

MICROPROCESSOR AND COMPUTER ARCHITECTURE

UNIT-3

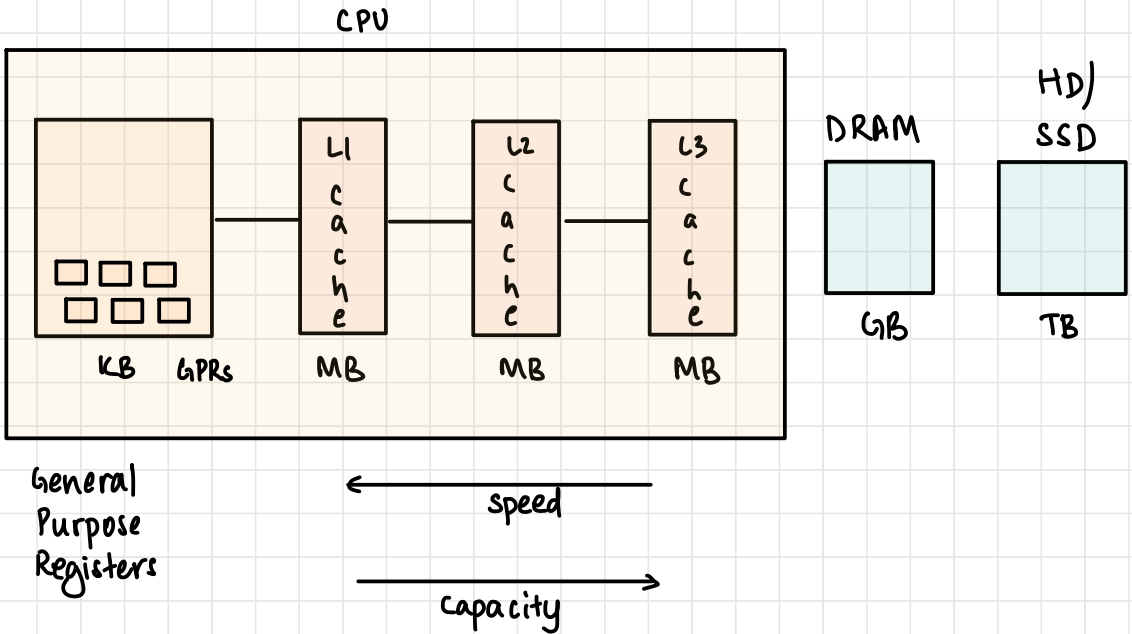
memory

feedback/corrections: vibha@pesu.pes.edu

VIBHA MASTI

MEMORY

- Processor speed vs memory - gap
- Registers: reside in CPU, faster to fetch from / write to; limited in number



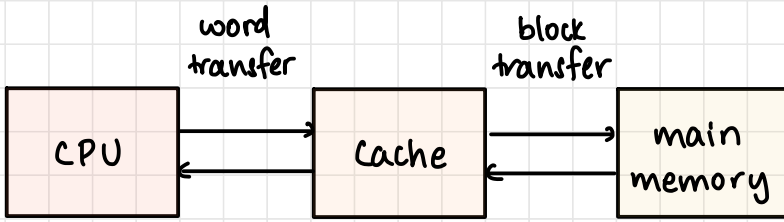
Locality of Reference

1) Temporal Locality

- set of instr likely to be referenced soon
- eg: looping

2) Spatial locality

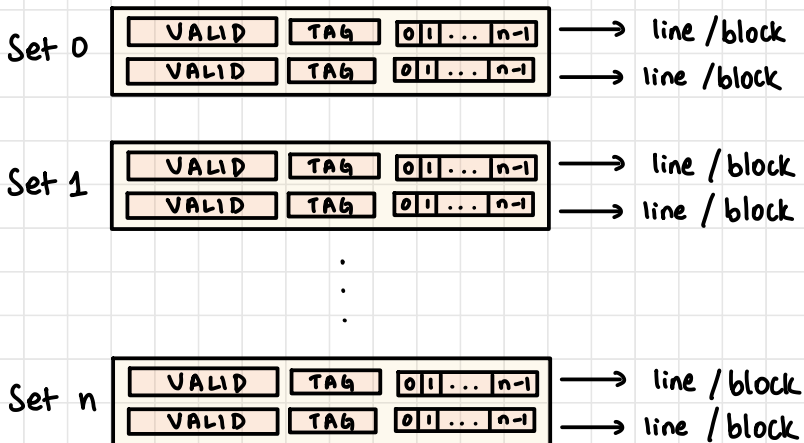
- if fetching from memory, fetch from other nearby locations
- eg: array elements



