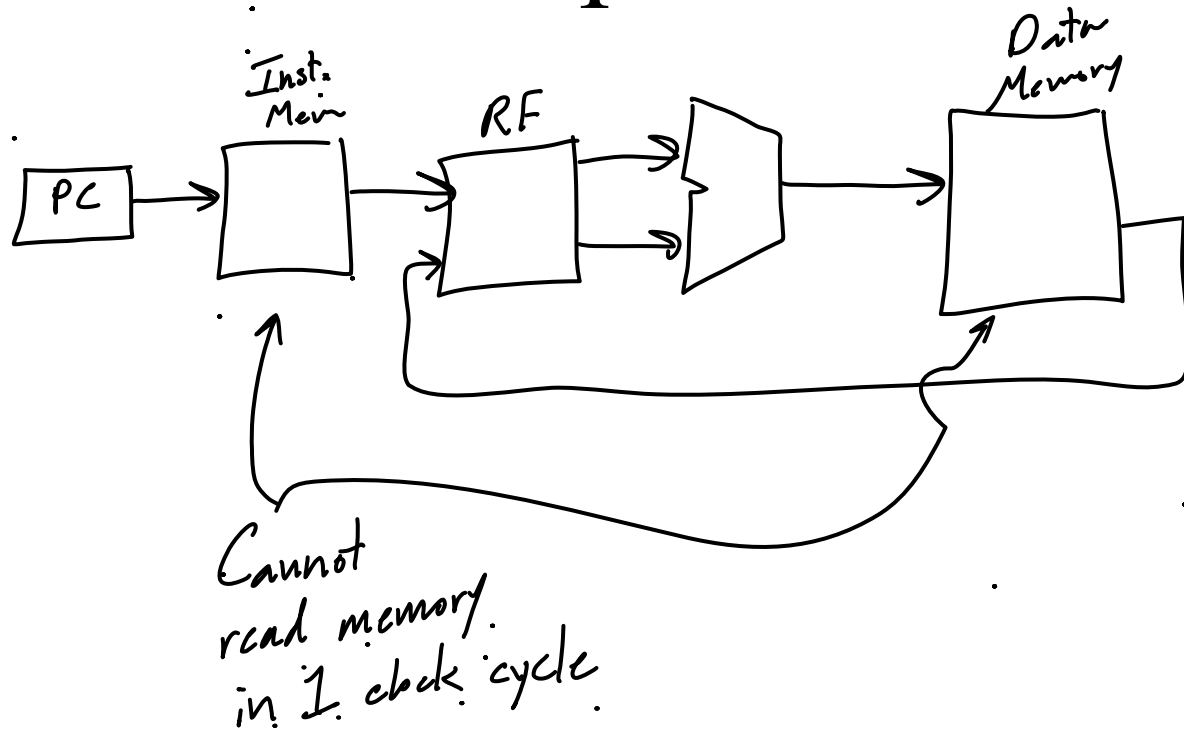# Caching

Ch. 5.
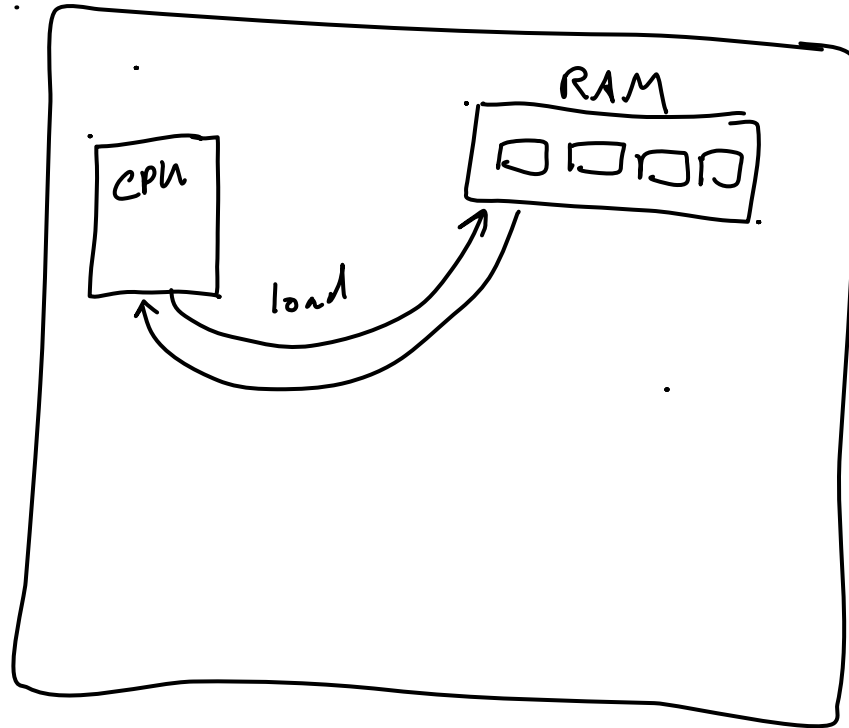
