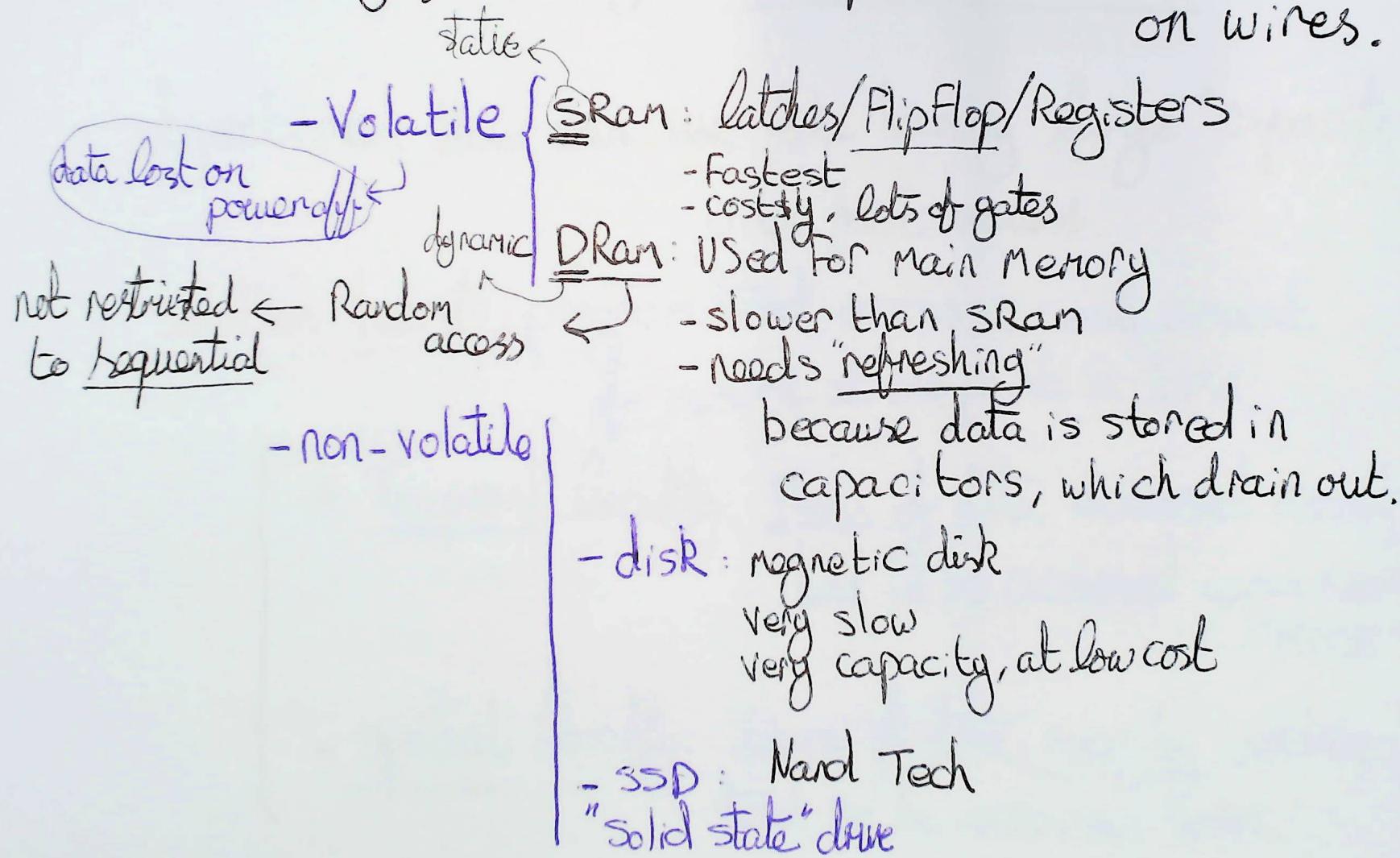


# chapter 5: The Memory

## Memory Technology

Latency of access  $\propto$  Capacitance  $\times$  Resistance  
on wires.



## Main feature of memory Tech

Larger "Capacity" is slower to access

smaller "capacity" is faster to access

objective: How can we fake having large "capacity"  
with fast access.

solution: heuristic program tend access a small amount  
of memory at one point in time.

