

Cache memory:



→ L₂

→ Last Level Cache (LLC)

→ Main memory

Principle of locality

```

graph LR
    P[Principle of locality] --> T[Temporal locality]
    P --> S[Spatial locality]
  
```

Important terms :

→ core request (word)

→ cache block

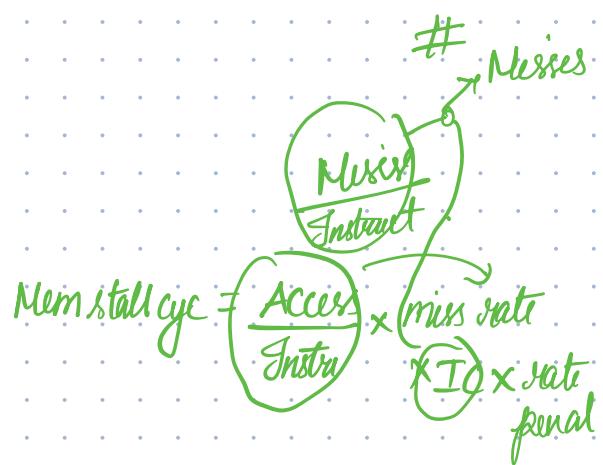
→ Cache hit / Cache miss

→ Cache replacement

CPU:

Instruction Execution cycle:

- IF : Instruction Fetch
- ID : Instruction Decode
- Ex : Execution
- MEM : Memory access
- WB : Write back



$$\text{CPU execution time} = (\text{clock cycle} + \text{memory stall cycles}) \times \text{clock cycle time}$$

Memory stall cycles = number of misses × miss penalty

$$= \frac{\text{misses}}{\text{instruction}} \times \text{IC} \times \text{miss penalty}$$

$$= \frac{\text{Access}}{\text{instruction}} \times \text{miss rate} \times \text{IC} \times \text{miss penalty}$$

$$\text{CPI} = \frac{\text{Clock cycle}}{\text{IC}}$$

$$\left| \begin{array}{l} \text{AMAT} = \text{Hit time} + \text{Miss rate} \times \text{miss penalty} \\ \text{CPU time} = (\text{CPU cycle} + \text{Mem stall}) \times \text{clock time} \end{array} \right.$$

$$\text{Mem stall} = \text{IC} \times \frac{\# \text{misses}}{\text{Instruction}} \times \text{penalty}$$

$$\# \text{misses} = \text{Access} \times \text{miss rate}$$

$$\text{CPT} = \frac{\text{Clock}}{\text{IC}}$$

$$\text{AMAT} = \text{Hit time} + \text{miss} \times \frac{\text{rate}}{\text{penalty}}$$

- Cache memory
- 1) $L_1 - I$
 - 2) $L_1 - D$
 - 3) LLC
 - 4) Main memory
 - 5) Secondary memory (Not discuss)

→ Principle of Locality

* Temporal locality : If a core request for a data, it will request it again in future.

* Spatial locality : If a core request for a data, it will request nearby data in future.

→ Core will request for word. But we will always fetch a block.

→ The first fetch is always compulsory miss.

32-bit



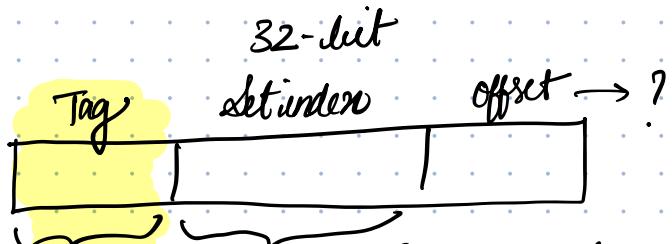
Cache memory → SRAM technology

→ Multiple words stored in block

Managing data in Cache → Block searching
→ Block placement
→ Block replacement

Average Memory Access Time (AMAT) = Hit time + Miss rate × Miss penalty.

Types in Cache → Direct mapped Cache
→ Fully associative Cache
→ Set-associative Cache



Word
↓
Block
↓
Set

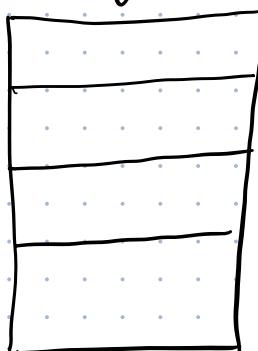
Direct mapped cache:

Way 1

Inside a set, the tag is used to search if the word already exists.