# Single-cycle and pipelined datapaths
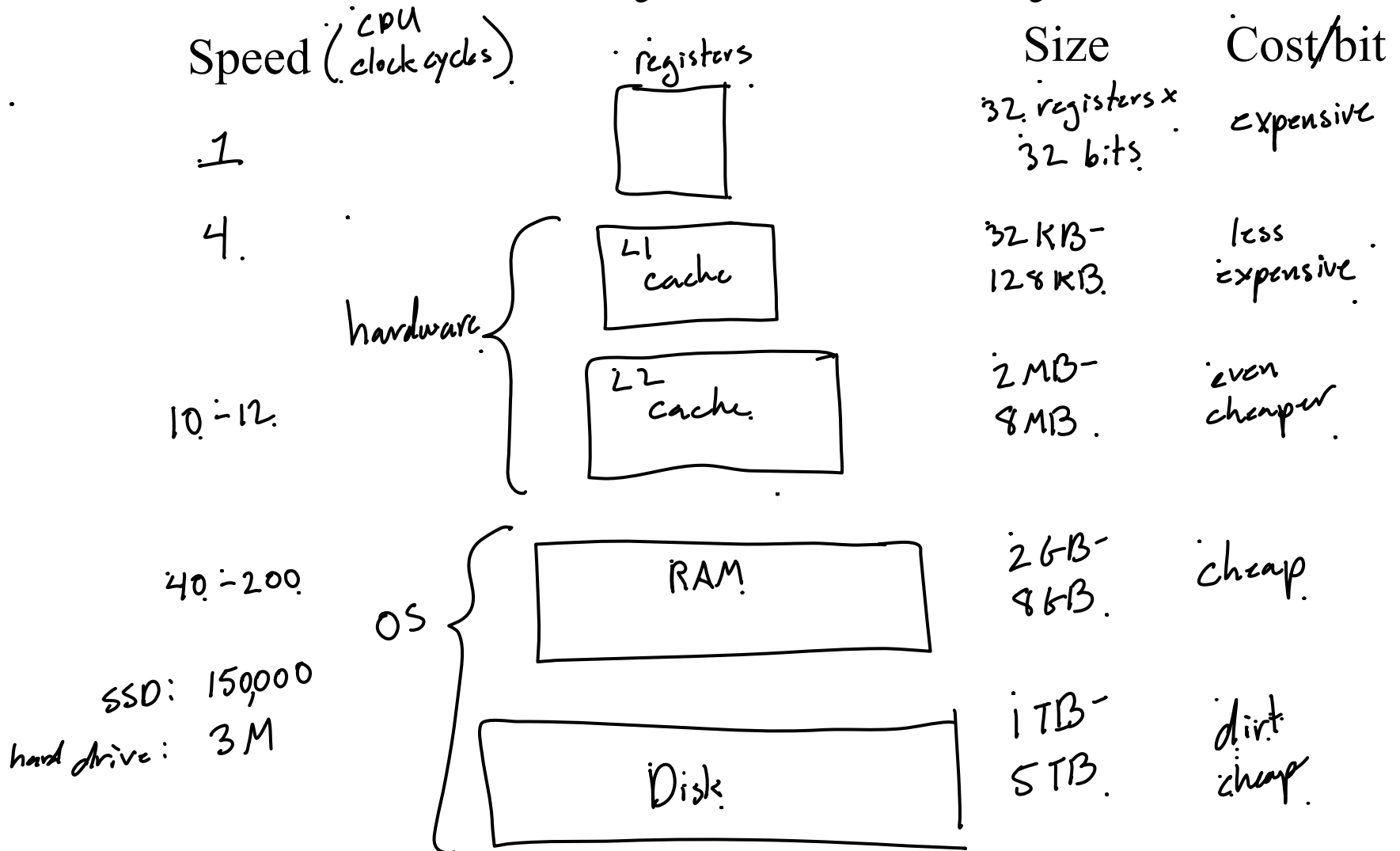
# CPU speed vs. Memory speed