Cache mapping technique

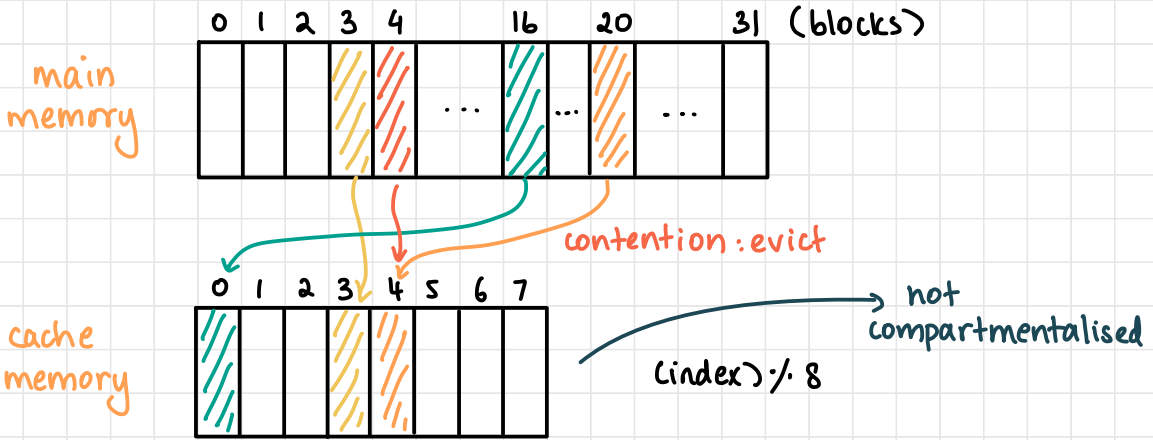
- 1) Direct Mapped Cache
- 2) Set Associative Mapping
- 3) Fully Associative Mapping

GENERAL ORGANISATION of CACHE



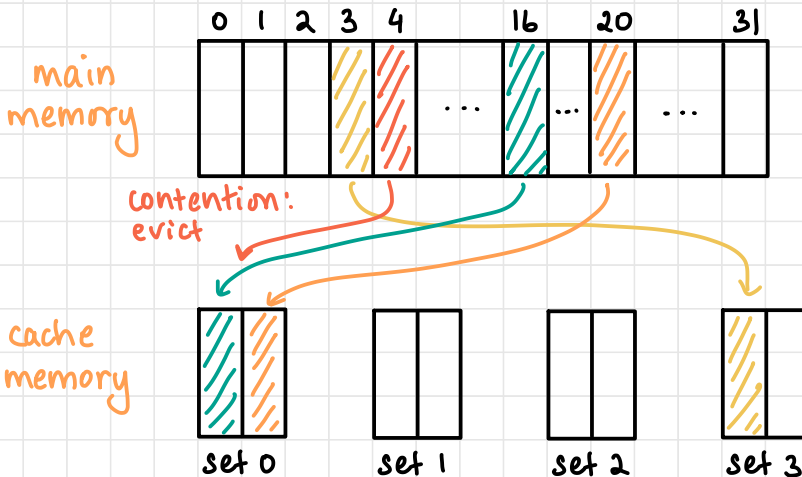
- compartmentalised into sets of 2 (rows)

