

Q:- How is a block identified if present in the cache?

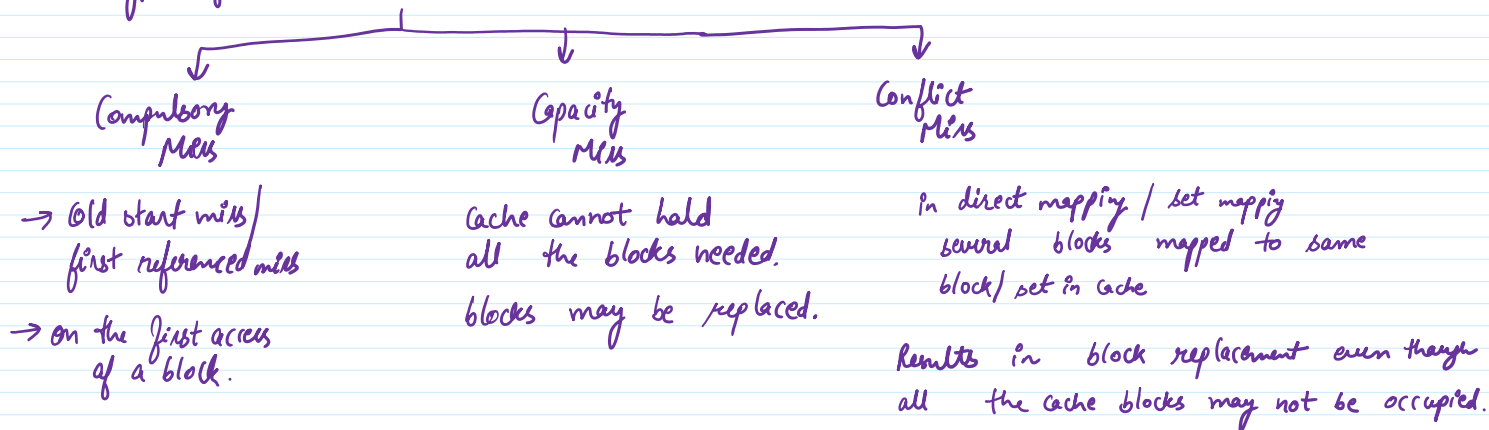
Caches include a TAG associated with each block.

— The TAG of every cache block where the block being requested may be present needs to be compared with the TAG field of Main Memory Address.

How many comparisons?

- 1) Direct Mapping :- One Comparison
- 2) Associative Mapping :- Full associative search over all TAGs of cache blocks.
- 3) Set associative Mapping :- Limited associated search over TAGs (only the selected SET).

Types of Cache Misses :-

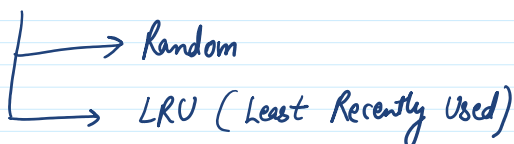


Q:- Which block should be replaced on a cache miss?

Direct Mapping :- trivial, no choice

Associative & Set Associative Mapping :-

there can be several blocks to choose from for replacement when a miss occurs.