2 GHz - 3 GHz

1066 MHz

CPU

RAM

load

# Memory Hierarchy

| Speed (CPU clock cycles) | | Size | Cost/bit |
|---|---|---|---|
| | registers | 32 registers × 32 bits | expensive |
| 1 | | | |
| 4 | L1 cache | 32 KB - 128 KB | less expensive |
| | hardware { | | |
| 10 - 12 | L2 cache | 2 MB - 8 MB | even cheaper |
| 40 - 200 | OS { RAM | 2 GB - 8 GB | cheap |
| SSD: 150,000 | | | |
| hard drive: 3 M | Disk | 1 TB - 5 TB | dirt cheap |

# Latency Numbers Every Programmer Should Know

■ 1 ns

■ L1 cache reference: 0.5 ns

■ Branch mispredict: 5 ns

■ L2 cache reference: 7 ns

■ Mutex lock/unlock: 25 ns

= ■ 100 ns

■ Main memory reference: 100 ns

= 1 μs

Compress 1 KB with Zippy: 3 μs

= ■ 10 μs

■ Send 1 KB over 1 Gbps network: 10 μs

SSD random read (1 Gb/s SSD): 150 μs

Read 1 MB sequentially from memory: 250 μs

Round trip in same datacenter: 500 μs

= ■ 1 ms

■ Read 1 MB sequentially from SSD: 1 ms

Disk seek: 10 ms

Read 1 MB sequentially from disk: 20 ms

Packet roundtrip CA to Netherlands: 150 ms

# Locality

↳ Programs tend to use
the same data/inst. that
were used recently

```
loop: --- -X
