# MEMORY SYSTEM

29/10

→ DRAM

- For modern memory, one memory access takes 150 - 200 cycles :'(
  (not 1 as before)
  - Not counting multiple memory accesses (just one, to fetch ir), each instruction will take ~200 cycles :(
  - Pipelining makes no sense, as a fetched instruction will take only #stages cycles more to complete.

- One possible solution :- 'Cache'
  Between processor and memory, add a very small, but fast block of memory. This has, say, only 8 bytes (worth 2 instructions) and is made from more expensive SRAM. SRAM memory is accessible in single cycle (must fetch consecutively written instructions)

  - Fetch 8 bytes (2 instructions) from memory into SRAM (~200 cycles
    Processor fetches $I_1$ ($200^{st}$ cycle)
    $I_2$ ($201^{nd}$ cycle)  ---- pipelining
    $I_3$ ($402^{rd}$ cycle)
    $I_4$ ($402^{th}$ cycle) ---- pipelining

  - Performance almost doubles by introducing SRAM, as at least 2 instructions are being pipelined.

eg | Eg Program                           for i in range (n):
                                              $I_1$
Assume instruction                            $I_2$
do not access memory                          $I_3$
                                              $I_4$

If SRAM is 16 bytes, 4 instructions will be fetched in 200 cycles and stored in SRAM.
For every subsequent iteration, we don't need to fetch again and instructions can be pipelined as before
                    CPI → 1          :) as we are using only 4 instruction multiple times

eg for i in range (100):-
    if i % 2 == 0 :-        $I_1$ $I_2$ $I_3$ $I_4$
    else                    $I_5$ $I_6$ $I_7$ $I_8$

- With 16 bytes of SRAM, we will have to fetch and replace the four instructions again and again.

Total no. of cycles = $(200 + 1 + 1 + 1) \times 100$       ☹

- With 32 bytes of SRAM, we can get CPI → 1, but only if all 8 instructions are consecutively written. Otherwise same as before.

$(200 + 1 + 1 + 1 \cancel{+ 1 + 1}) \times \cancel{50} + (1 + 1 + 1 + 1 \cancel{+ 1 + 1 + 1 + 1}) \times 49$

- With 2 SRAM chunks of 16 bytes each, we can get CPI → 1. The instructions in each individual chunk must be consecutive.

$(200 + 1 + 1 + 1 + 200 + 1 + 1 + 1) + (1 + 1 + 1 + 1 \cancel{+ 1 + 1 + 1 + 1}) \times 48$

• The SRAM chunks are hidden from programmer and known as 'Cache'.

20/10

→ Bandwidth of Memory :-

$$\frac{Data}{sec} = \frac{Data}{Instruction} \times \frac{Instruction}{Cycle} \times \frac{Cycle}{Sec}$$

$$= \left(\frac{\# Bytes}{Instruction} + \frac{\# Bytes\ of\ data}{Instruction}\right) \times IPC \times freq$$

$$= (4B + 4B) \times 1 \times 3 \times 10^9$$

$$\approx 24\ GBps.$$

→ Final CPI depends on no. of 'stalls' required by the program
- Instructions may stall, for reading from or writing to memory, for 200 cycles.
- Stalling also happens to fetch instructions not present in cache

• Instructions can be 'hit' or 'miss'

present in cache                          not present
No penalty                                200 cycle penalty.

• No. of stalls = No. of misses × penalty.

$$= \text{\# Instruction} \times \frac{\text{\#miss}}{\text{Instruction}} \times \text{penalty}$$

$$= \text{\# Inst} \times \frac{\text{\#Memory references}}{\text{Inst}} \times \frac{\text{\# miss}}{\text{\# Mem ref}} \times \text{penalty}$$

$$= \underset{(IC)}{\text{Inst Count}} \times \frac{\text{\# Mem ref}}{\text{Inst}} \times \underbrace{\text{Miss rate} \times \text{penalty}}_{\text{In the hands of memory designee}}$$

→ Miss Rate reduction

1  Increase no. of blocks in SRAM

2  Replacing contents of cache :-
Preferentially, do not remove instructions that are consecutively likely to follow current instruction (spatial likeliness) or that are likely to be executed again at some later point in time (temporal likeliness)

- Increasing the number of blocks in SRAM beyond a certain point will increase access time from 1 cycle to 2 cycles. ☺

→ Questions to be answered

1. In which block of cache to place?
2. How to identify where required instruction is placed in cache?
3. Which & cache instruction to replace when overwriting?
4. When to write ~~to cache~~ & data (obtained from processor) from cache to memory?

Cache                                          Memory

A1, A2  TAG

8 blocks of 32 B each

                                    1024 blocks of 32B each

In memory, every byte can be addressed as

| Block No. | offset |
|-----------|--------|

10 bits (mem)   5 bits
3 bits (cache)

cache

- We will only be transferring blocks as a whole from memory to cache. After transferring, offset part of each byte remains same. Block no. will be changed accordingly

- We maintain a mapping of which block no. of memory is ~~map~~ transferred to which block no. of cache.

- Simplest Mapping = 'Direct Mapping'

(Mem Block No.) % 8 = (Cache Block No.)

eg- Mem Block 45 is transferred to Cache block 5.     P.T.O.

i.c - ~~Ta~~ ~~Mem bloc~~

(Cache block no.) = Last 3 bits of (Mem block no.)

∴ (Cache address) = Last 8 bits of (Mem address)

- The first seven bits of Mem address (which were lost in transferring to cache) will be stored in an adjacent array to cache. (Tag)
  - TAG is an 8-element array of 7-bit elements each (one element for each block)

T