1) Direct Mapped Cache



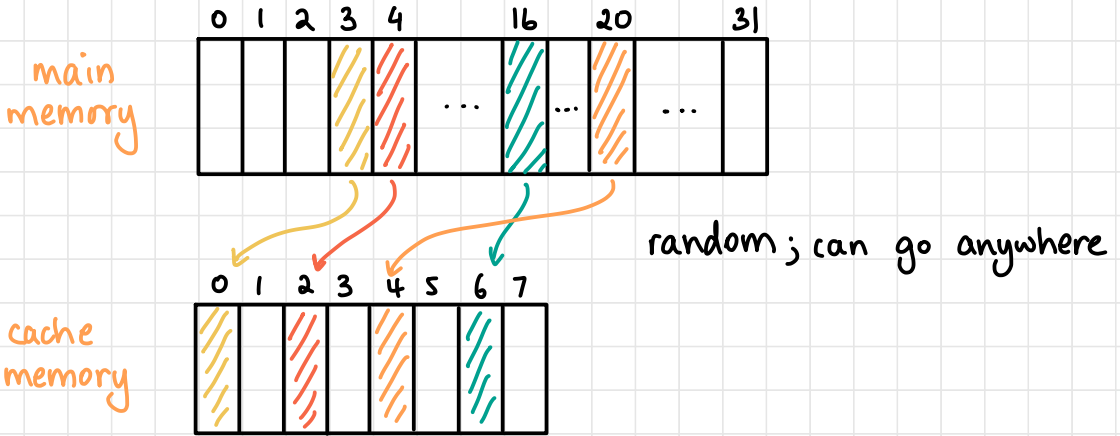
- (block no.) mod (no. of sets/blocks)
- leads to contention
- each line treated as block/set
- aka line mapping technique
- more than one block of memory mapped to one block in cache

2) Set Associative Mapping



- improvement over direct
- flexibility to place block anywhere in set
- better utilisation of space
(block no.) mod (no. of sets)

3) Fully Associative Mapping

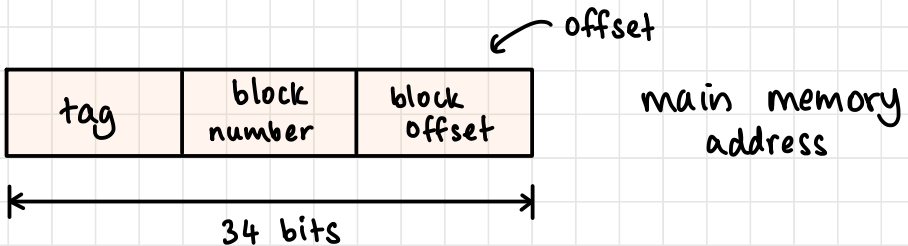


- any block of main memory mapped anywhere on cache
- search is slower

BITS NEEDED TO ADDRESS MAIN MEMORY

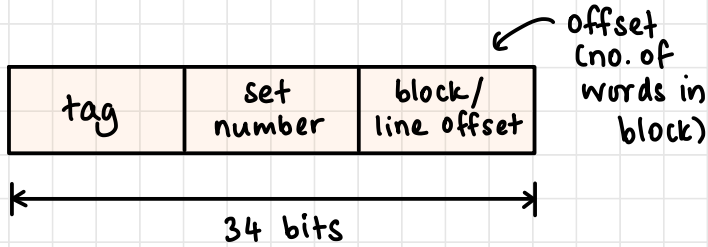
- eg: 16 GB RAM : $2^4 \times 2^{30} = 34$ bits

Direct Mapped Cache

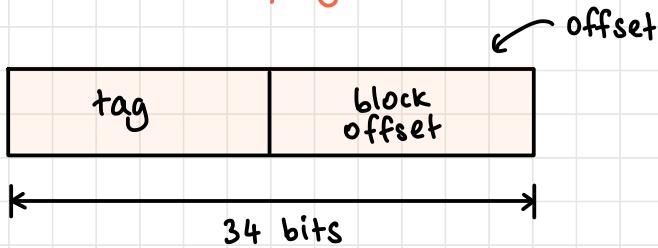


- valid 0: stale/junk data

Set Associative Mapping



Fully Associative Mapping



Q: A 4-way set associative cache memory unit with a capacity of 16 KB is built using a block size of 8 words. The word length is 32 bits. The size of the physical address space is 4 GB. How many bits are required for the tag field?