- { - Temporal Locality: Items of data accessed recently  
tend to be accessed again soon.  
"reuse"
- spatial Locality: items of data near by (address)  
tend to be accessed soon.

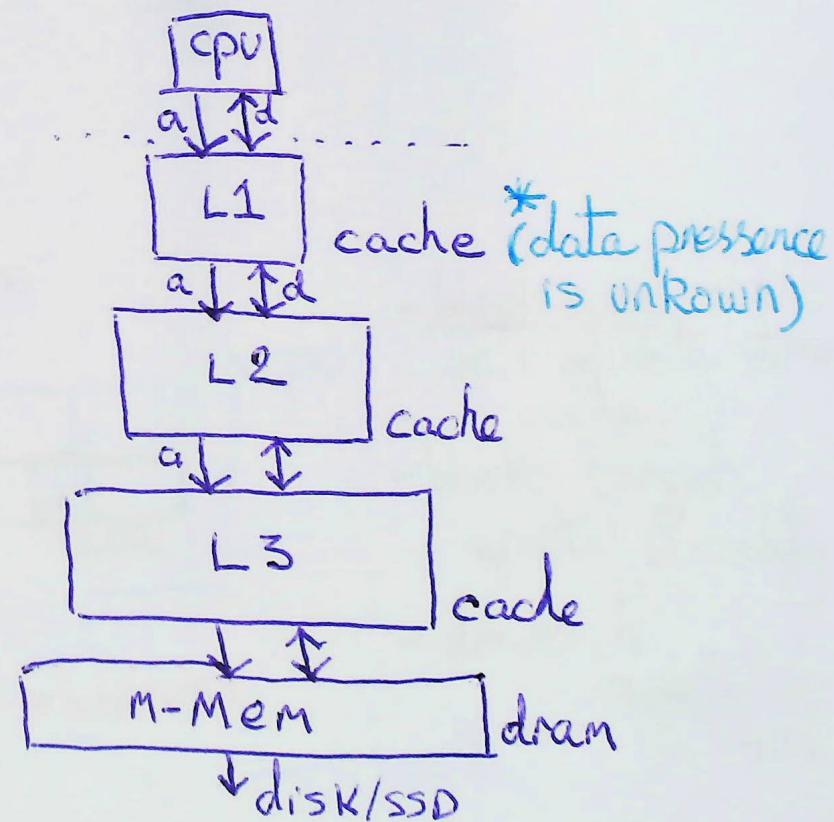
# How to take advantage of data locality?

## Memory Hierarchy

- Copy and Reep recently accessed and nearby data in memory that is smaller - closer to the processor.

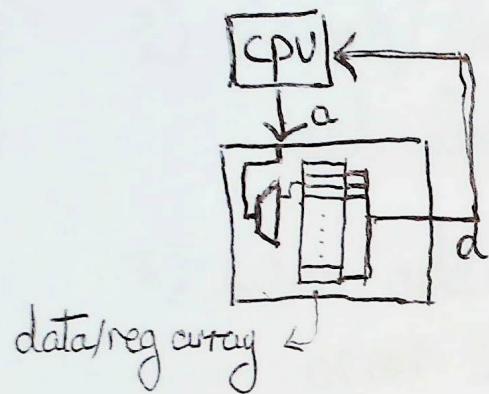
on average data appears to be accessed fast, but with a high total capacity to the memory.

③ How do levels of memory tell if they have a certain address?

