

Cache controllers

My CPU when the L1 cache misses



This little maneuver is gonna cost us 3 nanoseconds

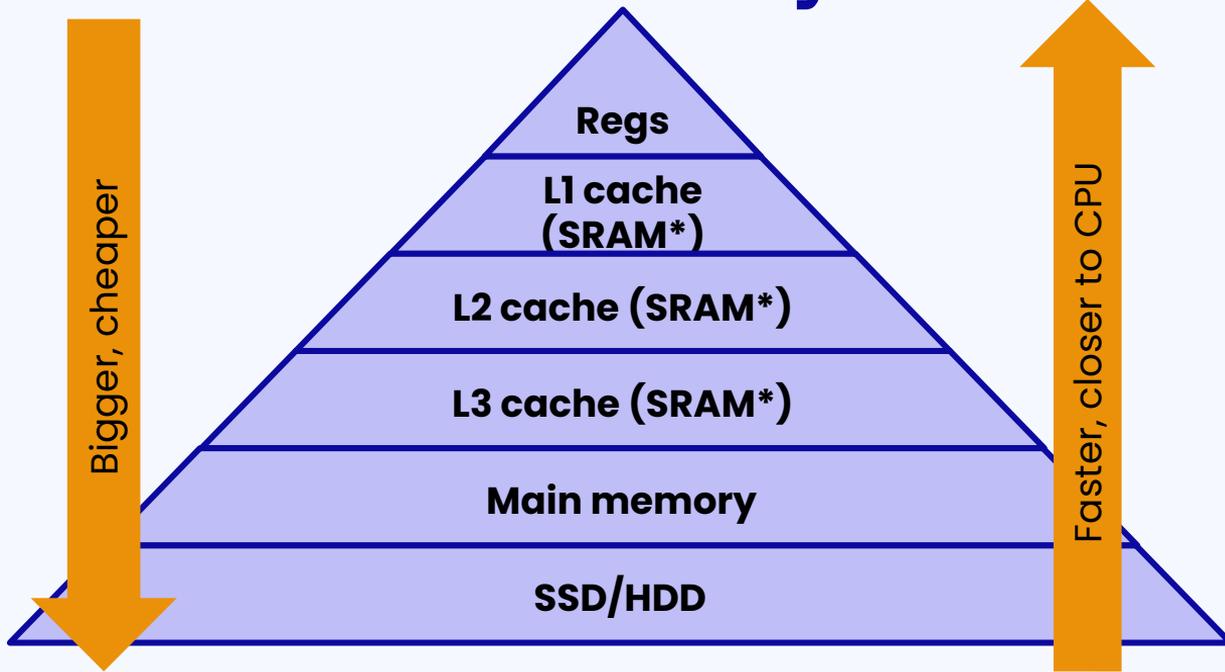


Async (Paris Arc 🇫🇷)

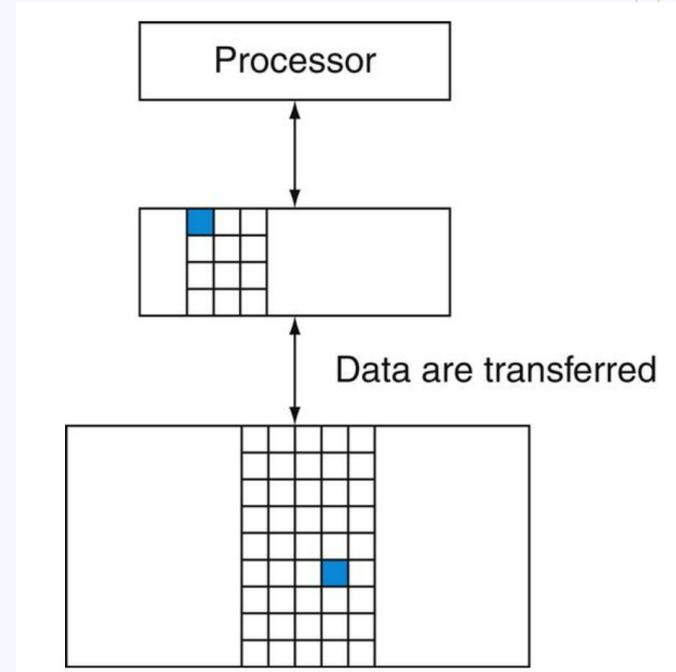
@0xAsync

No mom it's not a "messy pile of clothes on my chair" it's an L1 cache for fast random access to my frequently used clothes in $O(1)$ time. It needs to be big to avoid expensive cache misses (looking in my closet). I NEED to be minimizing latency, this is important to me. Please.

Mem. hierarchy review



Each level stores (**caches**) a subset of the one below it, for faster access of specific data



P&H 5.2

Cache design questions

- How do we decide what goes where in a cache?
- How do we decide what data to evict when our cache gets full?
- How do we maintain consistent copies of data across the hierarchy?
- What does the cache controller need to keep track of?

