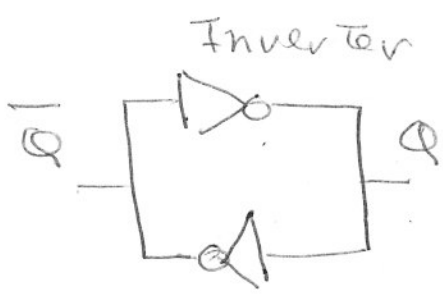
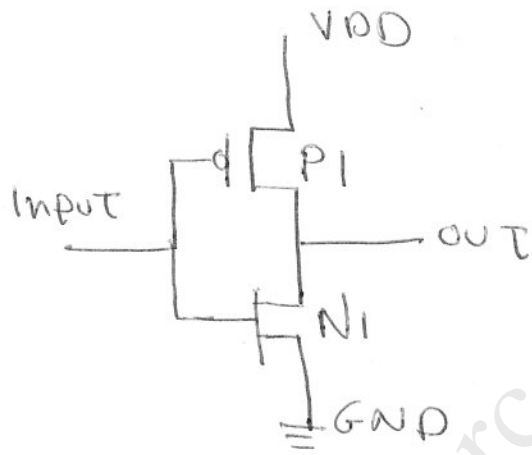


SRAM Organization

Latch



Inverter



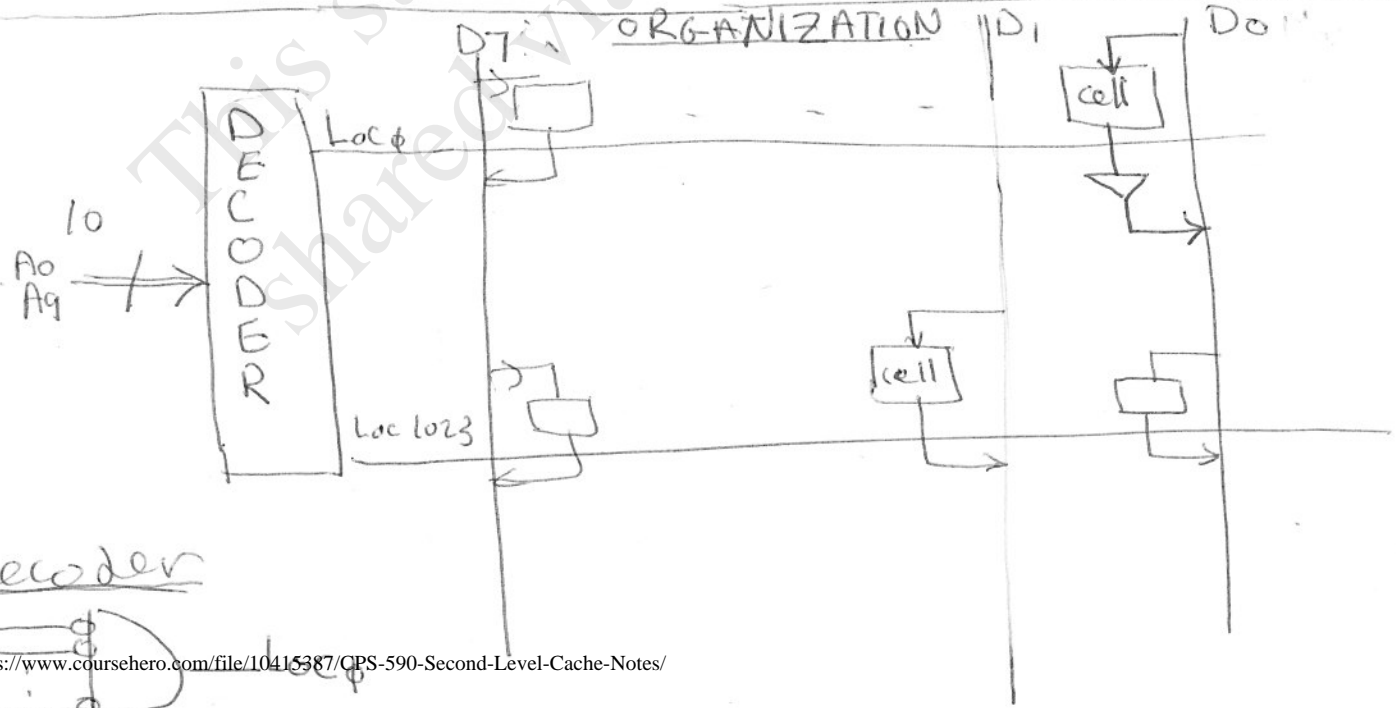
Input = 0

P1 conducts, N1 off
OUT = VDD = 1

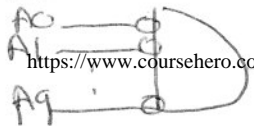
Input = 1

N1 conducts, P1 off
OUT = GND = 0

ORGANIZATION



Decoder



Method 4: Non Blocking Caches

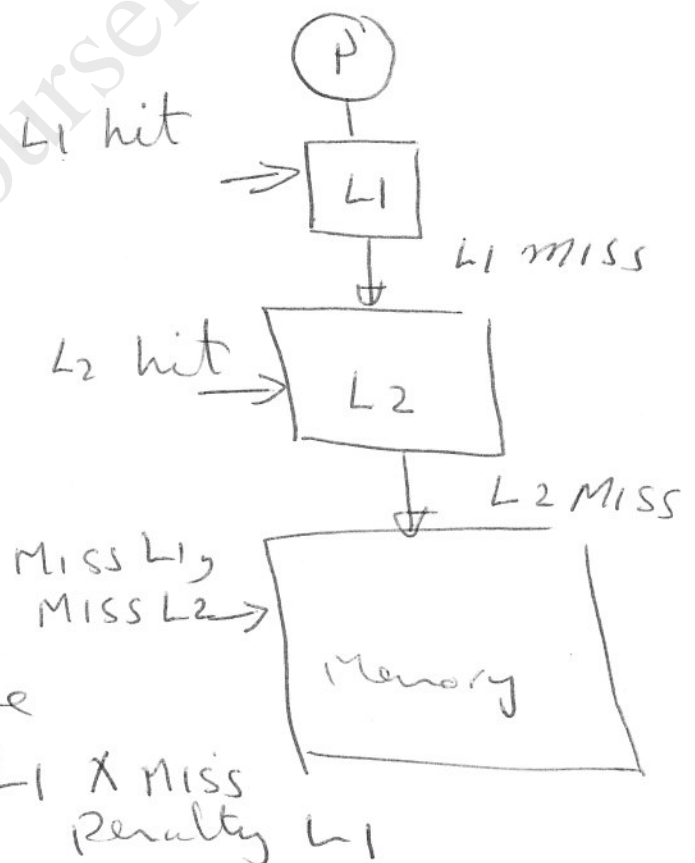
For pipelined machines that allows out of order execution and completion, CPU could continue fetching instructions from instruction cache while waiting for data cache miss. It also allows data cache to supply data while waiting for a miss data to arrive from memory.

- need scoreboard or Tomasulo control.

Method 5: Second Level Caches

on a first level cache miss, can use a faster level 2 cache and not the memory.

