

Lecture 23: Cache Examples

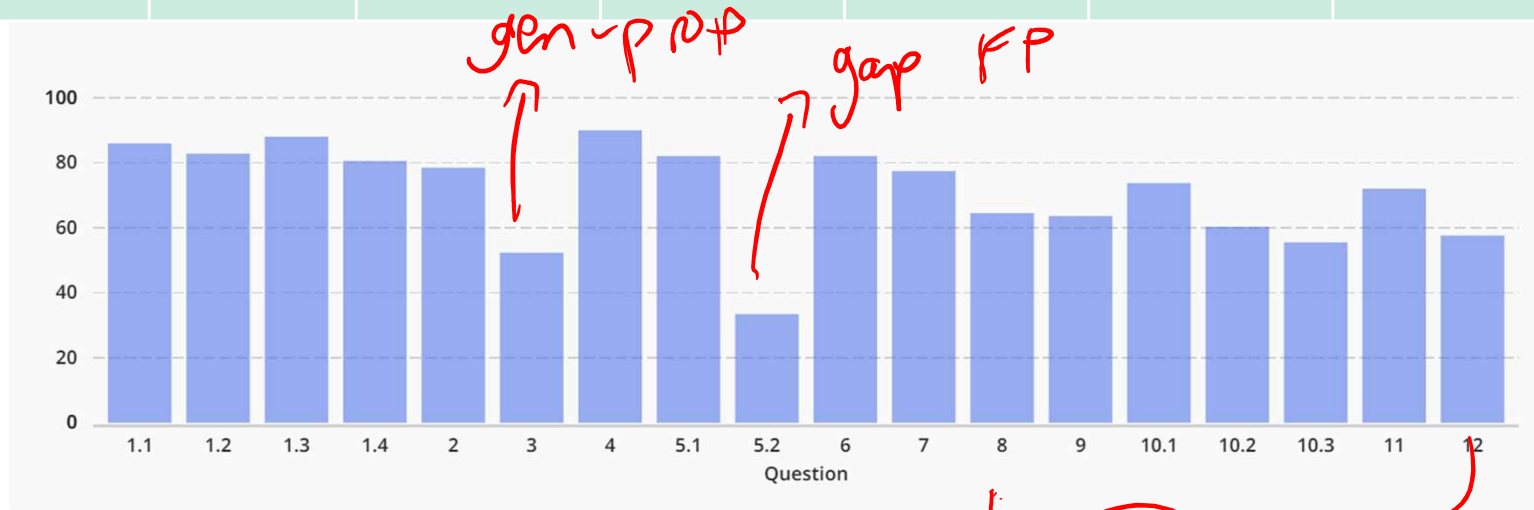
- Today's topics:

- Cache access
- Example problems in cache design
- Caching policies

Midterm Stats

Rank or
Total = 30% M + 30% H + 40% F

	A	A-	B+	B	B-	C+/C/C-	D+ or worse
Percentile	Top 16%	32%	45%	58%	71%	86%	
Rank	47	94	133	172	211	256	295
Midterm Score	83.5	76.25	69.5	64	57	48	