main mem = 4 GB

cache size = 16 KB = 2^{14} bytes

block size = 8×32 bits = 32 bytes = 8 words

$$\text{size of mem} = 4 \times 2^{30} = 2^{32}$$

\therefore 32-bit address for main mem

$$\text{size of set} = 4 \times 8 \times 4 = 2^7 \text{ bytes}$$

$$\therefore \text{no of sets} = \frac{16 \text{ KB}}{4 \times 8 \times 4} = \frac{16 \times 2^{10}}{16 \times 2^3} = 2^7$$

\rightarrow block

$$\text{block size} = 8 \times 4 = 32 \text{ bytes} = 2^5 \text{ bytes}$$

$$\text{block offset} = \text{word} = 5 \text{ bits}$$

\therefore 7 bits for sets

5 bits for word

\therefore 20 bits for tag



Q: Size of MM = 64 K words, cache = 128 blocks, 16 words/block, 1 word = 4 bits, 2-way set associative mapping. Find no. of bits for tag, set, word.

$$\text{memory} = 2^6 \times 2^{10} = 2^{16} \text{ words}$$

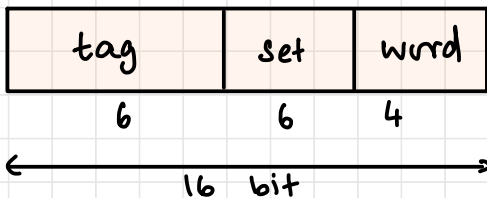
$$1 \text{ set} = 2 \text{ blocks} = 2 \times 16 \text{ words} = 32 \text{ words}$$

$$1 \text{ block} = 16 \text{ words} \Rightarrow 4 \text{ bits offset}$$

$$\begin{aligned} \text{main memory} &= 64 \times 2^{10} \text{ words} = 2^{16} \text{ words} \\ &= 16 \text{ bit addressing} \end{aligned}$$

$$\text{cache} \rightarrow 128 \text{ blocks} = 64 \text{ sets} = 2^6 \text{ sets}$$

\therefore 6 bits for set



Q: Consider direct mapped cache 512 kB, block size 1kB, 7 bits in tag. Find size of main memory. Word = 1 byte.

$$\text{cache size} = 512 \text{ kB} = 2^9 \text{ bytes}$$

$$\text{block size} = 2^{10} \text{ bytes} \Rightarrow 10 \text{ bits for offset}$$

$$\text{no of blocks} = 512 = 2^9 \text{ blocks}$$

$$\Rightarrow 9 \text{ bits for block}$$

$$\text{main memory address} = 7 + 10 + 9 = 26 \text{ bits}$$

$$\therefore \text{size of main memory} = 2^{20} \times 2^6 = 64 \text{ MB}$$

Q: Consider a 4-way SAC with block size = 4 KB. Size of MM = 16 GB. 10 tag bits. Find cache size.

$$\text{size of MM} = 16 \times 2^{30} = 2^{34} \text{ bytes}$$

$$\therefore \text{address} = 34 \text{ bit}$$

$$\text{tag bits} = 10 \text{ bits}$$

$$\text{size of block} = 4 \text{ KB} = 2^2 \times 2^{10} = 2^{12} \text{ bytes}$$

$$\therefore \text{block offset} = 12 \text{ bits}$$

$$\begin{aligned} \therefore \text{Set no. bits} &= 34 - (10 + 12) \\ &= 34 - (22) = 12 \text{ bits} \end{aligned}$$

$$\therefore \text{no of sets} = 2^{12}$$

$$\begin{aligned} \therefore \text{size of cache} &= 2^{12} \times 4 \times 4 \text{ KB} \\ &= 2^{12} \times 2^2 \times 2^2 \times 2^{10} \\ &= 2^{26} \text{ bytes} \\ &= 64 \text{ MB} \end{aligned}$$

$$\text{size of cache} = 64 \text{ MB}$$

Q: Cache = 512 KB, tag = 10, set = 8 blocks / set
Main memory = ?

cache size = no. of sets \times 8 \times size of block

$$\text{no of sets} = 2^x$$

$$\text{size of block} = 2^y$$

x: set bits

y: offset bits

$$2^9 \times 2^{10} = 2^{x+y+3}$$

$$19 = x+y+3 \Rightarrow x+y = 16$$

size of memory address = $10 + x + y$

$$= 26 \text{ bits}$$

$$\therefore \text{size of main memory} = 2^{26} = 64 \text{ MB}$$

Q: A computer system uses 16 bit memory, direct mapped 2KB cache, 64 bytes/block. Assume word is 1 byte.