      --- X
      --- X
      --- X
beq---, ---, loop.
```
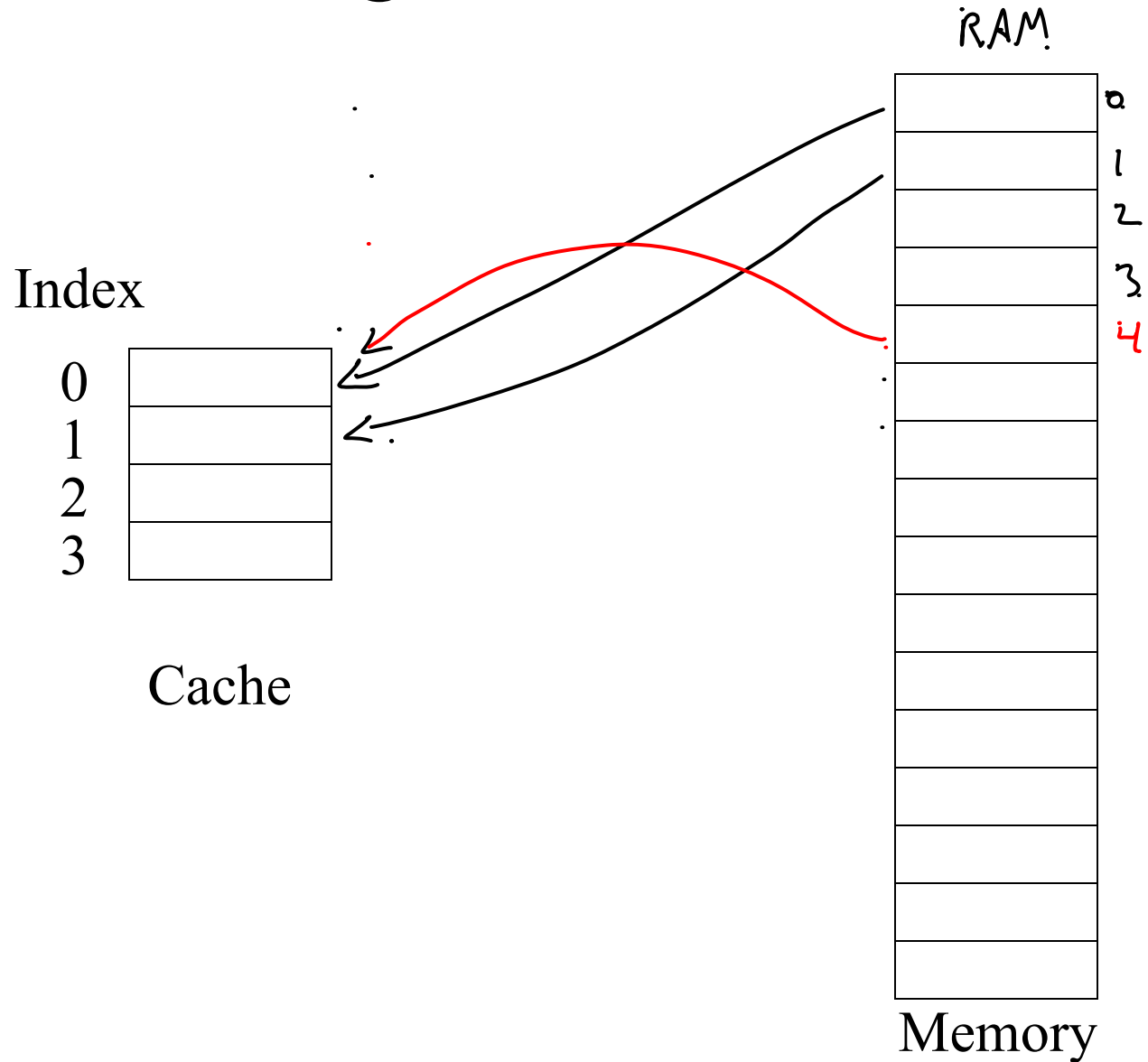
# Program Characteristic

- Temporal Locality – items accessed recently are likely to be accessed in the near future
  (Time)

- Spatial Locality – if an item is accessed, a program will likely access its neighbor
  (Space)

# Cache Design to use locality

- Temporal Locality


- Spatial Locality

# Placing data in a cache

RAM

Index

0
1
2
3

Cache

Memory

# Definitions

- Block (Line) — the unit of data transferred into a cache
- Hit — data is found in the cache
- Miss — data is not found


- Hit time (Access time) — how long to get the data


- Miss Penalty — time to receive the data from a lower level

# Accessing a cache

Index = Address % 4

Address stream:

Index

| | |
|---|---|
| 0 | 8 |
| 1 | 1 |
| 2 | |
| 3 | 7. |

Cache

Addr:   1   4   7   1   5   1   4   7   8   7

Index:  1   0   3   1   1   1   0   3   0   3

        M   M   M   H   M   M   H   H   M   H

Hit Rate = $\frac{4}{10}$ = 40%

Miss Rate = 60%

— HW. due Friday.
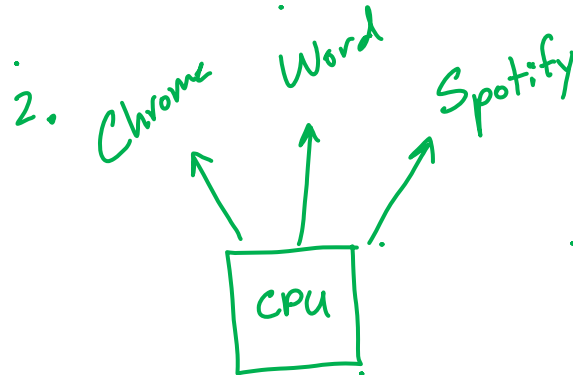
# Addressing into a cache