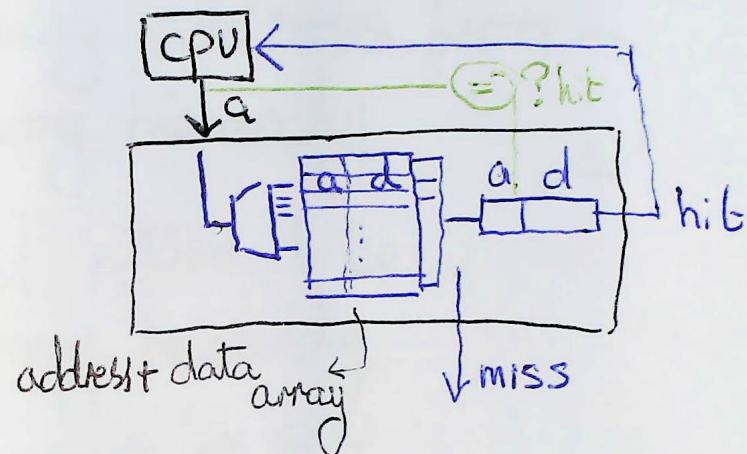


# Inside memory design

## Simple Ram

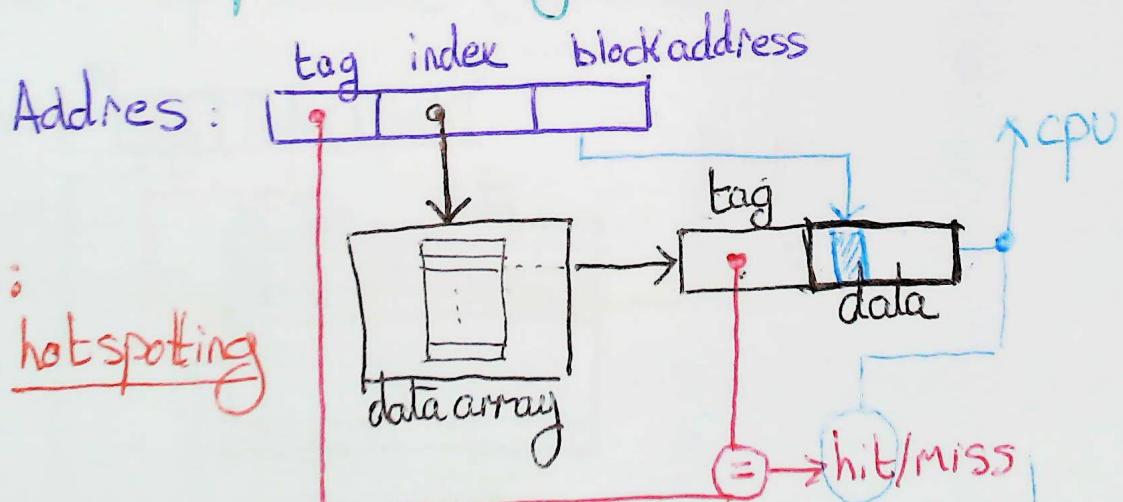


## Simple cache



## Direct Map Caching

problem:  
index hotspotting



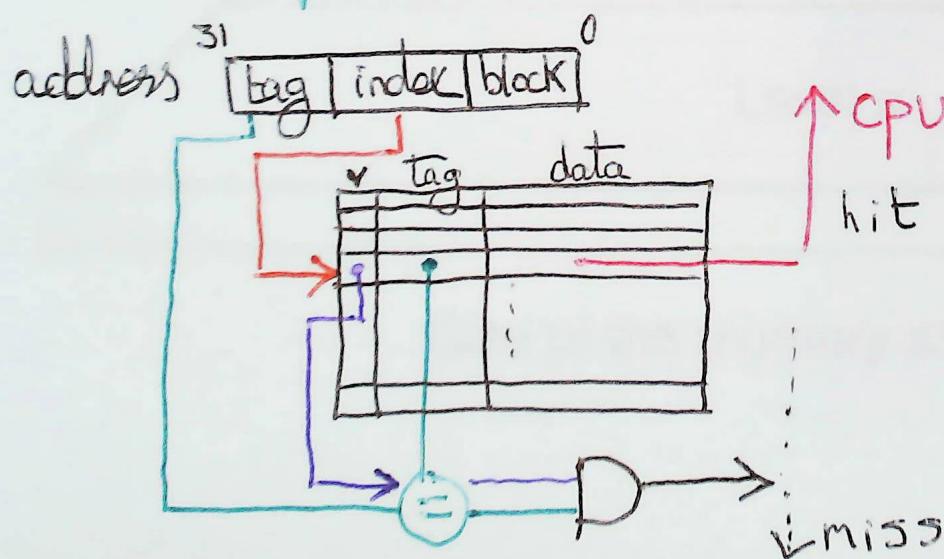
- block size is the unit of data stored in cache.
- blocks target spatial locality.
- downside: overfetch on miss

