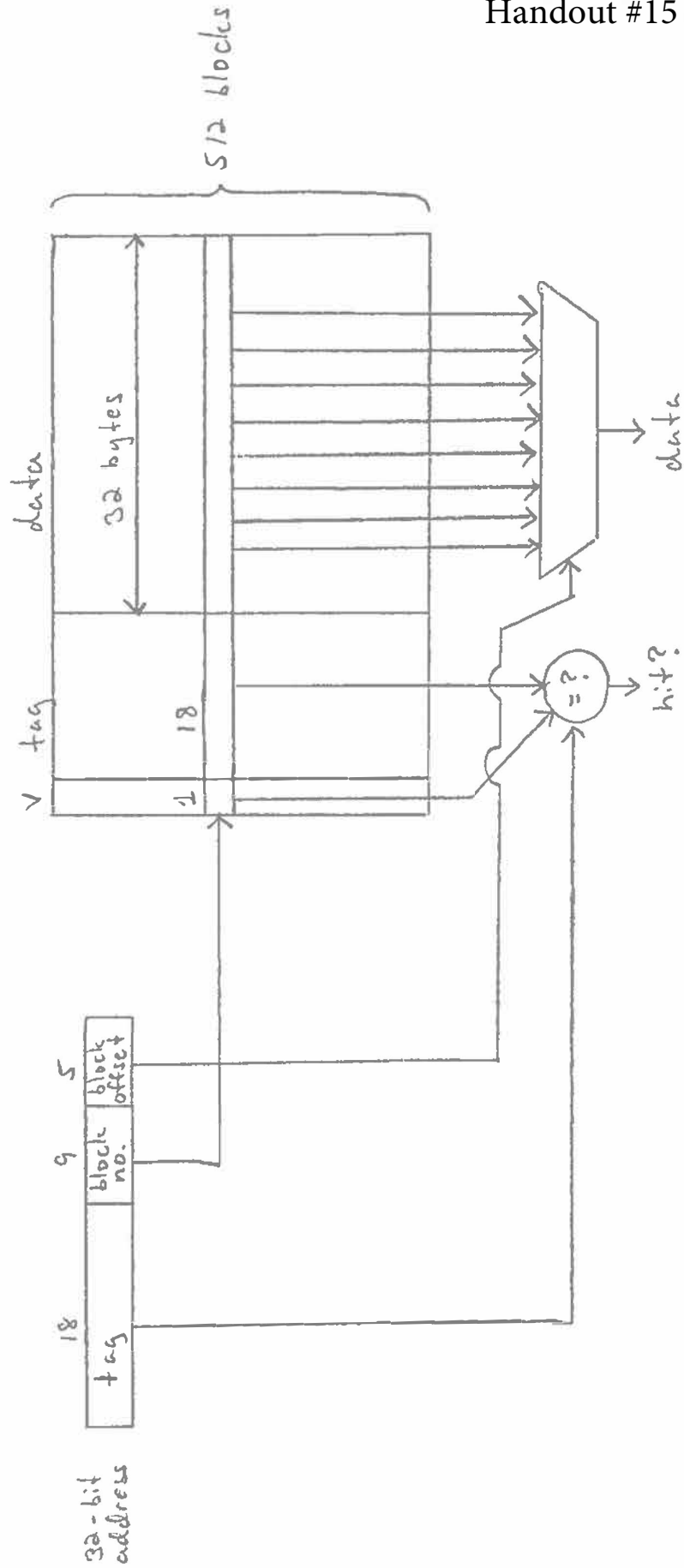


# Direct-Mapped Cache Lookup

Ex: Cache Size = 16 KB  
 Block Size = 32 bytes  
 $\# \text{ blocks} = \frac{16 \text{ KB}}{32} = 512$