32-bit address

| | | 1 | 0 |

Tag
(bits remaining)

Index
(which cache block)

Byte offset
(which byte in a word)

**64 KB cache**

16-bit Tag | 14-bit index | 1 0

Byte offset

Valid ← Tag → 16-bits = 2B → Data 1 word = 4 bytes →

Cache Block

$2^{14} = 2^4 \cdot 2^{10} = 16 \cdot 1K = 16K$ entries

= —Tags = 32 KB

hit

# Valid bit

2 Reasons:

   1. On power up, the CPU clears all valid bits.

   2. Chrome    Word    Spotify



CPU

switching programs
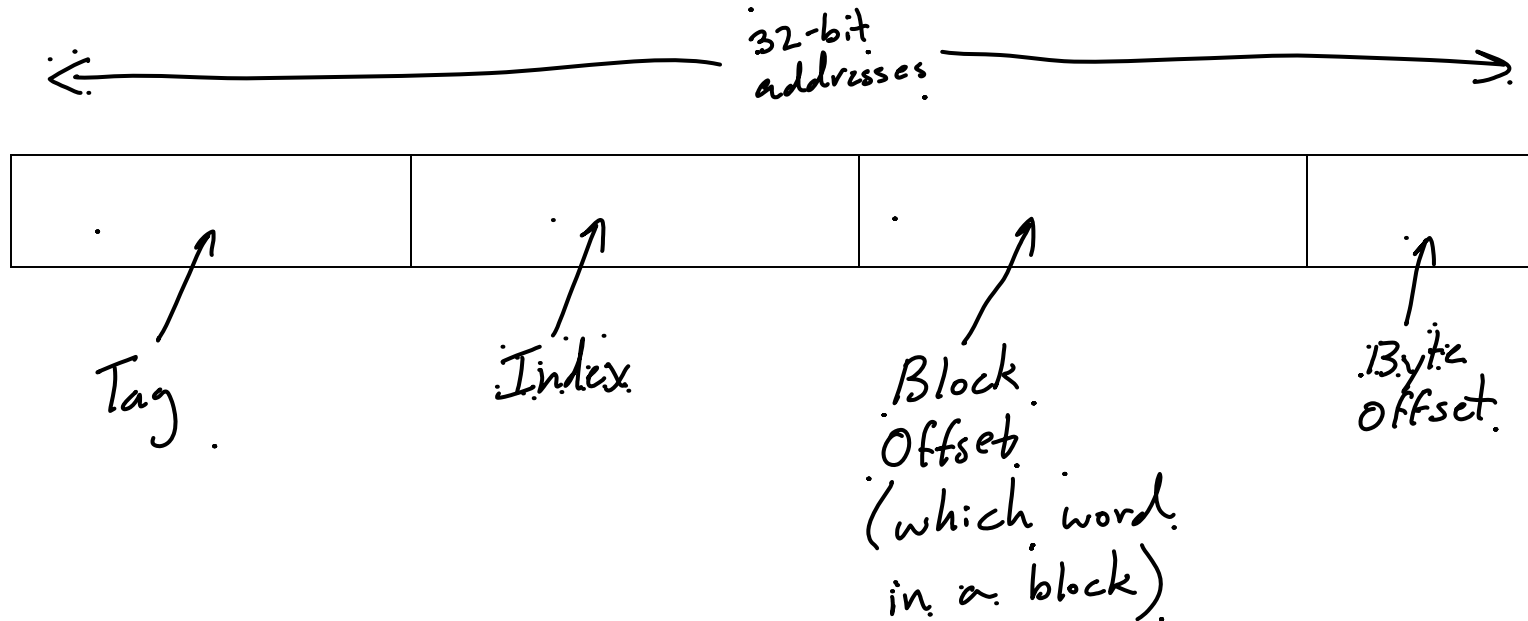
- When a context switch, the CPU can invalidate the cache
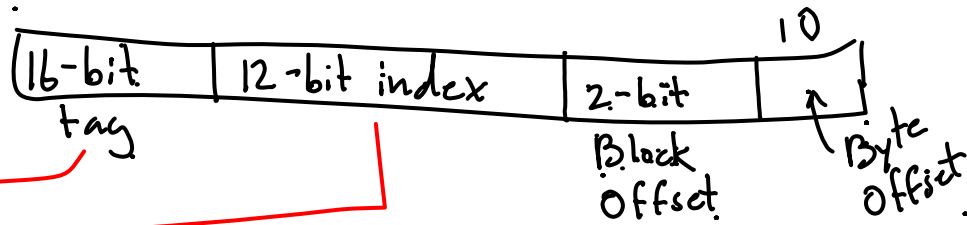
# Taking advantage of spatial locality