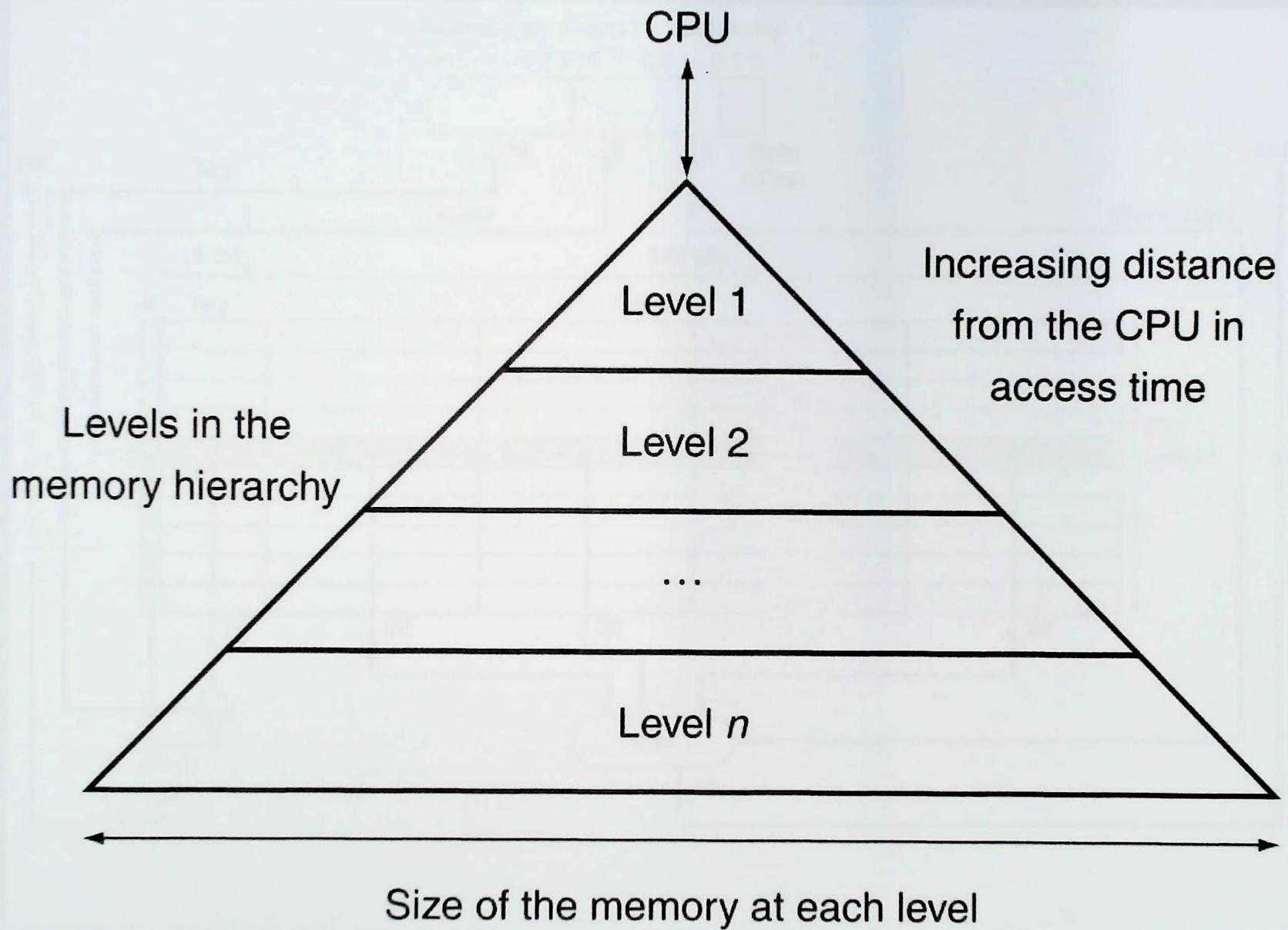
miss handle

## What happens on a Cache miss?

- 1) stall the cpu
- 2) fetch/lookup missing address from next level of memory hierarchy. ↗
- 3) when address is found, return it to prior level of memory hierarchy.  
→ to CPU
- 4) unstall the cpu

