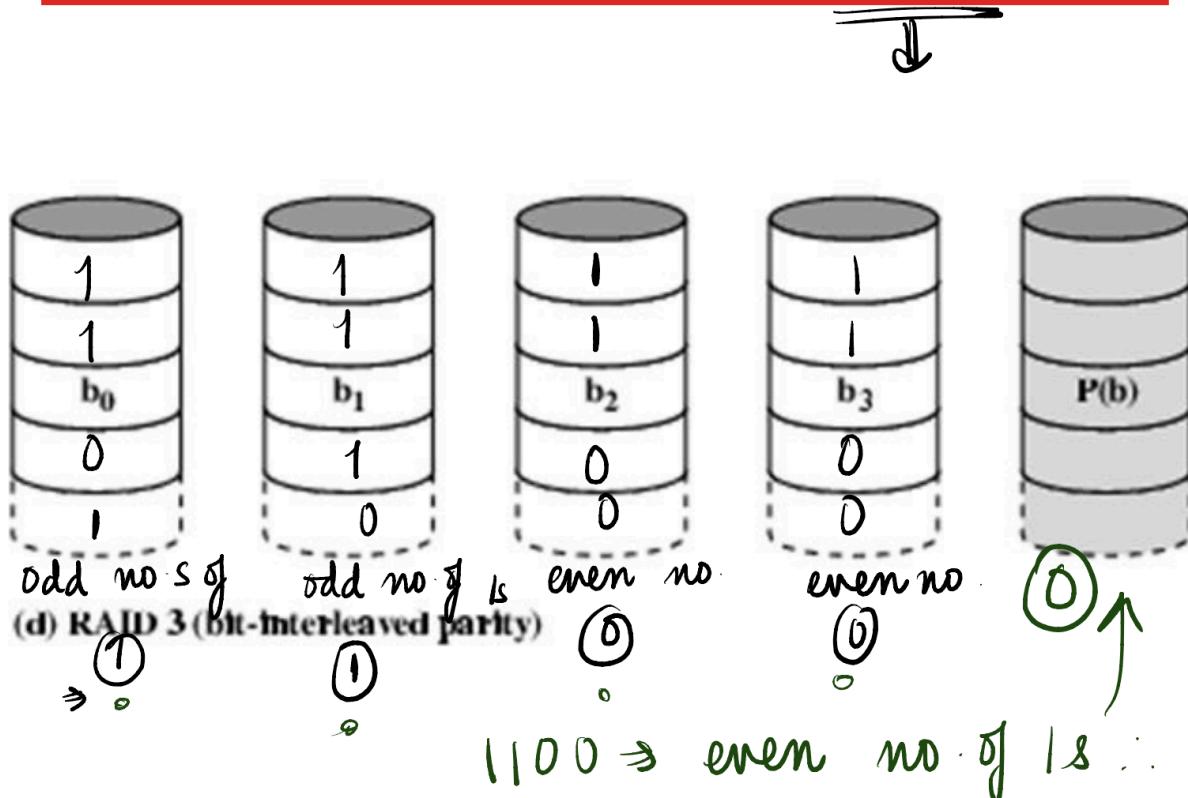
RAID Level 2- Hamming Code



RAID Level 3 - Bit Interleaved Parity:

- RAID 3 uses **bit-level striping**, where data is broken into bits and distributed across disks, with one dedicated disk for **parity**. This parity disk holds error-checking information, allowing the system to recover from a single disk failure. The main drawback is the potential bottleneck of the parity disk during write operations.

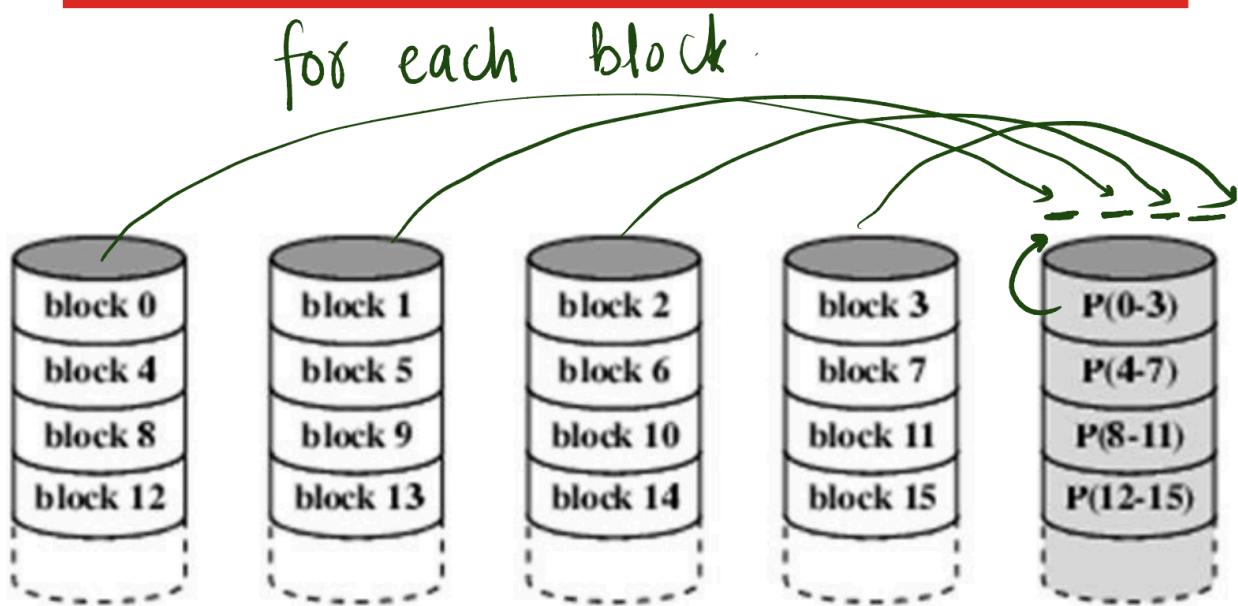
RAID Level 3 –Bit Interleaved Parity



RAID Level 4 - Block Level Parity:

- RAID 4 uses **block-level striping** where data is divided into fixed-size blocks and distributed across multiple disks, while **parity** is stored on a dedicated disk. This allows fault tolerance, but write performance can be slow due to the dedicated parity disk, which creates a bottleneck when updating data.

