

How do caches even work?

Key terms:

Temporal Locality:

- data that was recently accessed is likely to be accessed again in the near future.

Spatial Locality:

- Data that is physically close in memory tends to be accessed around the same time.

Cache memory structure

- cache organizes data into LINES, smallest chunk of memory a cache can transfer

TAG DATA

Lines are $32/64/128$ bytes

- data has tags with a number representing where the data came from

- then, its divided into SETS. Each SET has multiple cache lines



n-way set associative cache

PROJECT 3: CACHE

Tag Set index Offset

cache line

- offset bits are $\log_2(\text{block size})$

- set index is $\log_2(\# \text{ sets})$

- tag bits are whatever's left

tag bits = address width - offset bits - index bits

1 or 4 or n-way

- Cache is a 2D array

- Cacheline cache [numsets] [numways]

$$(\text{total cache size})^* = (\# \text{ sets})^* \cdot (\# \text{ ways})^* \cdot \text{block size}^*$$

or the other way, since we know cache size at execution

$$\# \text{ sets} = \frac{(\text{total cache size})}{(\text{associativity} \cdot \text{block size})}$$

instructions are fixed at 8 bytes.

num lines = 64
block size = 16 bytes

$$\text{ways} = \frac{\text{total lines}}{\# \text{ sets}} \rightarrow \text{we know sets now and lines}$$

no fractional ways, check for that

$$\text{ways} = \frac{64}{64 \text{ sets}} \rightarrow \text{direct mapped, 1 way}$$

$$\text{ways} = \frac{64}{1 \text{ set}} \rightarrow \text{fully associative, 64 ways}$$

Addr width

$$\text{offset bits: } \log_2(16) = 4 \text{ bits}$$

$$\text{set index: } \log_2(\# \text{ sets}) \rightarrow \text{will be known at runtime}$$

$$\text{tag bits} = 32 - (4 + \log_2(\# \text{ sets}))$$

So direct mapped:
64 cache lines
16 byte blocks

64 "slots"

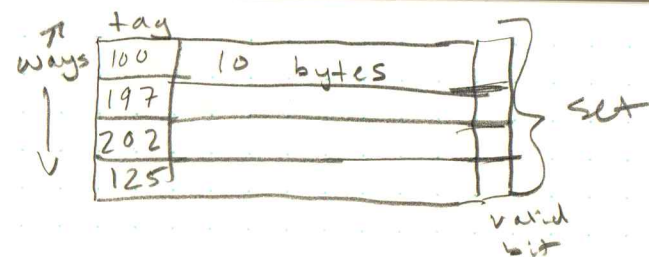
each slot has:

Valid bit 1 bit
tag 22 bits
data block 128 bits (16 bytes)

addresses are split as
[31...10] [9...4] [3...0]
tag index offset
(22b) (6b) (4b)

offset picks 1 of 16 blocks
index picks which ~~way~~
tag must match

bytes
cache[0][b] → {valid, tag(22b), data[0...15]}
⋮
cache[63] → {valid, tag(22b), data[0...15]}



TAG WAY
125 28 6

TAG → Is my line the one you want?

INDEX → Which row of the 2D array do we look in?

WAY → Which column in that row has the winner?

OFFSET → Which Byte/Word inside the 16-byte line?

Important to note: The cache is a 2D array of Cache Lines,
where the Size of the total cache is $\frac{1}{2}$

↳ [Number of sets] [Associativity]

keeping track of LRU ~~array~~ (counter, update when last used)