

13.3 Main memory updating

(datum)

DATUM

Cache coherency

→ reads/writes

→ Statistics: Reads are dominating

- Instructions are read.

- Data are read/written

- Writes represent $\leq 10\%$ of mem. accesses.

Andahl's Law: make the common part faster.

We cannot neglect writes!

a) Write Through: write in the cache and MM at the same time

b) Write Back: write in the cache only

- Update in MM when the block is replaced

- 1 bit dirty/clean

What happens in case of a Write miss?

In the case of read misses ~~not always~~ allocate!

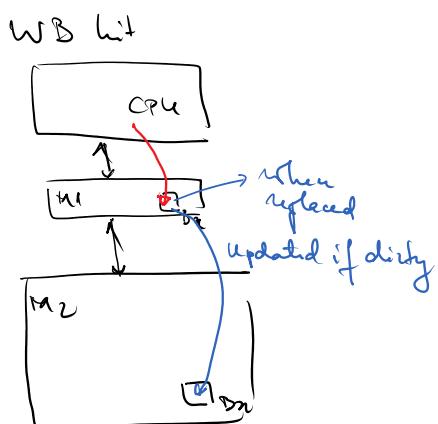
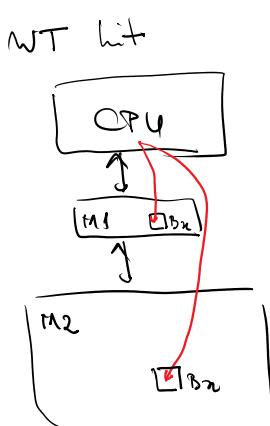
α) Write Allocate

β) Write No-Allocate

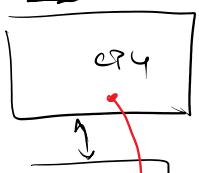
2 Natural policies

WT + WNA + Write in MM

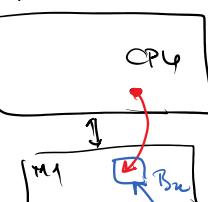
WB + WA + Write in Cache



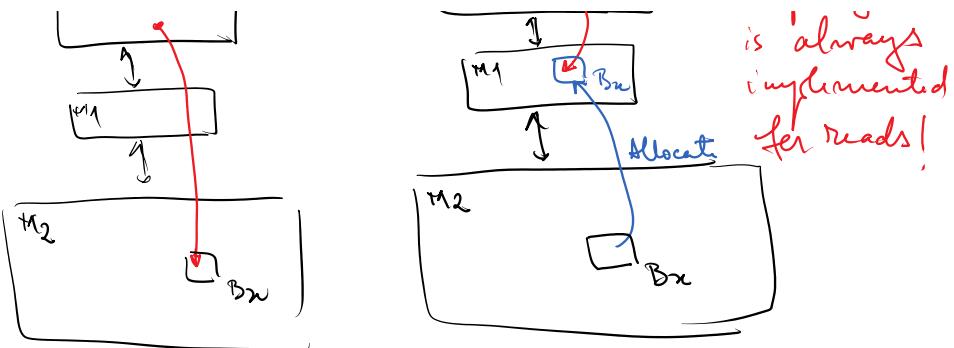
WT + WNA (Miss)



WB + WA



This policy is always implemented



3.4 Memory hierarchy performance

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

$$\begin{aligned} \text{Memory stall clock cycles} &= \frac{\text{Reads}}{\text{Program}} \times \text{Read Miss Rate} \times \text{Read Penalty} + \\ &+ \frac{\text{Writes}}{\text{Program}} \times \text{Write Miss Rate} \times \text{Write Miss Penalty} \end{aligned}$$

$$\text{Memory stall clock cycles} = \frac{\text{Memory accesses}}{\text{Program}} \times \text{Miss Rate} \times \text{Miss Penalty}$$

Misses / Instruction

$$\text{CPU time} = IC \times (\text{CPU execution} + \frac{\text{Memory accesses}}{\text{Instruction}} \times \text{Miss Rate} \times \text{Miss Penalty}) \times \text{Clock cycle time}$$

Example 1 CPU execution = 8.5 c.c.