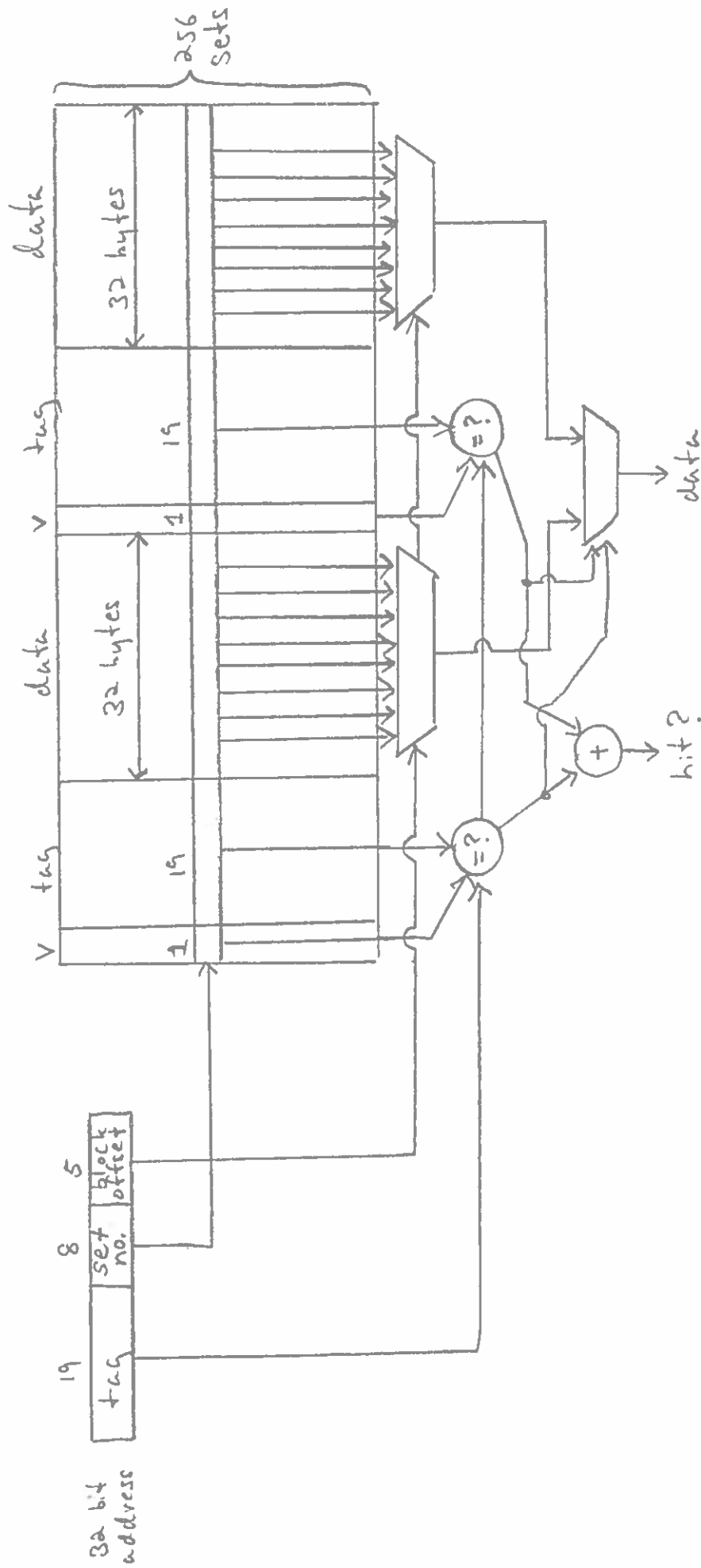


# (2-way) Set Associative Cache Lookup

Ex: Cache Size = 16KB

Block Size = 32 bytes

$$\# \text{ Sets} = \frac{16K}{32 \cdot 2} = 256$$

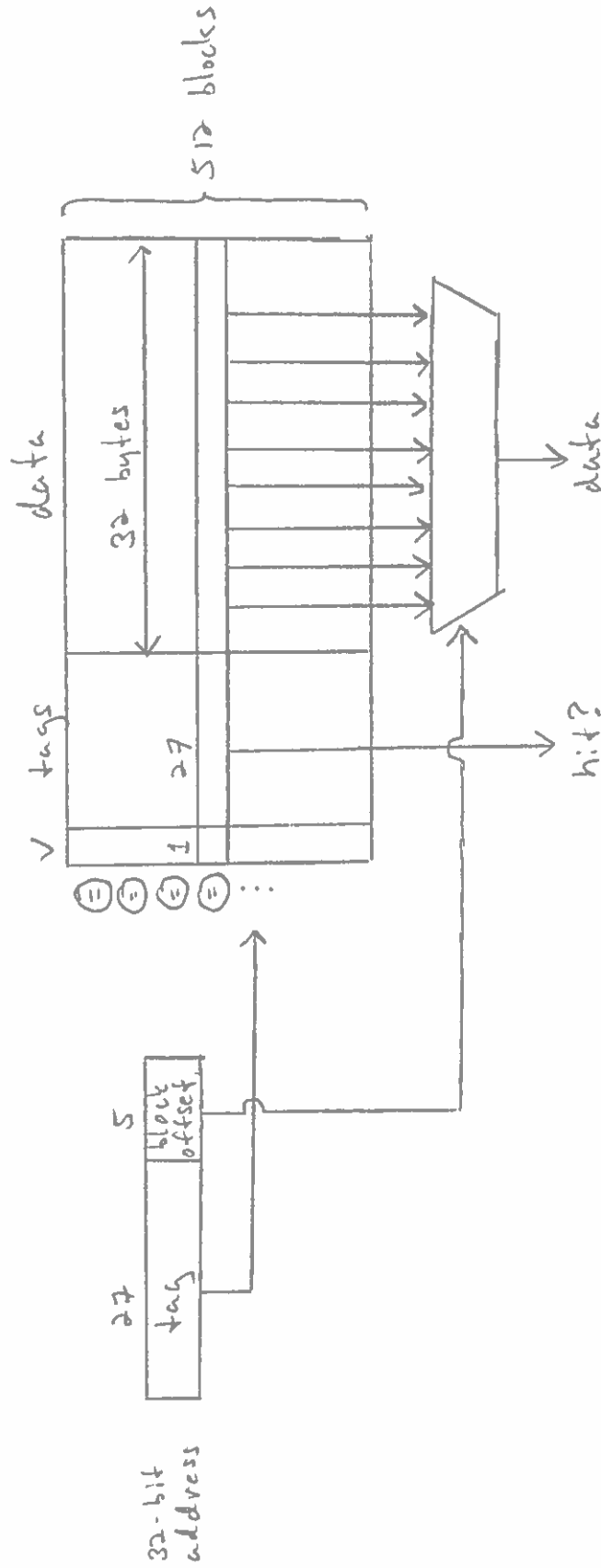


# Fully Associative Cache Lookup

Ex: Cache Size = 16 KB

Block size = 32 bytes

$$\# \text{ blocks} = \frac{16 \text{ K}}{32} = 512$$



# Victim Cache

Motivation:

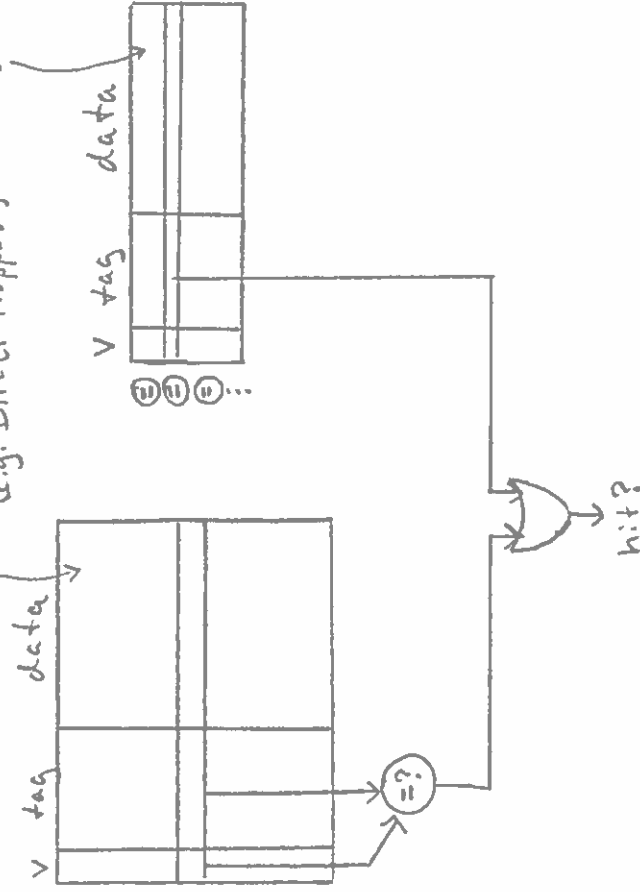
Conflict misses usually isolated to a few sets

Set	1	2	3	4	5
0	X				
1	X	X			
2	X				
3	X	X	X	X	X
4	X				
5	X				
6					
7	X	X			

↑  
number of blocks that map to each set.

Primary Cache:  
• Large (eg. 16KB)  