calculate tag, block, word bits

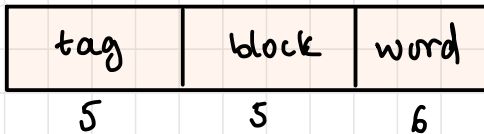
block size = 64 bytes = 2^6 bytes

\therefore block offset = 6 bits = word

$$\text{no. of blocks} = \frac{2 \times 2^{10}}{2^6} = 2^5$$

\therefore block = 5 bits

$$\text{tag} = 16 - (5 + 6) = 5 \text{ bits tag}$$



Q: A computer system uses 16 bit memory, 2-way set associative, 2KB cache, 64 bytes/block. Assume word is 1 byte.

(a) calculate tag, block, word bits

(b) Processor reads data sequentially from the following addresses: 128, 144, 2176, 2180, 128, 2176. Indicate hits and misses

(a) cache = $2 \times 2^{10} = 2^{11}$ bytes

set size = 128 bytes = 2^7

no. of sets = $\frac{2^{11}}{2^7} = 2^4 = 16$ sets

$$\therefore \text{set bits} = 4$$

$$\text{no. of words/block} = 64 \text{ words} = 2^6$$

$$\therefore \text{word bits} = 6$$

$$\therefore \text{tag bits} = 6$$



(b) 128, 144, 2176, 2180, 128, 2176 in base -10

	tag	set	offset/word
1. $(128)_{10} =$	000000	0010	000000
2. $(144)_{10} =$	000000	0010	010000
3. $(2176)_{10} =$	000010	0010	000000
4. $(2180)_{10} =$	000010	0010	000100
5. $(128)_{10} =$	000000	0010	000000
6. $(2176)_{10} =$	000010	0010	000000

1. will be miss; tag for first block of set 0010 is set to 000000
2. hit; set 0010, block 1 tag is 000000
3. miss; set 0010, block 2 tag = 000010
4. hit
5. hit
6. hit

VALID BIT

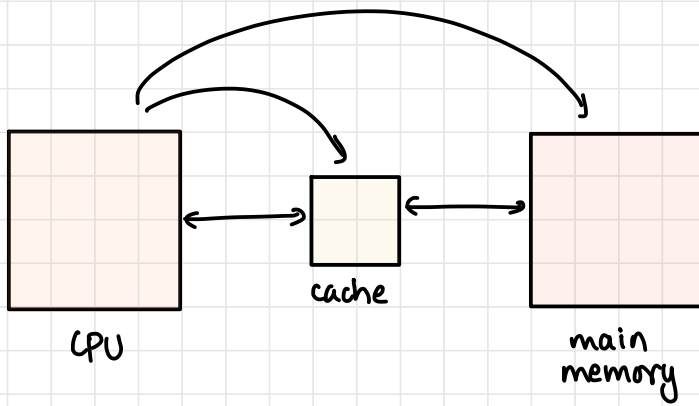
- Provided for each block
- Valid=0 when power turned on (stale/invalid memory)
- Valid=1 when block loaded
- If data on disk changes, main memory updated, cache valid bit set to 0

write hit

- (1) Write Through Protocol
- (2) Write Back Protocol

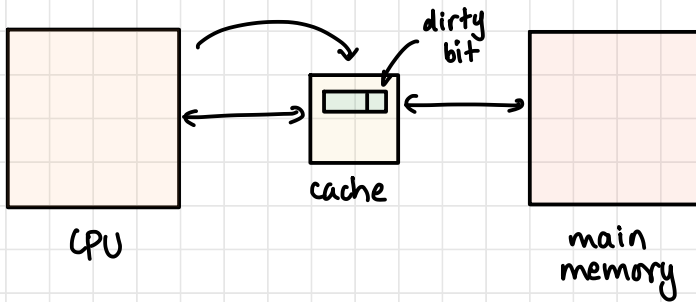
(1) Write Through Protocol

- both main memory & cache updated simultaneously
- ensures consistency
- increased latency



