

Introduction to Cache Memory

Cache Memory

- ❖ Cache Memory is an important component of modern computer systems that helps improve their performance.
- ❖ Cache memory is a small and fast memory that stores data and instructions that are frequently accessed by the CPU.
- ❖ The main purpose of cache memory is to reduce the access time of data from the main memory.
- ❖ The cache memory is placed between the CPU and main memory, and it stores a copy of the most frequently used data and instructions. When the CPU needs to access data, it first checks the cache memory.
- ❖ If the data is found in the cache memory, it is retrieved from there, which is much faster than accessing the data from the main memory.

Architecture of Cache Memory

- Cache memory is placed between the CPU and main memory and is divided into a set of blocks or lines. Each block in the cache memory corresponds to a block or line in the main memory. The cache memory uses a tag to store the address of the block in the main memory. The tag is used to identify the data or instruction stored in the cache memory. The cache memory also has a valid bit to indicate whether the block in the cache memory is valid or not.

Working Principle of Cache Memory

- When the CPU needs to access data, it first checks the L1 cache memory. If the data is found in the L1 cache memory, it is retrieved from there, which is much faster than accessing the data from the main memory. If the data is not found in the L1 cache memory, the CPU checks the L2 cache memory. If the data is found in the L2 cache memory, it is retrieved from there. If the data is not found in the L2 cache memory, the CPU retrieves the data from the main memory and stores a copy of it in the cache memory for future use.
- Cache memory works on the principle of locality of reference, which means that data and instructions that are accessed frequently in a short period are likely to be accessed again in the near future.

Types of Cache Memory

- There are two types of cache memory, namely, level 1 (L1) cache and level 2 (L2) cache. L1 cache is the fastest and smallest cache memory, which is built into the CPU itself. It is the first memory the CPU checks for data and instructions. The L1 cache size ranges from a few kilobytes to a few megabytes, depending on the processor. L2 cache, on the other hand, is larger and slower than L1 cache, and it is located outside the CPU. It is shared among all the cores of the processor and has a size ranging from a few megabytes to several megabytes.

Some More info about types of cache memory that are used:-

1.Level 1 (L1) Cache:

L1 cache is the fastest and smallest cache memory that is integrated into the CPU. It is used to store frequently accessed data and instructions that the CPU needs to access quickly. L1 cache is split into two parts: data cache and instruction cache. The size of L1 cache is typically between 8 KB to 64 KB.

2.Level 2 (L2) Cache:

L2 cache is larger than L1 cache and is located outside the CPU. It is used to store frequently accessed data and instructions that the CPU needs to access quickly. L2 cache is slower than L1 cache but still faster than main memory. The size of L2 cache is typically between 256 KB to 8 MB.

3.Level 3 (L3) Cache:

L3 cache is even larger than L2 cache and is also located outside the CPU. It is used to store frequently accessed data and instructions that the CPU needs to access quickly. L3 cache is slower than L2 cache but still faster than main memory. The size of L3 cache is typically between 2 MB to 64 MB.

4.Unified Cache:

Unified cache is a type of cache memory that is used to store both data and instructions. It is typically found in systems where L1 cache is too small to hold both data and instructions. Unified cache provides faster access to both data and instructions than main memory.

5.Write-Through Cache:

Write-through cache is a type of cache memory in which any write operation to the cache is also written to the main memory immediately. This ensures that the data in

the main memory is always up to date. Write-through cache is slower than write-back cache, but it ensures data consistency.

6. Write-Back Cache:

Write-back cache is a type of cache memory in which any write operation to the cache is not immediately written to the main memory. Instead, the data is written to the main memory only when the cache line is replaced or when the system is shut down. Write-back cache is faster than write-through cache, but it can cause data inconsistency in case of system failure.

7. Inclusive Cache:

Inclusive cache is a type of cache memory in which the contents of a higher level cache (e.g., L2 cache) include all the contents of the lower level cache (e.g., L1 cache). This ensures that any data or instructions that are in L1 cache are also present in L2 cache. Inclusive cache helps to reduce cache misses and improve system performance.

8. Exclusive Cache:

Exclusive cache is a type of cache memory in which the contents of a higher level cache (e.g., L2 cache) do not include the contents of the lower level cache (e.g., L1 cache). This ensures that any data or instructions that are in L1 cache are not duplicated in L2 cache. Exclusive cache helps to reduce cache conflicts and improve system performance.

Page Replacement Policies

Page replacement is an essential component of the virtual memory system in modern computer systems. It is used when the operating system needs to bring a page from the secondary storage (hard disk) into the main memory to be used by the processor. When the main memory is full, the operating system needs to select a page to be evicted from the main memory to make room for the new page. The process of selecting a page to be evicted is known as page replacement. In this article, we will discuss the various page replacement policies used in modern computer systems.

When the cache memory is full, and the CPU needs to access a new block of data or instruction, the cache controller needs to replace an existing block in the cache with the new block. The process of selecting which block to replace is known as page replacement. There are various page replacement policies used in the cache memory, and each policy has its advantages and disadvantages.

Some of the commonly used page replacement policies are 😊

1.First In First Out (FIFO): As mentioned before, this policy removes the oldest page in memory when a new page is to be brought in. FIFO is simple to implement and requires no additional bookkeeping. However, it does not consider the frequency or importance of a page, which can result in thrashing. Thrashing occurs when the operating system spends more time swapping pages in and out of memory than executing processes, which can severely impact system performance.

2.Least Recently Used (LRU): LRU is a popular page replacement policy that removes the least recently used page in memory when a new page is to be brought in. LRU takes into account the history of page usage and tries to keep the most frequently used pages in memory. LRU requires maintaining a timestamp or counter for each page to determine which page was least recently used. Implementing LRU can be challenging, especially in large systems with many pages, as it requires constant updates to the timestamp or counter for every page accessed.

3.Clock or Second Chance: The Clock or Second Chance algorithm is an improvement over the FIFO algorithm. Instead of removing the oldest page, it gives each page a "second chance" before being removed. Pages are initially marked as not referenced, and when a page needs to be removed, the algorithm scans a circular buffer and gives a second chance to the first page it encounters that has not been recently referenced. If all pages have been recently referenced, the algorithm continues the scan until it finds a page that has not been referenced recently.

4.Least Frequently Used (LFU): LFU is a page replacement policy that removes the page that has been used the least number of times. LFU requires maintaining a counter for each page to determine how many times it has been

accessed. Pages that are not frequently used are more likely to be replaced, freeing up memory for more frequently used pages.

5. Most Recently Used (MRU): MRU is a page replacement policy that removes the page that has been used most recently. MRU is useful when there are pages that are frequently used and need to be kept in memory. MRU requires maintaining a timestamp or counter for each page to determine which page was most recently used. MRU is simple to implement and can be useful in systems where frequently used pages are more likely to be accessed again soon.

6. Optimal Page Replacement (OPT): OPT is a theoretical page replacement policy that removes the page that will not be used for the longest time in the future. OPT requires knowledge of future page accesses, which is not possible in practice. However, OPT is useful for evaluating the performance of other page replacement policies, as it provides an upper bound on their effectiveness.

Conclusion

Page replacement is an essential component of the virtual memory system in modern computer systems. The choice of page replacement policy affects the performance of the system, and different policies have their advantages and disadvantages. The random, FIFO, LRU, clock, and optimal policies are some of the commonly used page replacement policies. The operating system needs to select a policy that balances the efficiency and simplicity of implementation.

Hopes You love the Documentation and if there is any mistake in there please do let me know i would be more than happy to correct that!!!

// Contributed By Rahul Gupta <https://github.com/RahulGupta403>