Midterm 100.0 points

Minimum	Median	Maximum	Mean	Std Dev ?
28.0%	69.0%	98.5%	67.51%	14.91%

Internet Browser analogy
Direct mapped - 1 block
per row



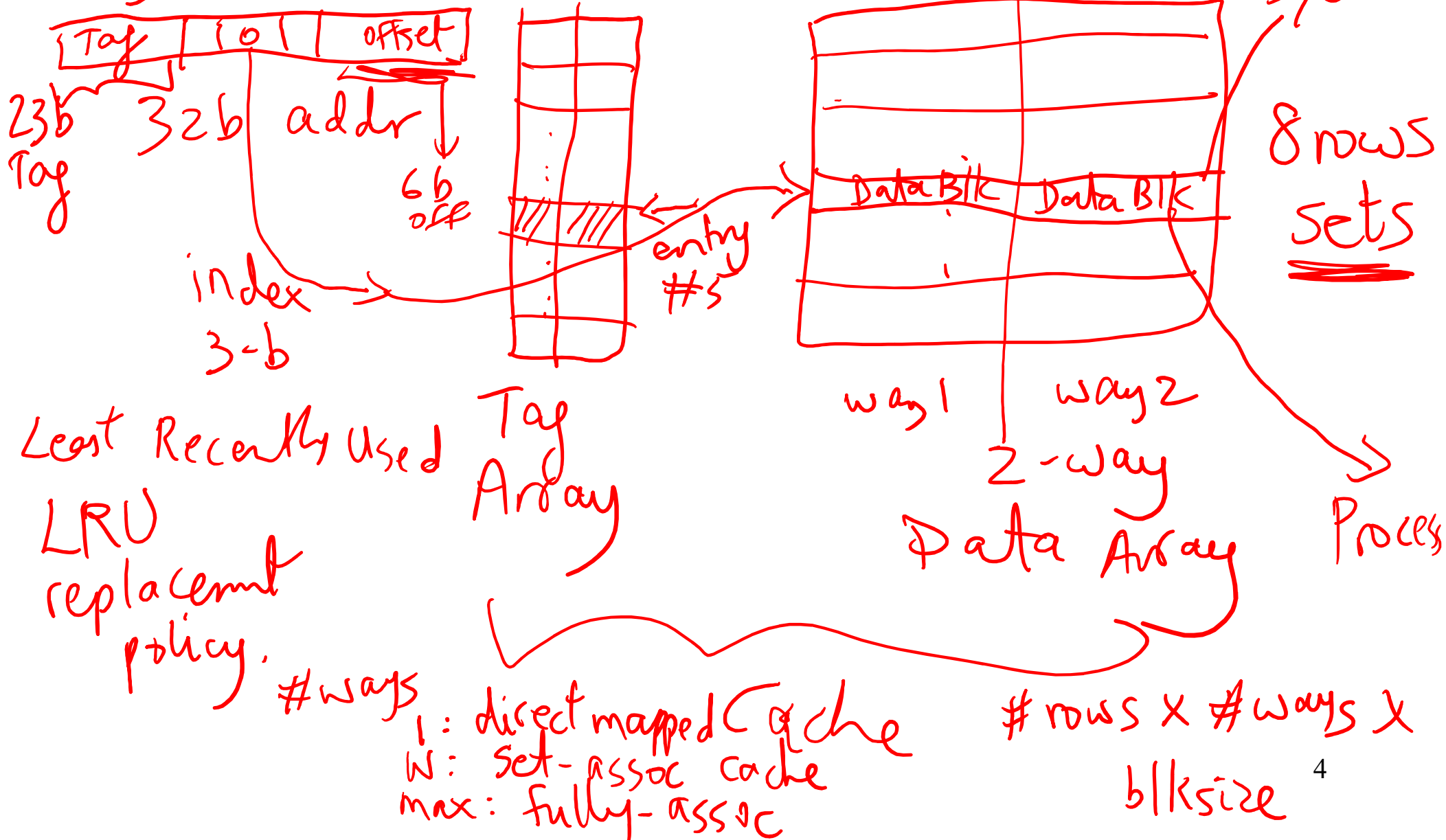
Accessing the Cache

lw, \$t1, 8(\$gp)