(2) Write Back Protocol

- change only cache
- dirty bit / modified bit flag in cache set to 1
- while sending out dirty block (victim block), main memory is updated



ADVANTAGE OF WRITE BACK

- Not all write operations need to access memory; lower latency
- Several writes in same cache block: force memory write only once at writeback time

Cache Miss

- Data not present in cache
- Miss penalty
- Latency: Time req. to retrieve first word of block
- Bandwidth: time req. to retrieve rest of block

read miss

- Load through/early restart - do not wait for entire block to be transferred

write miss

- Fetch and then overwrite in cache

(1) Write No Allocate

- Write directly to memory without affecting cache
- good if same location not needed soon
- valid bit set to 0

(2) Write Allocate

- load newly written data into cache
- easily accessible in cache
- same data needed again

Types of Misses

1) Compulsory Miss

- When cache initially empty, compulsory miss

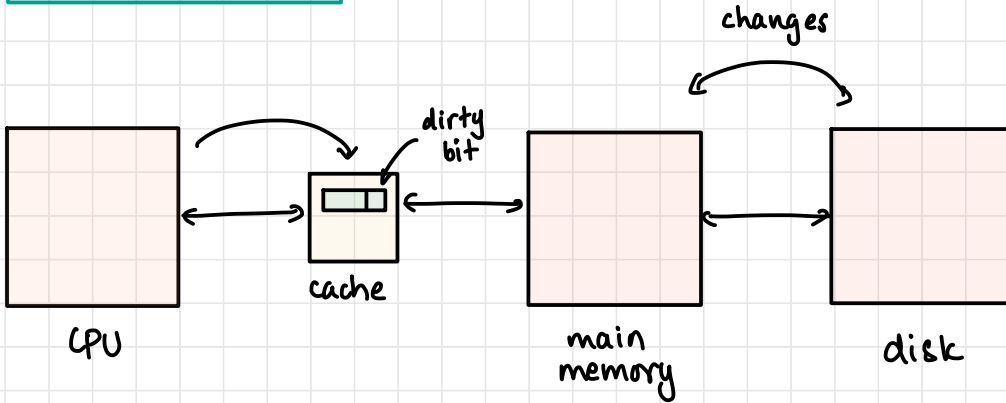
2) Capacity Miss

-

3) Conflicting Miss

- Set associative mapping
- Due to constraints, even if blocks empty

CACHE COHERENCE



- disk and memory data changes and writeback protocol used
- cache data might have also changed (dirty bit)
- copies of data different
- option: force writeback before main memory is updated from disk

Write Buffer

- processor waiting for writing into main memory: time consuming
- processor places write request into buffer, continues execution
- future access to data: access from buffer
- write through

Q: Consider DMC with 8 cache blocks. Memory block requests are 4, 3, 25, 8, 19, 6, 25, 8, 16, 35, 45, 22, 8, 3, 16, 25, 7, which memory blocks will be present in the cache at the end of the sequence? Also calculate hit/miss ratio.

0	8 16 8 16	8 miss → 8 hit → 16 miss → 8 miss → 16 miss
1	25	25 miss → 25 hit → 25 hit
2		
3	3 19 35 3	3 miss → 19 miss → 35 miss → 3 miss
4	4 45	4 miss → 45 miss
5		
6	6 22	6 miss → 22 miss
7	7	7 miss

hits = 3 misses = 14 hit ratio = 3/17

Q: Consider 2WSAM with 8 cache blocks. Memory block requests are 4, 3, 25, 8, 19, 6, 25, 8, 16, 35, 45, 22, 8, 3, 16, 25, 7, which memory blocks will be present in the cache at the end of the sequence? Also calculate hit/miss ratio. (LRU)

0	4 16	4 miss → 16 miss → 16 hit
	8	8 miss → 8 hit → 8 hit
1	25	25 miss → 25 hit → 25 hit
	45	45 miss
2	6	6 miss
	22	22 miss
3	3 35 7	3 miss → 35 miss → 7 miss
	19 3	19 miss → 3 miss

hits = 5 misses = 12 hit ratio = 5/17

Replacement Algorithms

- Direct mapped cache: memory block occupies fixed spot; replacement strategy trivial
- Fully associative, set associative need strategies