Based on this we decide which set it is mapping

Set 1 →
Set 2 →
:
Set n



Based on set index we will find which set the word belongs to

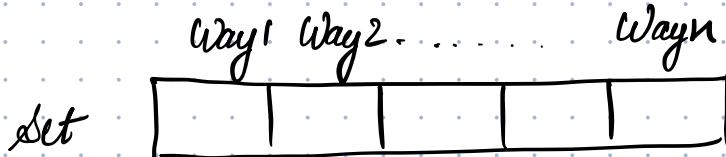


Cache

Problems :

- There is only one place to keep the block. Because it's 1 block per set.
- Conflict is more.

Fully associative Cache:

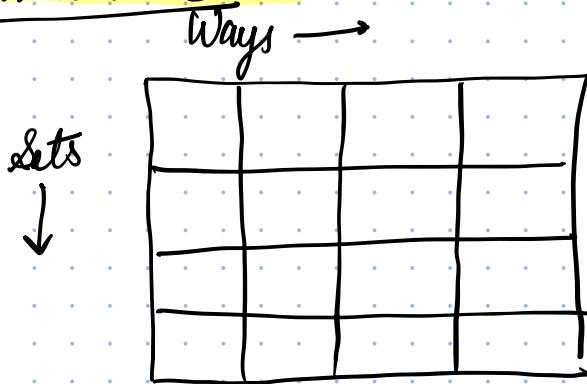


- Only 1 set
- But multiple ways
- No conflict miss
- High latency

Problems :

→ High latency because we have to search through all the ways by comparing the tags.

Set - associative Cache :



→ Best of both worlds

→ comparatively less conflict miss .

Example :

8-way set associative Cache :

$$\text{Size} = \text{ways} \times \text{sets} \times \text{blocks}$$

Types of misses :

Compulsory miss : First time accessing a block from main memory.

Conflict miss : A block is trying to occupy an existing cache line even if there are empty

Capacity miss : The miss occurs when the cache is full.

Multi processor :

- Multiprogram
- Multi threaded

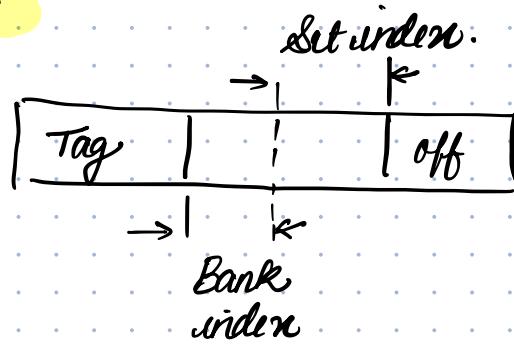
Cache in multi core:

- Shared Cache
- Private Cache

NUCA : Non-Uniform Cache Access :

- LLC is divided into banks (NUCA). ⇒ We need bank index.
- Each bank is set-associative.
- Banks near to the core have less latency

SNUCA : Static NUCA :



→ Easy search
→ More banks

→ If heavy use
bank is full,
latency would
be high
→ No bank migration

→ First find the bank → then set → Then search the ways using tag.

DNUCA : Dynamic NUCA :

- Allows migration of data inside bank-set.
- Introduces bank-sets : subgroup of banks.

→ Migrate closer
to core but
within the same
bank set.

Problem :

- Due to search,
Latency would be high because we have
search all the banks inside a bank set.
- Increases network congestion

NVLA in CMP:

CMP : Multiple cores & cache memory in a single chip

TCMP : Each core has its own small cache. But the L2 is not
Tile based
CMP

Cache coherence protocols :

→ MSI

→ MESI

$$CPU \text{ execution time} = (CPU \text{ clock cycle} + \text{Memory stall cycles}) \times \text{clock cycle time}$$

$$CPU \text{ clock cycle} = CPI \times IC$$

$$\text{Memory stall cycle} = \text{Number of misses} \times \text{penalty}$$

$$= \frac{\text{Number of misses}}{\text{Instructions}} \times IC \times \text{penalty}$$

$$= \frac{\text{Memory access}}{\text{Instruction}} \times \text{Miss rate} \times IC \times \text{penalty}$$