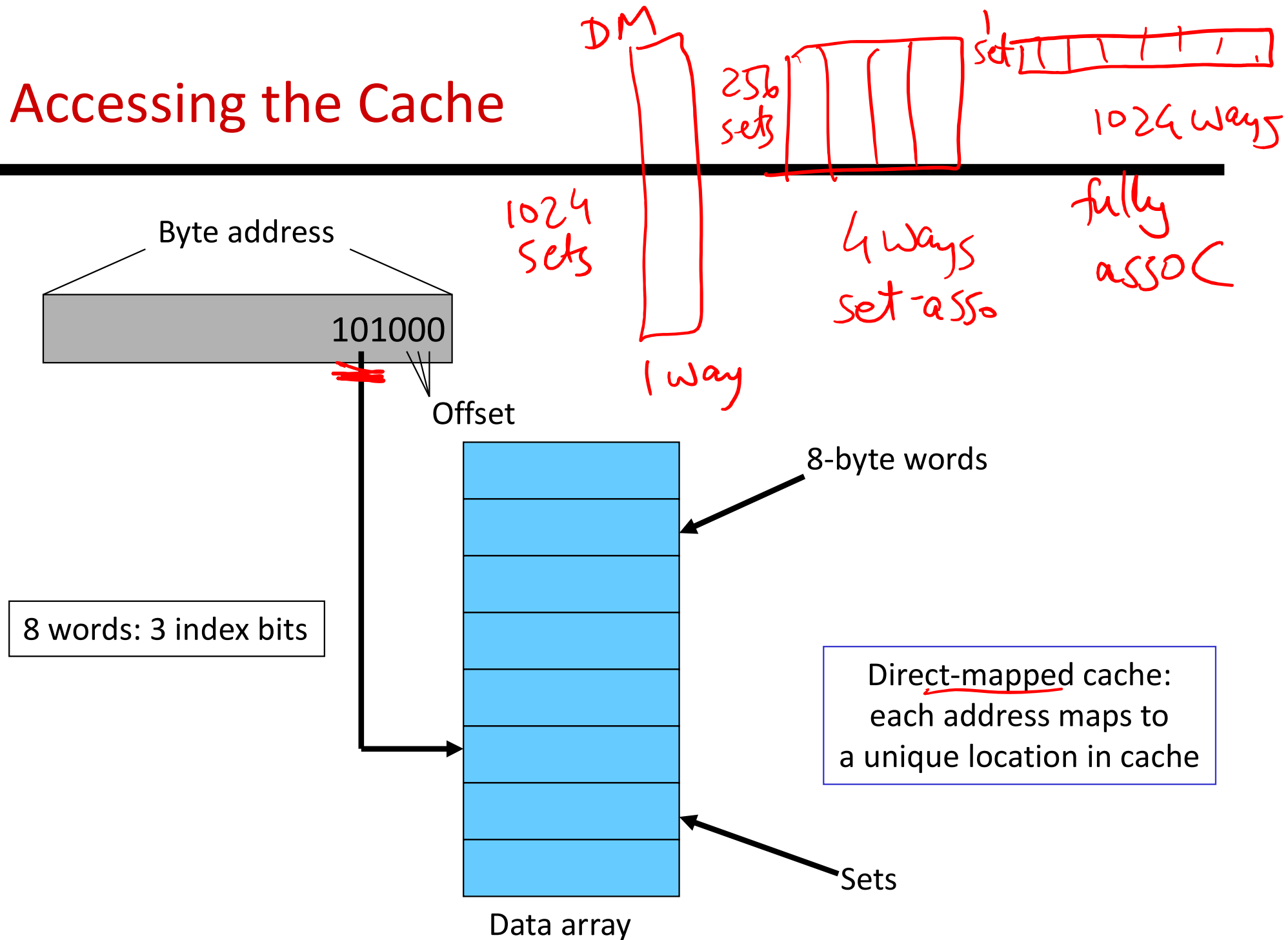
6-b
offset

query

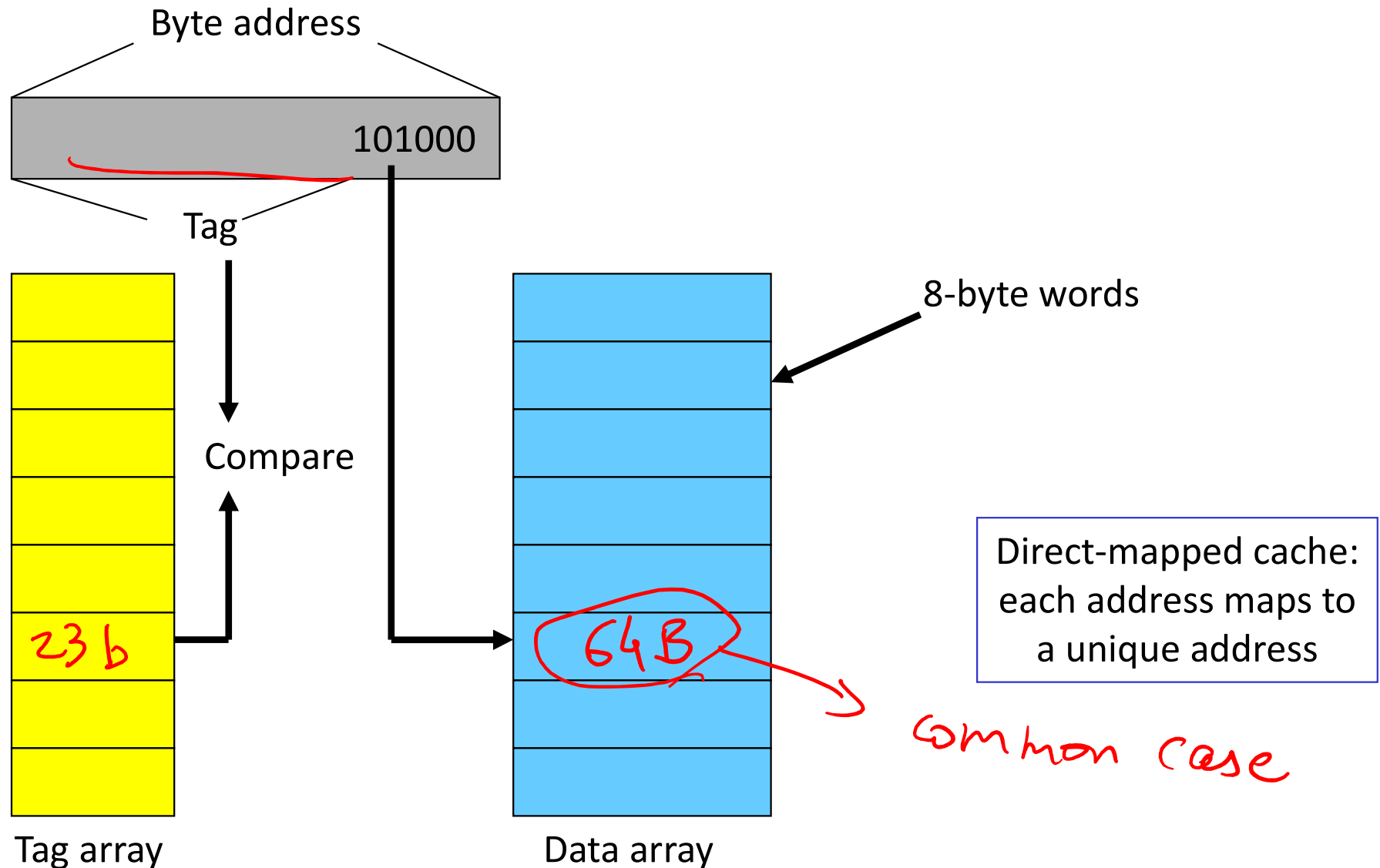
Query is the address



Accessing the Cache



The Tag Array

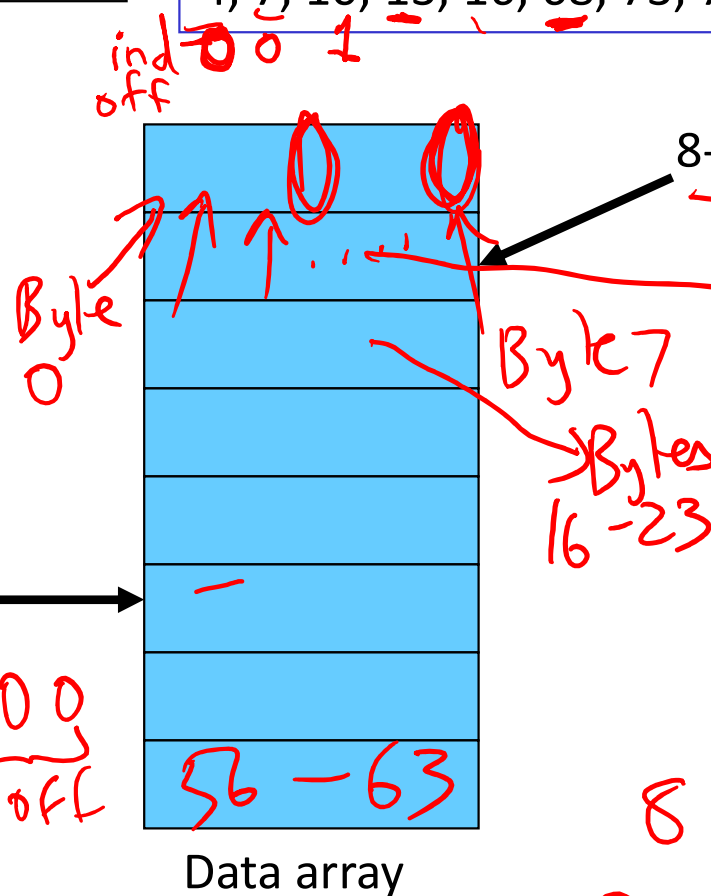
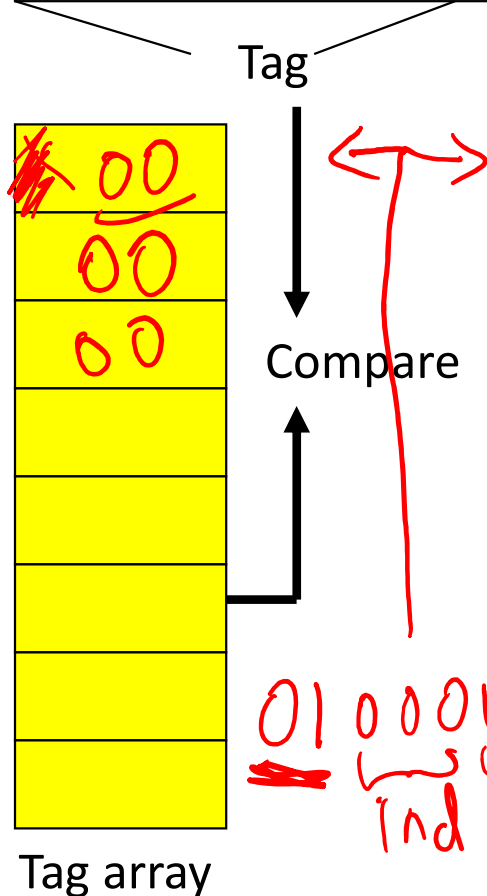


Example Access Pattern

00000111 8b addr
100000100
 2b 3b 3b offset 0-255
 tag index



Assume that addresses are 8 bits long
 How many of the following address requests
 are hits/misses?
 4, 7, 10, 13, 16, 68, 73, 78, 83, 88, 4, 7, 10...