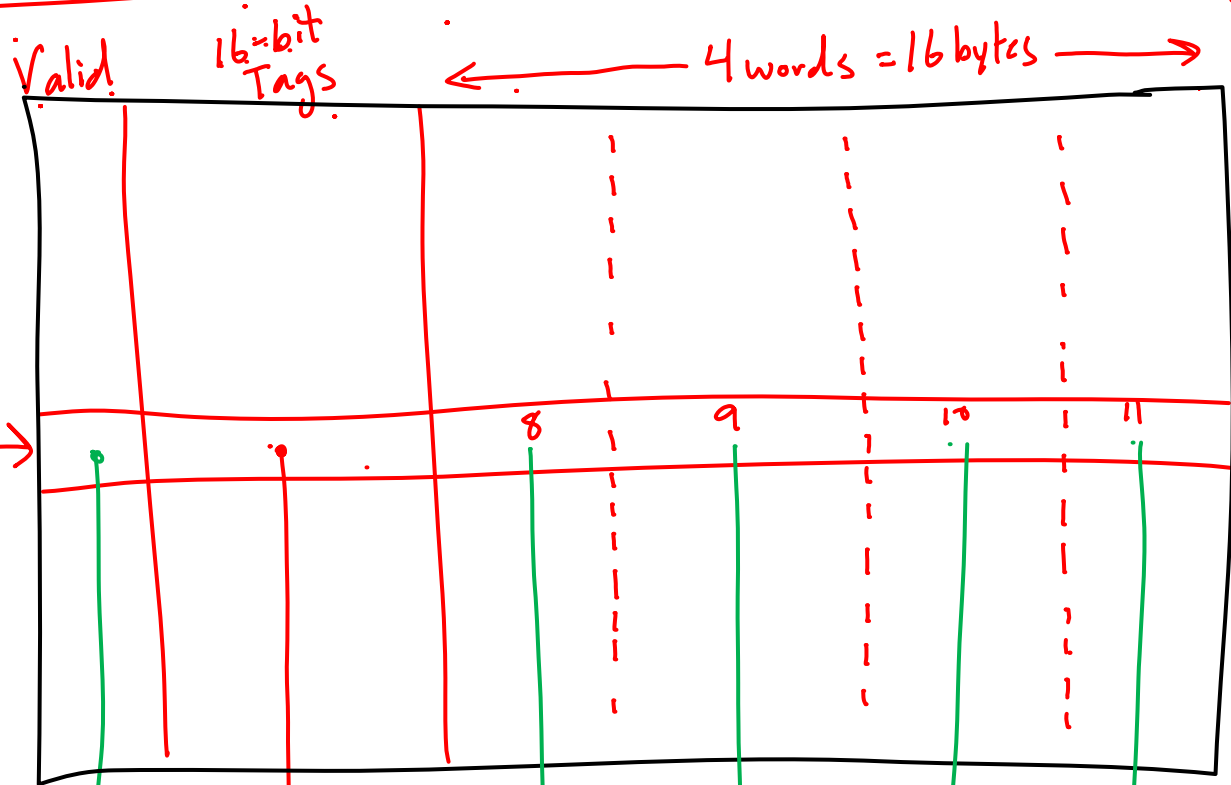Multi-Word Cache Block — store neighboring words in the same block