RAID Level 4- Block level parity



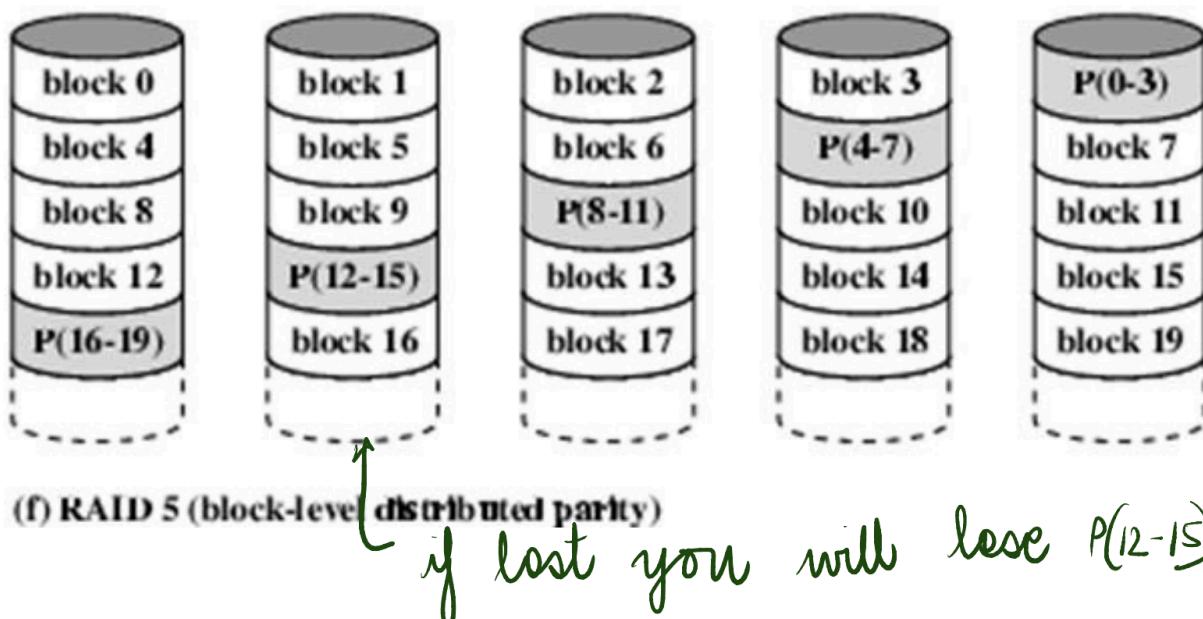
(e) RAID 4 (block-level parity)

RAID Level 5 - Block Level Distributed Parity:

- RAID 5 also uses **block-level striping**, but unlike RAID 4, the **parity data is distributed across all the disks** in the array. This improves performance since there is no dedicated parity disk, and the array can tolerate the failure of one disk without losing data. RAID 5 is commonly used for balancing fault tolerance and performance.

RAID Level 5- Block level Distributed Parity

distributed parity

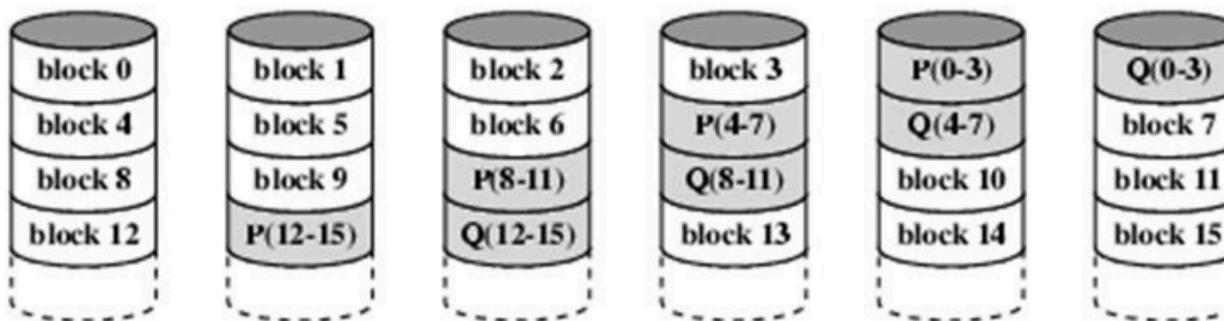


RAID Level 6 - Dual Redundancy:

- RAID 6 is similar to RAID 5 but adds **double parity**, storing two sets of parity data distributed across the disks. This provides **dual redundancy**, allowing the system to tolerate the failure of two disks simultaneously. It offers higher fault tolerance but comes at the cost of slower write performance due to the additional parity calculations.

RAID Level 6- Dual Redundancy

MIRRORED → RAID LEVEL 6



(g) RAID 6 (dual redundancy)