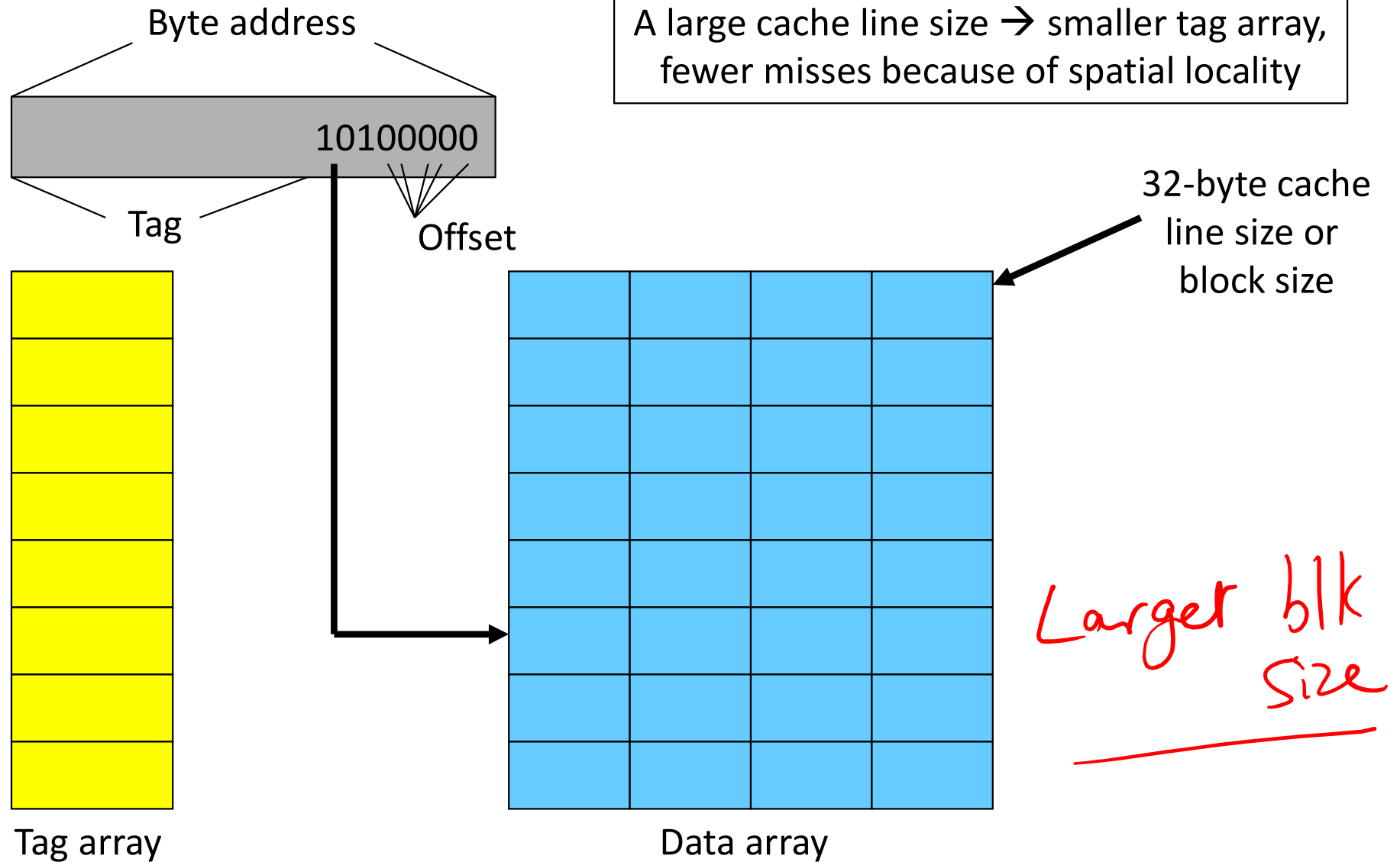


8-byte words Blk size
 Bytes 8-15

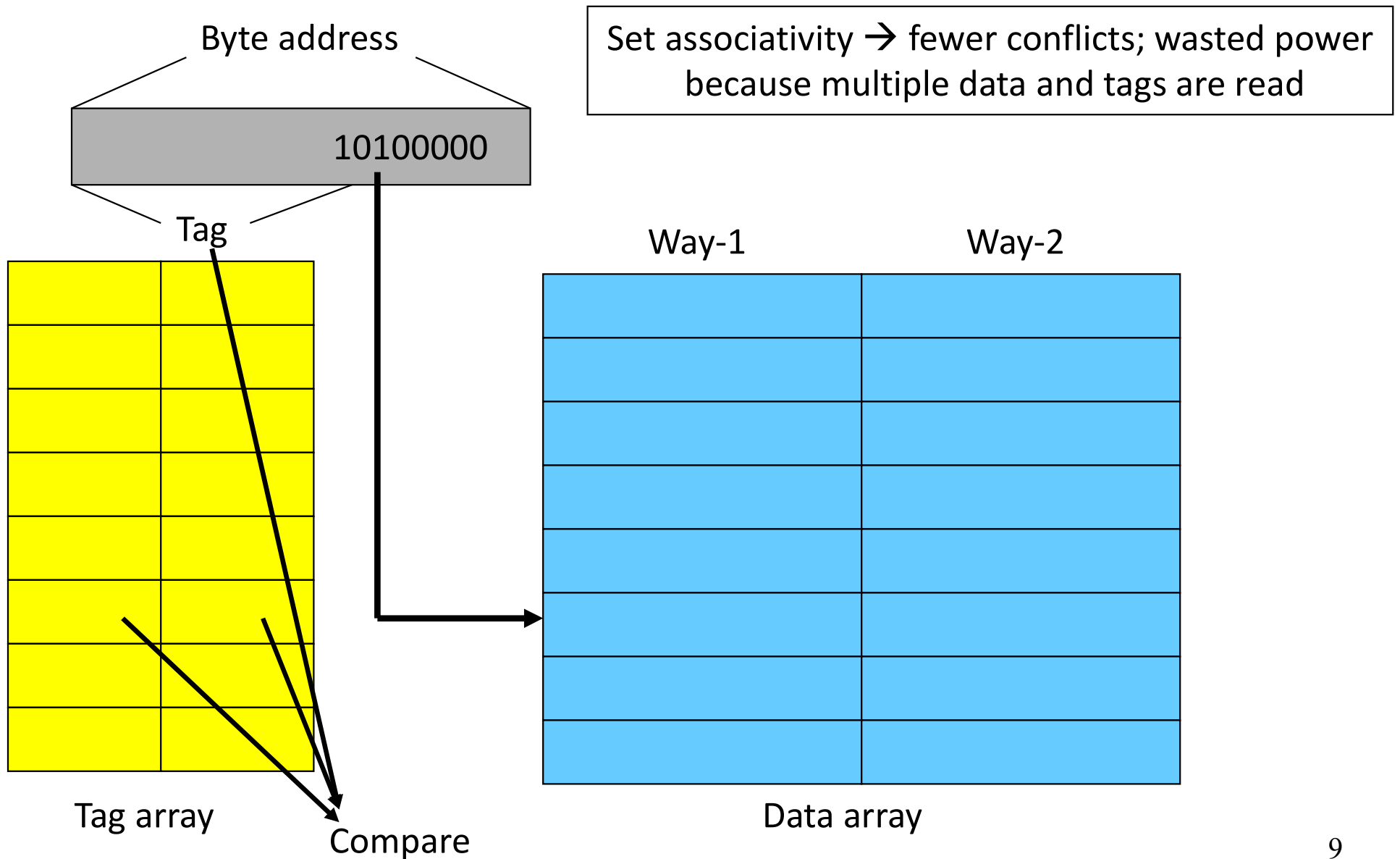
Direct-mapped cache:
 each address maps to
 a unique address