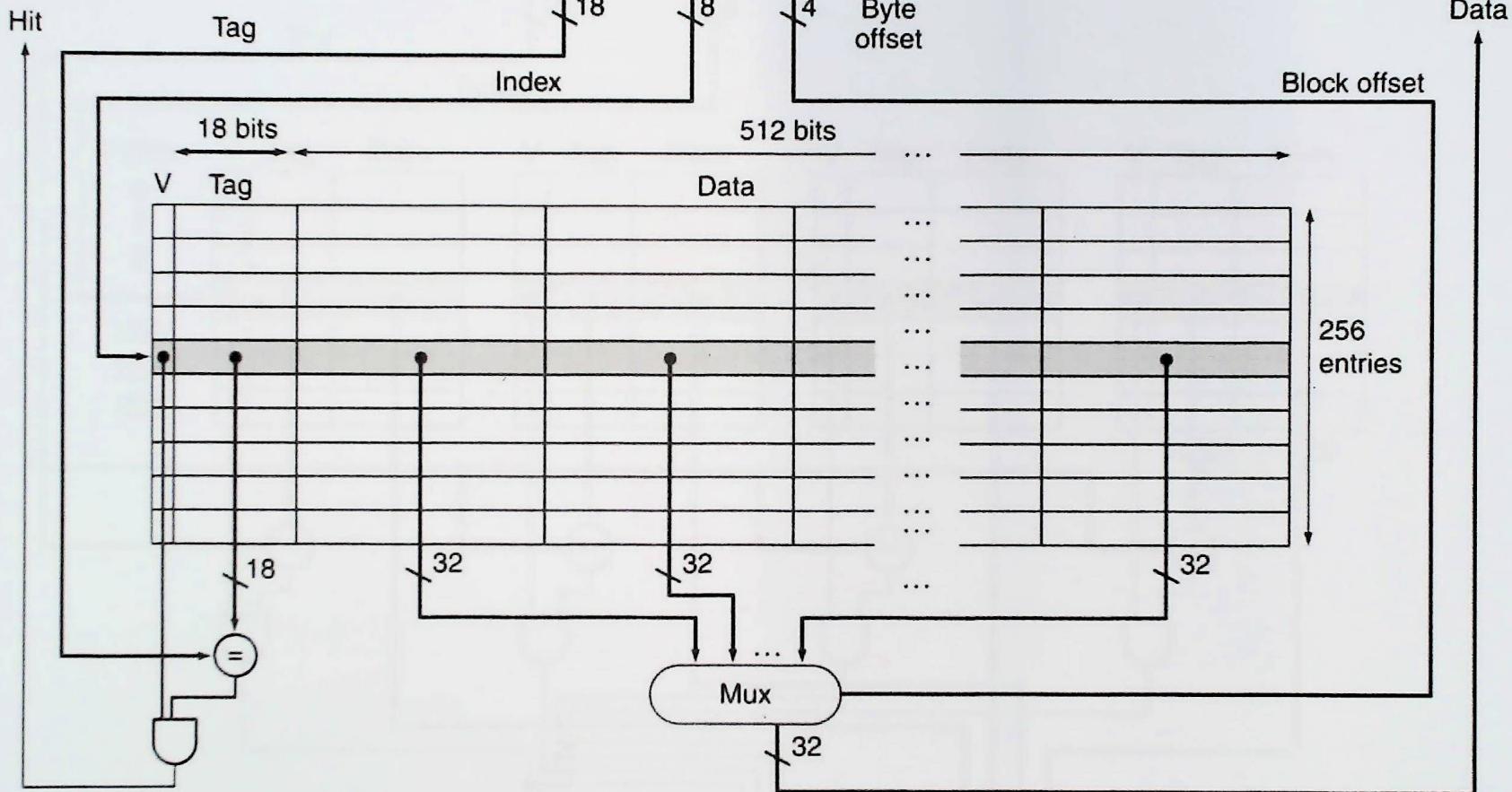
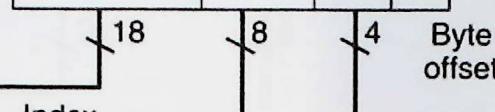
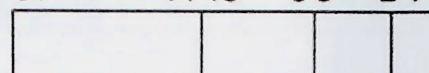
## Direct Map Cache (detailed)