(1) Random

- replace random block from cache

(2) Least Recently Used (LRU)

- replace block that has not been used for longest time
- expensive
- slight randomness introduced for better performance

(3) First In First Out (FIFO)

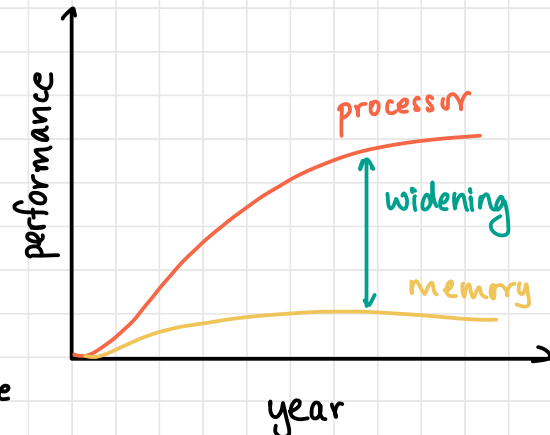
- evict block which has been in the cache longest

(4) Least Frequently Used (LFU)

- evict block that has been used least frequently

CACHE PERFORMANCE

- **Ideal:** all requests are hits
- **Realistic:** not 100%.
- How helpful cache is
- **Memory wall**
- **Locality of reference:** 10% of code takes up 90% of CPU time



$$\text{mem stall} = \text{IC} \times \frac{\text{mem accesses}}{\text{instruction}} \times \text{miss rate} \times \text{miss penalty}$$

CPI - Cycles per Instruction

$$\text{CPI} = \frac{\text{CPU clock cycles for program}}{\text{instruction count}}$$

Q: Assume we have a computer where $\text{CPI} = 1.0$ when all memory accesses hit the cache. Only data accesses are load and store, total about 50% of instructions. Miss penalty is 25 CC, miss rate 2%. How much faster would comp be if all instructions were cache hits?

$$\begin{aligned} \text{time CPU}_{\text{ideal}} &= (\text{CPU}_{\text{cc}} + \text{mem stall}) \times \text{clock cycle time} \\ &= (\text{CPI} \times \text{IC}) \times \text{clock cycle time} \\ &= \text{IC} \times \text{clock cycle time} \end{aligned}$$

$$\text{time CPU}_{\text{real}} = (\text{CPU}_{\text{cc}} + \text{mem stall}) \times \text{clock cycle time}$$

$$\text{mem stall} = \text{IC} \times \frac{\text{mem access}}{\text{instr}} \times \text{miss penalty} \times \text{miss rate}$$

$$= \text{IC} \times 1.5 \times 25 \times 0.02 = 0.75 \text{ IC}$$

$$\text{time CPU}_{\text{real}} = (1.75 \text{ IC}) \times \text{clock cycle time}$$

$$\therefore \text{Speedup} = 1.75 \text{ times}$$

8. Assume that the CPI for a computer is 1.0 and all memory accesses hit in the cache. If 30% of instructions are load/stores, miss penalty is 100 cycles, miss rate 5%, how much faster would the computer be if all instructions were cache hits?

CPU time with ideal cache = (CPU clock cycles + mem stall cycles) × clock cycle time

$$= (IC \times CPI + 0) \times \text{clock cycle time}$$

$$= IC \times \text{clock cycle time}$$

CPU time with imperfect cache = (CPU clock cycles + mem stall cycles) × clock cycle time

$$\text{mem stall cycles} = IC \times \frac{\text{mem accesses}}{\text{instructions}} \times \text{miss rate} \times \text{miss penalty}$$