8 sets \Rightarrow 3b index
 00001010 \leftarrow 7 \leftarrow 10
 Addr

Increasing Line Size



Associativity

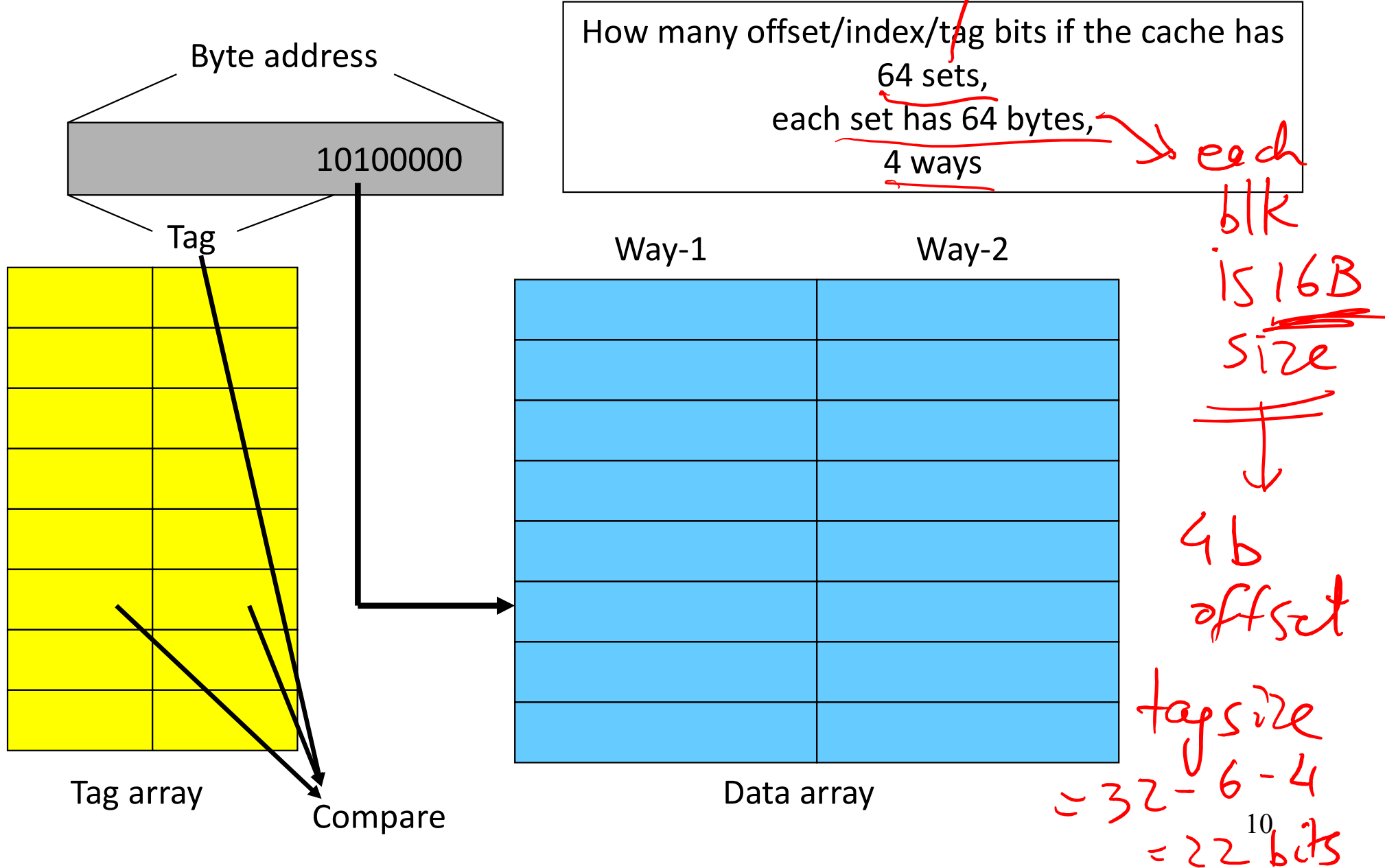


Associativity

$$\log_2 64 = 6$$

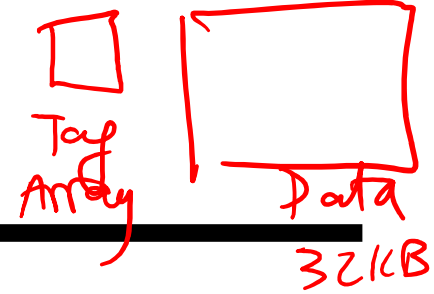
$$2^6 = 64$$

6 index bits



Example 1

Cache size = 32 KB
Data array = $2^5 \times 2^{10}$



ways = 4

- 32 KB 4-way set-associative data cache array with 32 byte line sizes
block size
- How many sets? 256
- How many index bits, offset bits, tag bits? 8 5 19 b
- How large is the tag array?

$$2^5 \times 2^{10} = \text{sets} \times 2^2 \times 2^5 = 32 = 2^5$$

$$\text{sets} = \frac{2^5 \times 2^{10}}{2^2 \times 2^5} = 2^8 = 256 \text{ sets}$$

$$= \# \text{ sets} \times \# \text{ ways}$$

$$\times \text{tag size}$$

$$= 256 \times 4 \times 19 \text{ b}$$

$$= 19 \text{ Kb} = 2.375 \text{ KB}$$

rows cols

$$\text{Cache size} = \# \text{ sets} \times \# \text{ ways} \times \text{blocksize}$$

$$\text{Index bits} = \log_2(\text{sets}) = \log_2(256) = 8$$

$$\text{Offset bits} = \log_2(\text{blocksize}) = \log_2(32) = 5$$

$$\text{Addr width} = \text{tag} + \text{index} + \text{offset}$$

Example 1

- 32 KB 4-way set-associative data cache array with 32 byte line sizes

cache size = #sets x #ways x block size

- How many sets? 256
- How many index bits, offset bits, tag bits?