only misses from L1 could be found in L2.



Average memory access time

$$= \text{Hit time } L1 + \text{Miss } L1 \times \text{Miss Penalty } L1$$

$$\text{Miss Penalty } L1 = \text{Hit time } L2 + \text{Miss } L2 \times \text{Miss Penalty } L2$$

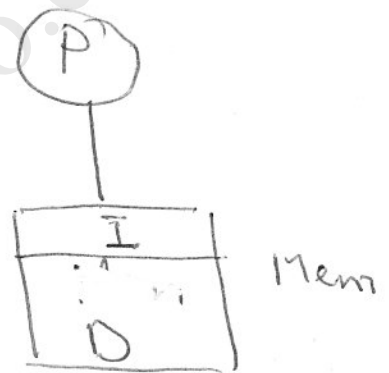
Main Memory

- Introduction (CPU & Mem)
- SRAM (Organization, Timing, Expansion)
- DRAM (Organization, Timing, Access modes, Interleaving, Expansion)

Introduction

Processor & Memory model

- Instructions stored in Memory.
Processor uses Fetch
- Data stored in Mem
Processor uses load/store



Each single instruction might access memory twice (Fetch, and load/store)

Example $lw R1, 100(R2)$
 $R1 \leftarrow M[100 + R2]$

- Using RAM = Random Access memory.
All accesses have same access time
- Computer system performance depends on memory speed (or CPU??)