Later:

- How do we build an efficient* memory hierarchy?
- How do we manage shared caches?
- Do caches expose security vulnerabilities?

Block size: 1 word (4 bytes)

Cache size: 1 KB ($256/2^8$ blocks)

Our example cache

Our program

Assume register a0 holds the value 0x20000004

```
lb t0, 1(a0)
```

```
lw t1, 4(a0)
```

```
lhu t2, 2(a0)
```

```
lw t0, 1024(a0)
```

```
addi t1, t1, 1
```

```
sw t1, 4(a0)
```

```
sh t1, 8(a0)
```

Terminology

Block: minimum unit of information that can be present/not present in a cache

Index: how we refer to locations in the cache; each block has a unique index and each (block-aligned) memory address maps to a unique index (using the lower bits of the block address)

Valid bit: indicates whether data has been pulled in to that block of the cache

Tag: the upper bits of an address, used to uniquely identify which data is in the cache

Write through vs write back

Write through: every time data is changed in cache, change is done to lower level in hierarchy

Pro: We always have a consistent view of data across levels

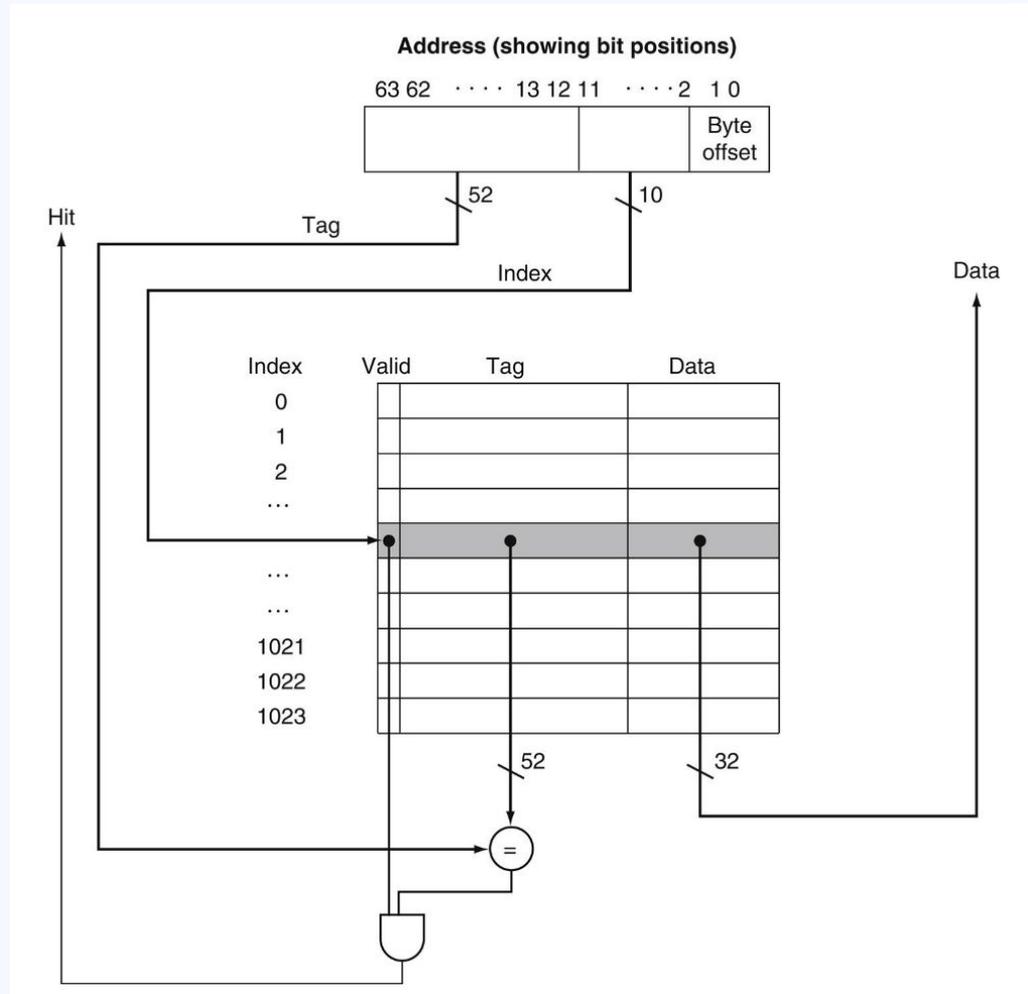
Con: Every write takes a long time

Write back: changes to lower level in hierarchy are only done when data is evicted from cache

Pro: Fast

Con: Consistency/needs a dirty bit

Cache controller



P&H fig. 5.10