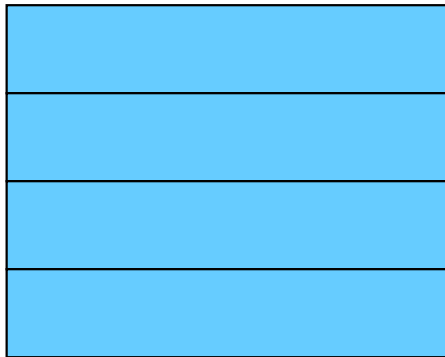
8	5	19
$\log_2(\text{sets})$	$\log_2(\text{blksize})$	addrsz-index-offset

- How large is the tag array?

tag array size = #sets x #ways x tag size
= 19 Kb = 2.375 KB

Example 2

Show how the following addresses map to the cache and yield hits or misses.
The cache is direct-mapped, has 16 sets, and a 64-byte block size.
Addresses: 8, 96, 32, 480, 976, 1040, 1096



Offset = address % 64 (address modulo 64, extract last 6)
Index = address/64 % 16 (shift right by 6, extract last 4)
Tag = address/1024 (shift address right by 10)

32-bit address

22 bits tag 4 bits index 6 bits offset

8:	0	0	8	M
96:	0	1	32	M
32:	0	0	32	H
480:	0	7	32	M
976:	0	15	16	M
1040:	1	0	16	M
1096:	1	1	8	M



Example 3

$$\text{Exec time} = 1000 + 600 + 3000 \\ = 4600 \text{ cyc}$$

$$\text{CPI} = \frac{4600 \text{ cyc}}{1000 \text{ instrs}} = 4.6$$

- A pipeline has CPI 1 if all loads/stores are L1 cache hits
40% of all instructions are loads/stores
85% of all loads/stores hit in 1-cycle L1
50% of all (10-cycle) L2 accesses are misses
Memory access takes 100 cycles
What is the CPI?

$$15\% \text{ of } 400 = 60$$

1000 instrs — 1000 cycles

400 lds/sts

340 L1 hits — no stalls (60

60 L1 misses — 600 stalls $\times 10$)

60 L2 load/store \uparrow

30 L2 misses (mem acc) 3000 stalls

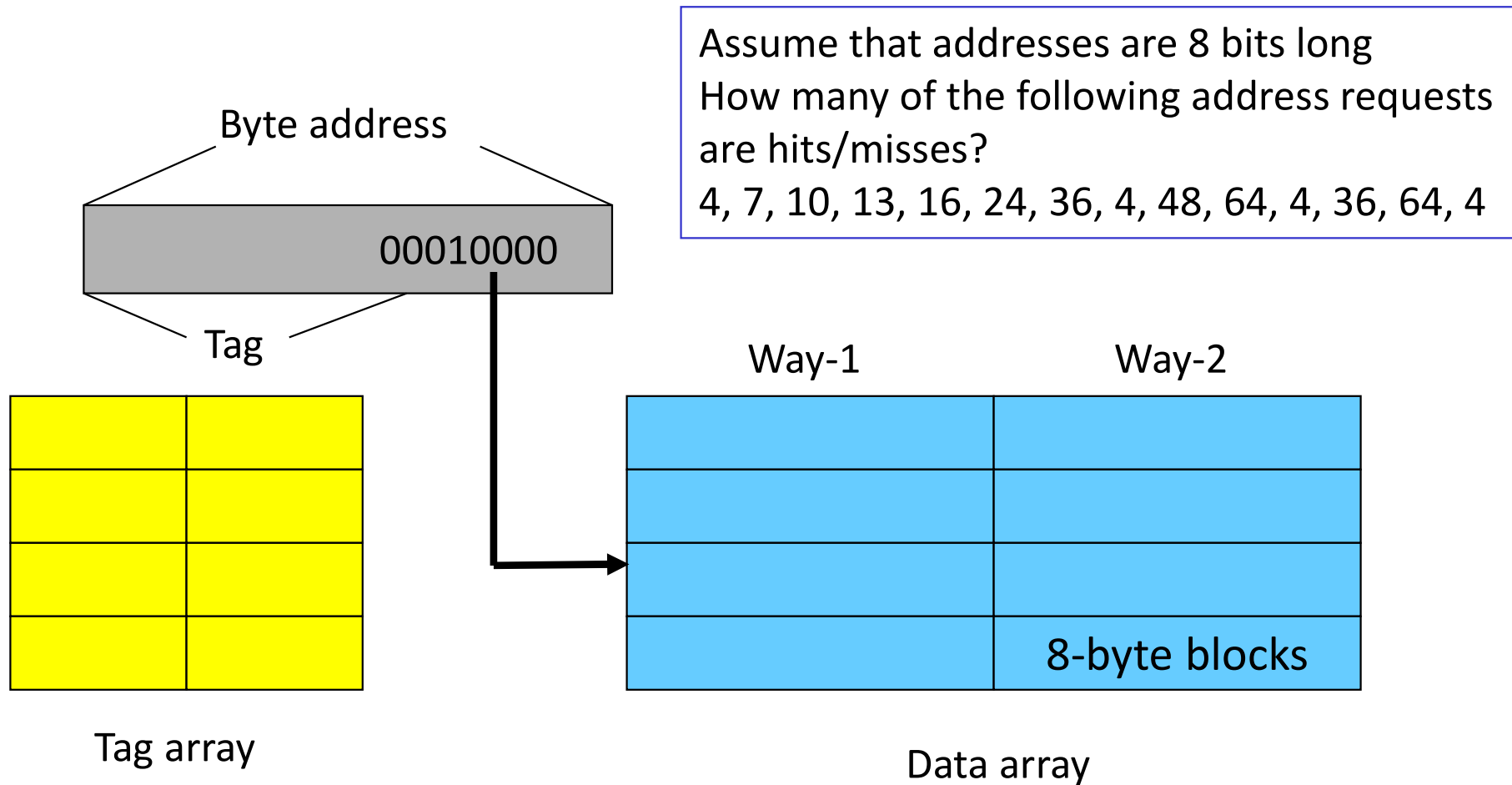
Example 3

- A pipeline has CPI 1 if all loads/stores are L1 cache hits
40% of all instructions are loads/stores
85% of all loads/stores hit in 1-cycle L1
50% of all (10-cycle) L2 accesses are misses
Memory access takes 100 cycles
What is the CPI?