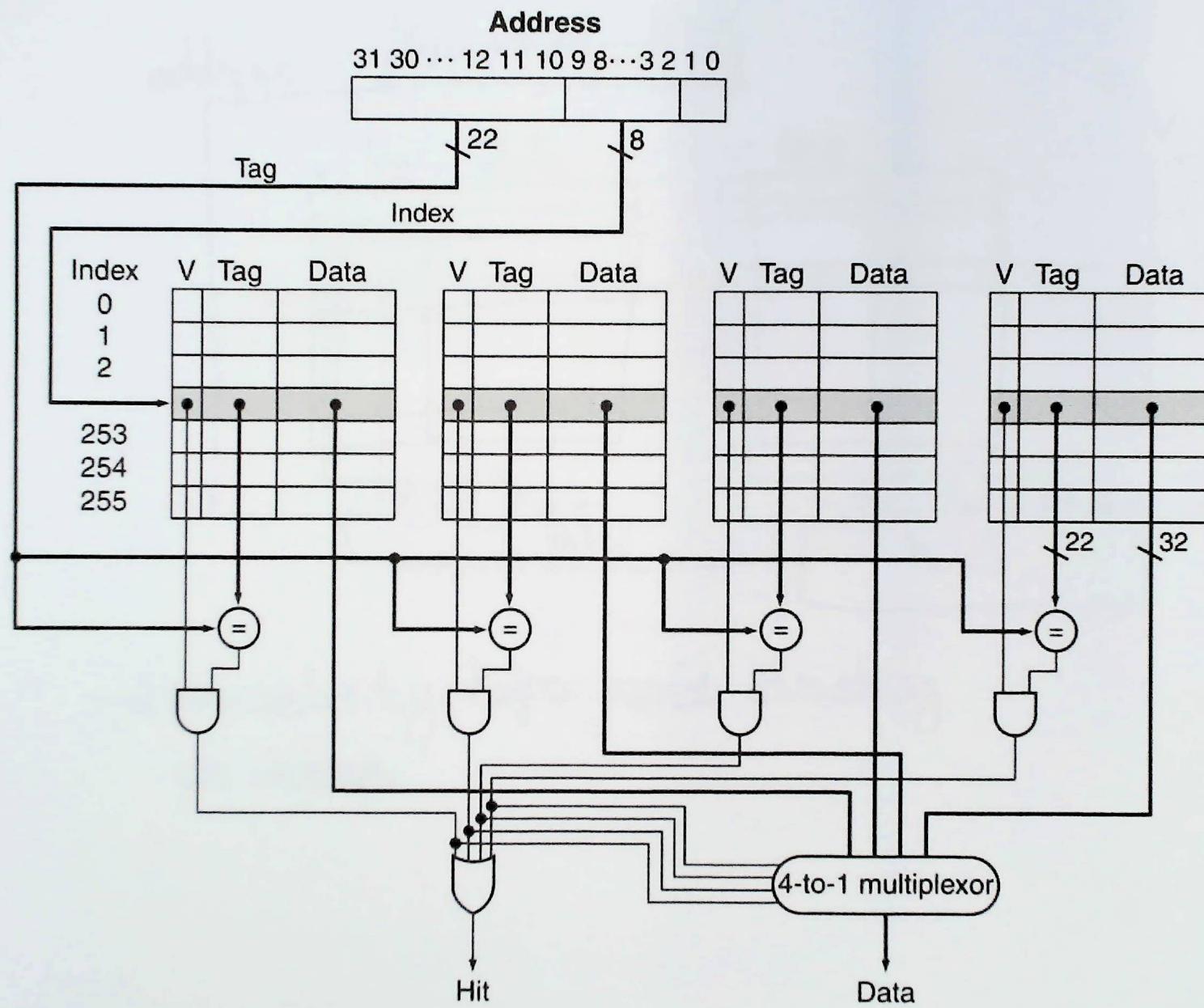




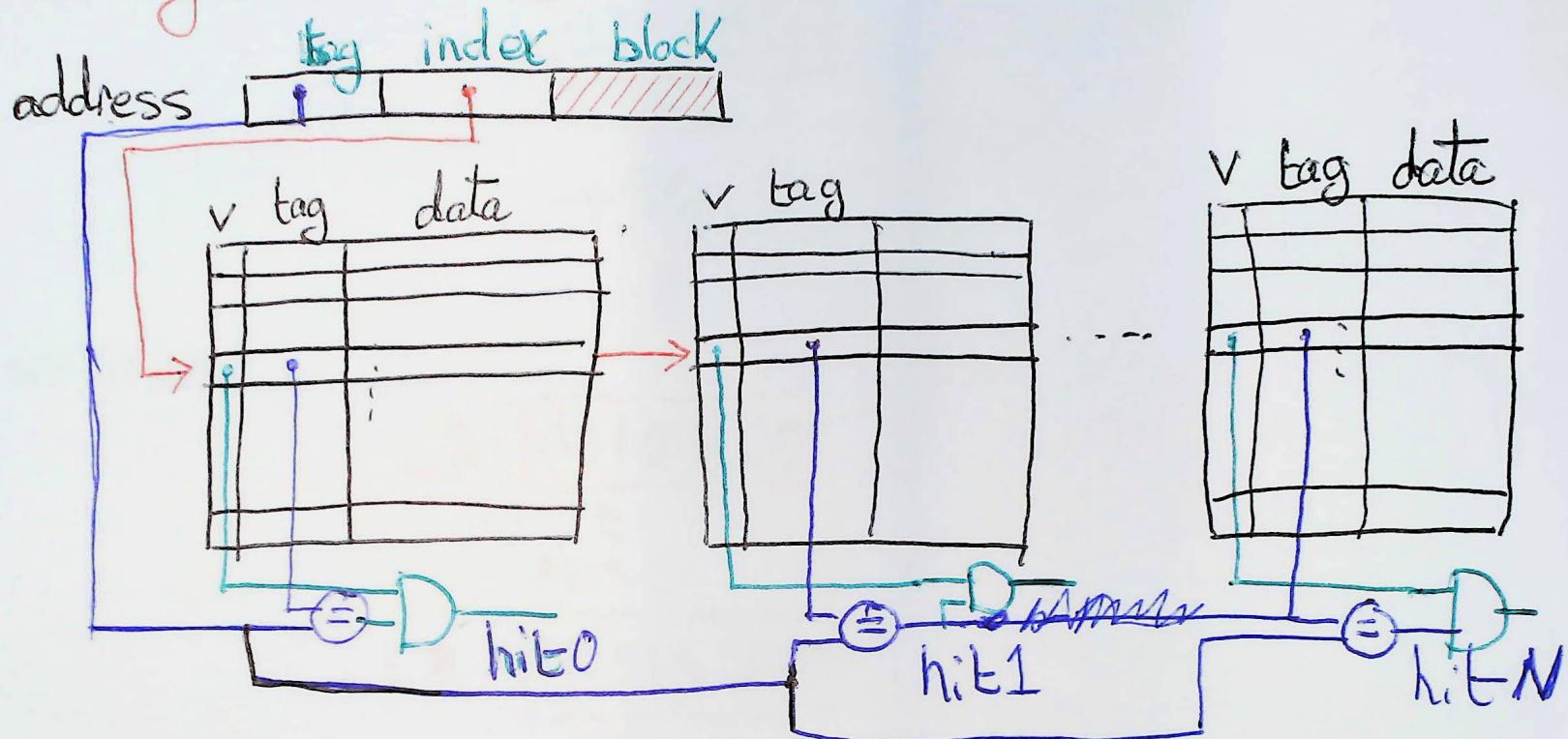
Address (showing bit positions)

31 ... 14 13 ... 6 5 ... 2 1 0





## Associativity in caches



- associativity helps avoid thrashing  
on indexes.

what happens to stores?

data comes from CPU.

nothing to fetch from further down hierarchy on miss.

Write consistency problem:

written blocks can be clobbered by later requests that land on same index.

then future reads of the data will come from further down hierarchy, which will be stale (data from before CPU store).