Memory references / instruction = 3
Miss Rate = 11% Miss Penalty = 6 c.c.

Influence of cache imperfection as %?

$$\text{CPU time perfect} = IC \times 8.5 \times \text{Clock cycle time}$$

$$\begin{aligned} \text{CPU time} &= IC \times (8.5 + 3 \times 0.11 \times 6) \times \text{Clock cycle time} \\ &\approx IC \times 10.5 \times \text{Clock cycle time} \end{aligned}$$

$$\frac{2 \text{ c.c. perf.} - 8.5 \text{ c.c. perf.}}{2 \text{ c.c. perf.}} \times 100$$

$$x = \frac{200}{8.5} = 23.5\%$$

Example 2

- 64KiB cache

- Clock cycle time = 200 ns

Memory accesses = 1.3
Instruction

- Clock cycles / instruction for a perfect cache = 1.5 c.c.

- For S.A. \Rightarrow clock frequency degradation with 8.5%

(A) DM cache Miss Rate = 3.5%

(B) 2-way SA Miss Rate = 3%

AMAT = Average Memory Access Time

• Miss Penalty for both solutions = 200 ns

miss cost = 2.210
③ 2-way SA miss Rate = 3%

Which solution is better?

AMAT = Average Memory Access Time
• Miss Penalty for both solutions = 200 ns
↳ usually is expressed as c.c.

$$\text{AMAT} = \underbrace{\text{tac} \times \text{Hit Rate}}_{1 \text{c.c.}} + (\underbrace{\text{tac} + \text{Miss Penalty}}_{\text{Time!}}) \times \text{Miss Rate} = \\ = \text{tac} + \text{Miss Penalty} \times \text{Miss Rate}$$

$$\text{AMAT}_A = 20 \text{ ns} + 0.039 \times 200 \text{ ns} = 27.8 \text{ ns}$$

$$\text{AMAT}_B = 21.7 \text{ ns} + 0.03 \times 200 \text{ ns} = 27.7 \text{ ns}$$

$$\text{Clock cycle time} = 1.085 \times 20 \text{ ns} = 21.7 \text{ ns}$$

2-way SA is better

Miss Penalty (time)

$$\text{CPU time}_A = \text{IC} \times \left(\text{CPI}_{\text{ideal}} + \frac{\text{Memory acc}}{\text{Hit}} \times \text{Miss Rate} \times \text{Miss Penalty} \times \text{Clock cycle time} \right) =$$

$$= \text{IC} \times \left(\text{CPI}_{\text{ideal}} \times \text{Clock cycle time} + \frac{\text{Memory acc}}{\text{Hit}} \times \text{Miss Rate} \times \text{Time Miss Penalty} \right) =$$

$$= \text{IC} \times (1.5 \times 20 \text{ ns} + 1.3 \times 0.039 \times 200 \text{ ns}) = [40.14 \text{ ns}] \times \text{IC}$$

$$\text{CPU time}_B = \text{IC} \times (1.5 \times 21.7 \text{ ns} + 1.3 \times 0.03 \times 200 \text{ ns}) = (7.8 + 32.55) \text{ ns} \times \text{IC} \\ = [40.35 \text{ ns}] \times \text{IC}$$

DMA is better!

Cacher performance is decisively influenced by misses!

What are the sources of misses?

- Compulsory → empty (cold start) misses ↗ bigger block size
- Conflict → DM (Collision misses) ↗ higher SA (maybe FA) ↗ degrees the tac
- Capacity → Too many references ↗ Compiler optimization (loop tiling)

3.5 Cache performance optimization

- Compiler optimization (Software)
- Hardware
 - ① - bigger block size
 - ② - higher associativity
 - ③ - Victim caches
 - ④ - Pseudo-associative caches