$$= IC \times (1 + 0.3) \times 0.05 \times 100$$

$$= IC \times 6.5$$

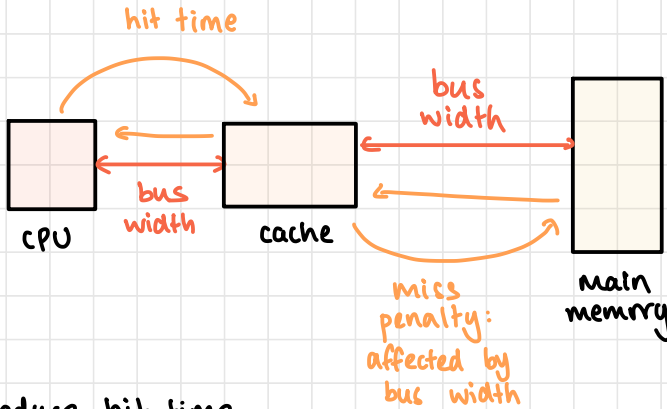
CPU time = (IC × 1 + IC × 6.5) × clock cycle time

$$= 7.5 \times IC \times \text{clock cycle time}$$

Speedup = 7.5 times

AVERAGE MEMORY ACCESS TIME

$$\text{AMAT} = \underbrace{\text{Hit time}}_{\text{constant}} + \text{miss rate} \times \text{miss penalty}_{\text{additional cc if miss incurred}}$$



- Reduce hit time
- Reduce miss rate
- Reduce miss penalty
- Miss penalty depends on bus width

Q: A certain processor uses fully associative cache of size 16 KB. Cache block size is 16 bytes. Assume byte addressable main memory, 32-bit addressing. Tag=? Index=?

$$\text{no. of blocks} = \frac{16 \times 2^{10}}{16} = 2^{10} \text{ blocks}$$

$$\text{words/block} = 16 \text{ words} = 2^4 \text{ words/block}$$

$$\text{tag} = 28 \text{ bits}$$

$$\text{index (block)} = 0 \text{ bits}$$

Q: The width of physical address is 40 bits. Width of tag field in 512 KB 8-way associative cache is?

Assume block size: 32 bytes

$$\text{no. of sets} = \frac{512 \times 2^{10}}{32 \times 8} = \frac{2^{19}}{2^8} = 2^{11} \text{ sets}$$

$$\text{size of block} = 32 \text{ bytes} = 2^5 \text{ bytes}$$

$$\therefore \text{tag bits} = 40 - (11 + 5) = 24 \text{ bits}$$

Q: Consider a 4WSA cache, 128 lines, line size 64 words, mem address 20 bits. Tag=? Line=? Word=?

$$\text{no. of sets} = \frac{128}{4} = 2^5$$

$$\text{set no} = 5 \text{ bits}$$

$$\text{line size} = 64 \text{ words} = 2^6 \text{ bytes}$$

$$\text{line bits} = 6$$

$$\therefore \text{tag bits} = 9$$

Q. Consider a fully associative mapped cache of size 16 KB with block size 256 bytes. The size of main memory is 128 KB. Find Number of bits in tag

$$\text{cache size} = 16 \text{ KB} = 2^4 \times 2^{10} = 2^{14} \text{ bytes}$$

$$\text{block size} = 256 \text{ bytes} = 2^8 \text{ bytes}$$

$$\text{main memory} = 128 \text{ KB} = 2^7 \times 2^{10} = 2^{17} \text{ bytes}$$

$$\therefore \text{addressing} = 17 \text{ bits}$$

fully associative



$$\text{block offset} = 8 \text{ bits}$$

$$\therefore \text{tag} = 9 \text{ bits}$$

8. Consider a 2-way set associative cache of size 16 KB with block size 256 bytes. The size of main memory is 128 KB. Find Number of bits in tag

$$\text{cache size} = 16 \text{ KB} = 2^{14} \text{ bytes}$$

$$\text{block size} = 2^8 \text{ bytes}$$

$$\text{set size} = 2^9 \text{ bytes}$$

$$\text{no of sets} = \frac{2^{14}}{2^9} = 2^5 \text{ bytes}$$

$$\therefore \text{set bits} = 5 \text{ bits}$$

$$\text{block offset} = 8 \text{ bits (size of block)}$$

$$\begin{aligned} \text{main memory} &= 128 \text{ KB} = 2^7 \times 2^{10} = 2^{17} \text{ bytes} \\ &= 17 \text{ bit addressing} \end{aligned}$$

$$\therefore \text{tag} = 17 - (8 + 5) = 4 \text{ bits}$$