# Address with multi-word blocks

32-bit addresses

| Tag | Index | Block Offset (which word in a block) | Byte Offset |

16-bit tag | 12-bit index | 2-bit Block Offset | 10 Byte Offset

64KB

same cache size as 1st ex.

Valid | 16-bit Tags | 4 words = 16 bytes

8 9 10 11

$2^{12} = 2^2 \cdot 2^{10}$

4K entries

= 

0 1 2 3

Block Offset

hit

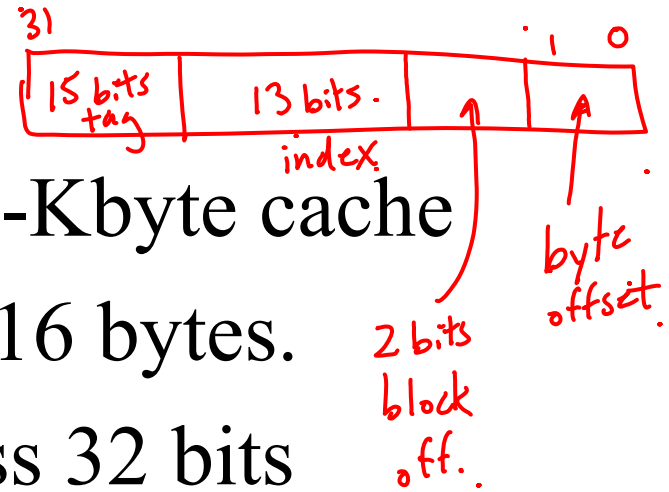data

# Example

31     1   0

| 15 bits tag | 13 bits index | | |

byte offset

2 bits block off.

- You are implementing a 128-Kbyte cache
- The block size (line size) is 16 bytes.
- Each word is 4 bytes, address 32 bits
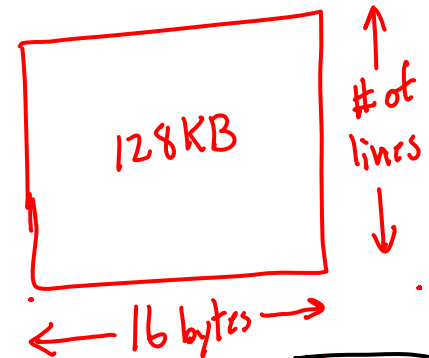- How many bits is the block offset?

line size = 16 bytes = 4 words

$2^2$ → 2 bits block offset

- How many bits is the index?

128 → $2^7 \cdot 2^{10}$ ← 1k

16 → $2^4$

$= 2^{13}$ → 13 bits for index

128KB    #of lines

- How many bits is the tag?

$32 - 13 - 2 - 2 =$ 15 bit tag

← 16 bytes →

Total # bits for tags $= 15 \cdot 2^{13}$ bits

# — Lab 5

— Correlating Branch Predictor
   60+%      90+%
   ↑          ↑
GHR (2, 4, 8)

mips > b
      ↗ branch predictor
        accuracy



→ Prediction

— Bresenham Algorithms

lines
circles

figure.asm →

plot(10, 3)
plot(15, 6)

| | |
|---|---|
| 0 | 10 |
| 1 | 3 |
| 2 | 15 |
| 3 | 6 |

→

10, 3
15, 6

O

Excel
x-y scatter
→