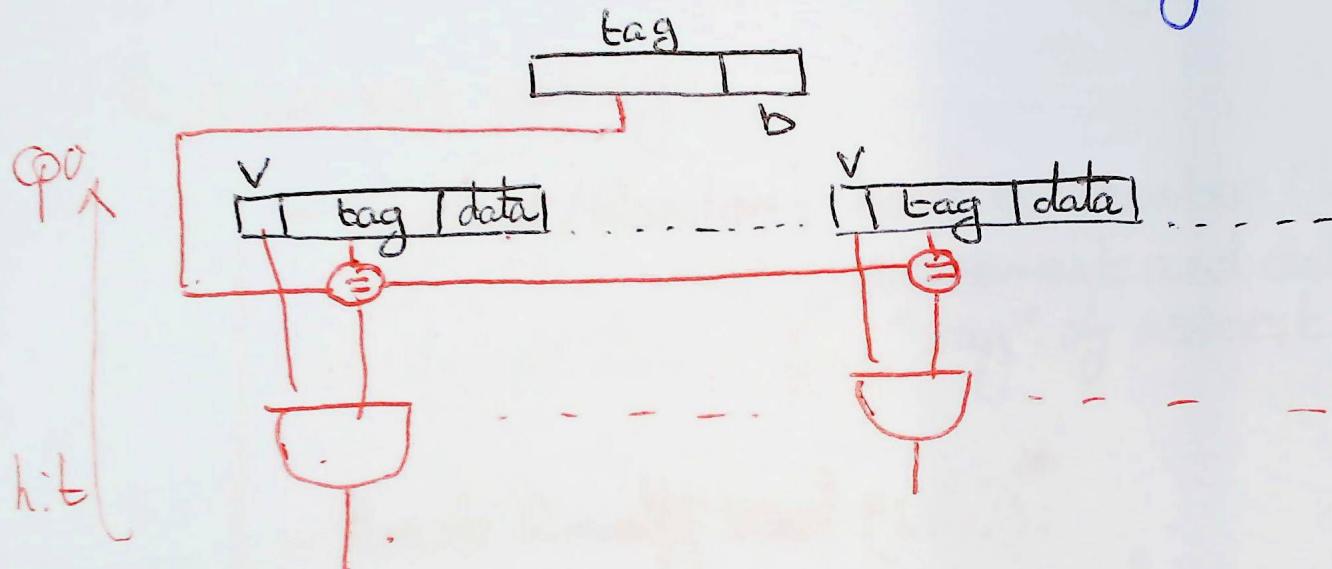
Solutions

{ - write-through: writes go all the way down the hierarchy, to make sure no stale version of data is stored anywhere.

- write-back: mark data stored by CPU as "dirty".  
only send writes to next level of memory if its dirty and needs to be kicked out.

## Fully associative cache

Extreme case of associativity



the index is shrunk to zero bits (only one index)

What to do on cache miss  
when we have associativity?

Replacement policy:

- Counter/Random: have a counter (hashed) to increment and determine next "way" of associativity to be Replaced.
- Least Recently Used (LRU):  
Keep track of "age" of each "way" of associativity with a counter that is reset on accesses to that way and incremented on access to other "way" in that index.  
That "age" tells you who is LRU.
- Data-specific replacement: provide mechanisms for inserting data with high age, make it kicked out soon.

Need to be able to quantify cache quality.

Cache hit rate: measure of cache performance

$$\text{hit rate} = \frac{\text{number of accesses that "hit"}}{\text{number of accesses}} \\ (\# \text{hits} + \# \text{misses})$$

Different types of cache misses

- Compulsory miss: First time accessing data/block/address.
- Capacity miss: not first access to the block but got "Ricked out" due to limited capacity.
- Conflict miss: due to competition for same index while total capacity still not exceeded.

# Cache hierarchy Average Access Time /Mem

$$ATT = t_1 + m_1 t_2 + m_1 m_2 t_3 + m_1 m_2 m_3 t_4 \dots$$

$t_n$  = access time to cache level n

$m_n$  = miss rate of cache level n (1 - hit rate)

# Design Choices in caches

choice	hit/miss rate	access time
#indexes	capacity ↑ hit rate ↑ capacity miss ↓ conflict ... ↓	access time ↑ larger cap. on wires.
associativity	capacity ↑ hit rate ↑ capacity conflict miss ↓	access time ↑↑ slower tag compare
Block size	capacity ↑ hit rate ↑ compulsory miss ↓	access time ↑ (loss) due to slower <u>Fill on miss</u>

	new			old
A	A			
B	B	A		
C	C	B	A	
D	D	C	B	A
E	E	D	C	B
A	A	E	D	C
B	B	A	E	D
A	A	B	E	D
B	B	A	E	D
C	C	B	A	E
D	D	C	B	A
E	E	D	C	B

Miss/hit

M  
m  
m  
M  
m  
m  
m  
h  
h  
M  
m

$$\text{hit rate} = \frac{2}{12} = \frac{1}{6}$$