• Low associativity (eg. Direct Mapped)

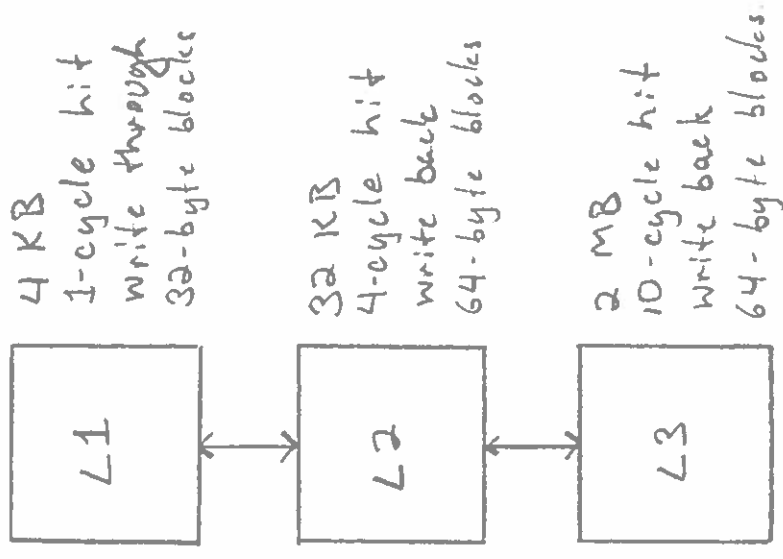
Victim Cache:  
• Small (eg. 16 blocks)  
• Fully Associative



- On access, check both caches (in parallel)
- Eviction from primary cache → victim cache
- Victim cache hit → primary cache
- If both caches miss, fetch from main memory → primary cache.

# Multi-Level Caches

- Goal: Fast + large cache
- As you go down hierarchy:
  - Larger cache
  - Larger block size
  - Higher associativity
- Inclusion



$$\text{Avg Mem. Acc. Time} = \text{hit}_{L1} + \text{miss rate}_{L1} (\text{hit}_{L2} + \text{miss rate}_{L2} (\text{hit}_{L3} + \text{miss rate}_{L3} \cdot \text{miss penalty}))$$