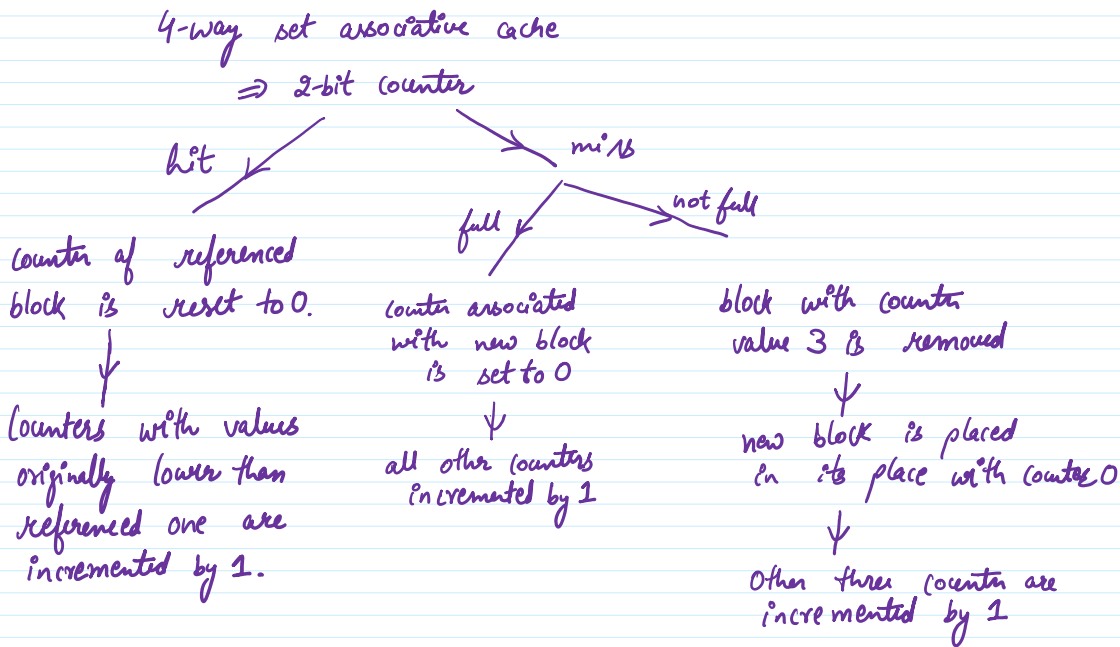
Replace the block which has not been used for the longest period of time.

Uses temporal locality:-

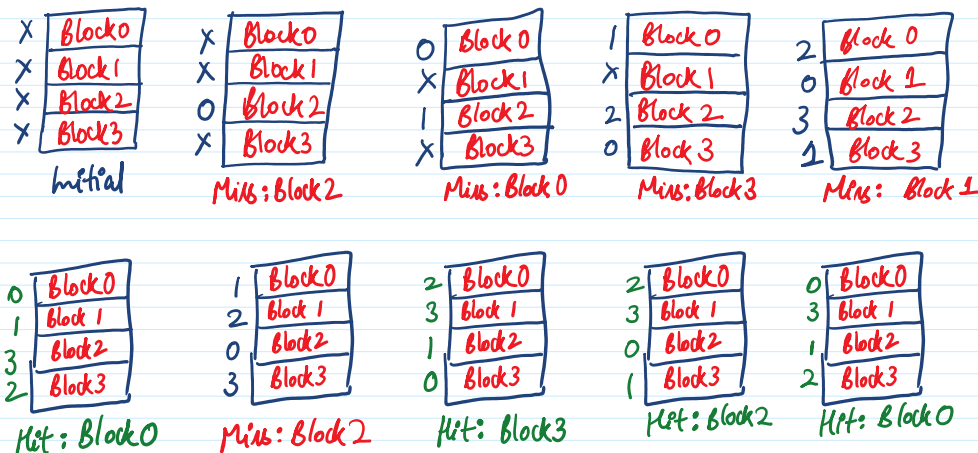
“If recently used blocks are likely to be used again, then best candidate for replacement is the LRU block.”

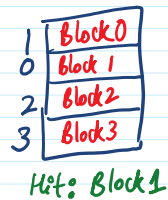
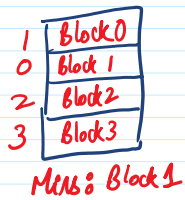
LRU cache block replacement policy:-

→ The cache controller tracks the LRU block using a counter.



Example:-





Improving Cache Performance :-

We know, for a two-level memory hierarchy, average memory access time,

$$T_{avg} = H_1 \times T_1 + (1 - H_1) \times T_2$$

Hit Time Miss Ratio Miss Penalty.

How can we improve the performance of cache memory?

Hit Time ↓ → avoiding the address translation when indexing the cache.

Miss Ratio ↓ → using larger block size, larger cache size, and higher associativity.

Miss Penalty ↓ → using multi-level caches.