



RECAP

Practical Caches

- Fully associative: great for temporal locality, but costly
- Direct Mapped: bad for TL but cheap
 - Compromise: set associative
- LRU: best cache replacement algorithm for temporal locality
 - Expensive tho
 - Cyclic, random... alternatives

Cache Control Bits

- Cache INITIALISATION?
 - need VALID bit for each entry to indicate meaningful data
- Dirty bit: \leftarrow cache value written BACK to memory
 - write through write back

Exploiting Spatial Locality

① complete!

- Under cache line stores

Multiple line Block Size

① complete!

Effect of Line Size

- Larger line \longrightarrow fewer misses ① more CTR

① Data unused

① Larger RAM accesses take longer

① Displaces other useful data

] if line becomes too big

temporal but more spatial

→ too small \longrightarrow less spatial

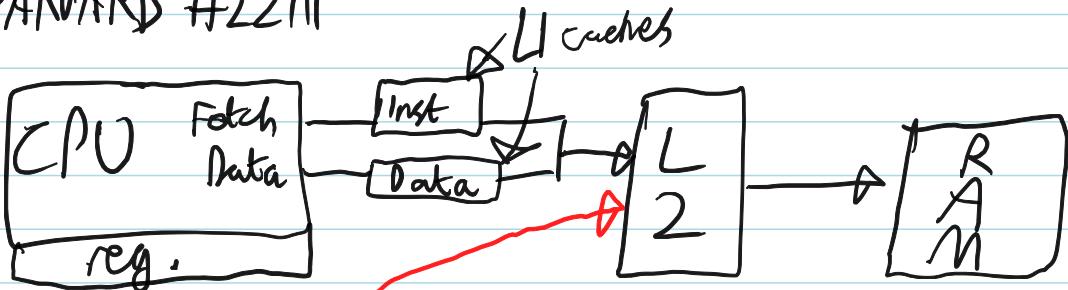
but more temporal

Separate I&D caches

! completed!

- Fetch per instruction
- Data fetch every 3 instructions

#HARVARD #2211



MULTIPLE level caches

- Bigger caches → Higher hit rates BUT slower speeds!
- Solution: decrease number of cycles per
- L2 cache bigger (slower than L1)
→ BUT L2 still faster than RAM
- If missrate in L2 and L1 were only 2% each:
→ memory access rate = 0.04%
- inclusive cache: can't have data in level ≥ 2 if data isn't at lower level ≤ 1
- L2 usually shared between I and D caches

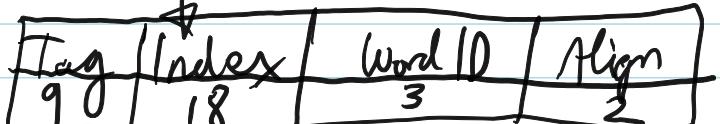
8 MB cache → 2^{23} bytes

32 bit 4 byte words, each line of entry = 8 words

∴ line = 32 bytes = 2^5 bytes

∴ entries $\approx 2^{23-5} = 2^{18}$ entries (≈ 256 K entries)

- looking at slide:



size and format of address

$\rightarrow 2 \times ((t \text{ number of cycles}))$

You need to determine if INITIALLY the cache has a miss or hit

$$h_i = \text{hit rate at level } i \\ m_i = 100\% - h_i$$

$t_i = \text{technology based access time}$
 $T_i = \underline{\text{perceived}}$ at level i
what we calculate

$$T_i = t_i \cdot h_i + m_i \cdot (t_i + T_{i+1}) \\ = t_i + m_i T_{(i+1)} \quad (\text{You do the math})$$