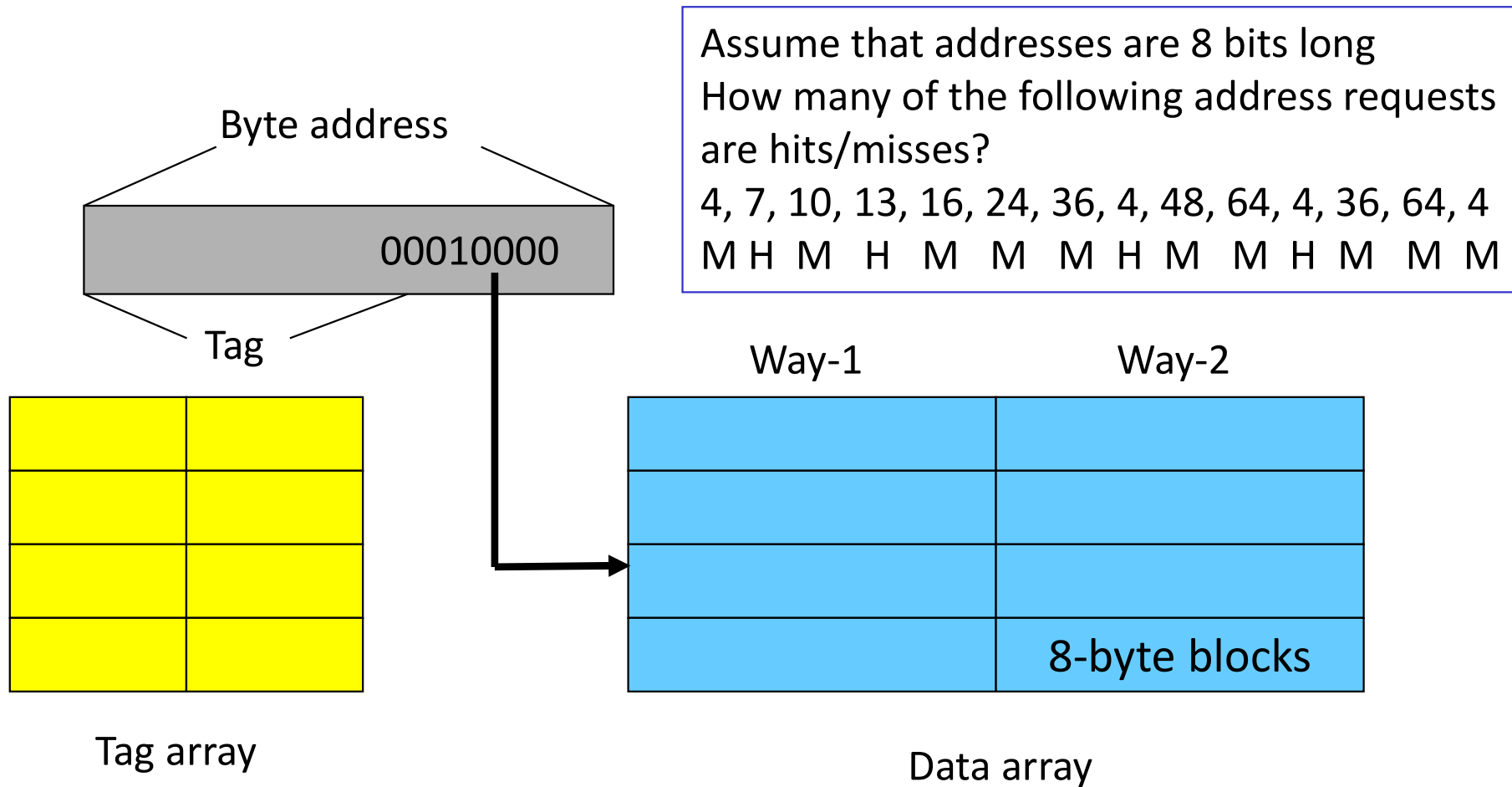
Start with 1000 instructions

1000 cycles (includes all 400 L1 accesses)
+ 400 (ld/st) x 15% x 10 cycles (the L2 accesses)
+ 400 x 15% x 50% x 100 cycles (the mem accesses)
= 4,600 cycles
CPI = 4.6

Example 4



Example 4



Cache Misses

- On a write miss, you may either choose to bring the block into the cache (write-allocate) or not (write-no-allocate)
- On a read miss, you always bring the block in (spatial and temporal locality) – but which block do you replace?
 - no choice for a direct-mapped cache
 - randomly pick one of the ways to replace
 - replace the way that was least-recently used (LRU)
 - FIFO replacement (round-robin)

Writes

- When you write into a block, do you also update the copy in L2?
 - write-through: every write to L1 → write to L2
 - write-back: mark the block as dirty, when the block gets replaced from L1, write it to L2
- Writeback coalesces multiple writes to an L1 block into one L2 write
- Writethrough simplifies coherency protocols in a multiprocessor system as the L2 always has a current copy of data

Types of Cache Misses

- Compulsory misses: happens the first time a memory word is accessed – the misses for an infinite cache
- Capacity misses: happens because the program touched many other words before re-touching the same word – the misses for a fully-associative cache
- Conflict misses: happens because two words map to the same location in the cache – the misses generated while moving from a fully-associative to a direct-mapped cache