

# CMPE 110: Computer Architecture

## Week 6

### Cache

Jishen Zhao (<http://users.soe.ucsc.edu/~jzhao/>)

[Adapted in part from Jose Renau, Mary Jane Irwin, Joe Devietti, Onur Mutlu, and others]

### Reminder

---

- Quiz 2 is due today midnight

## Review: Temporal and Spatial Locality

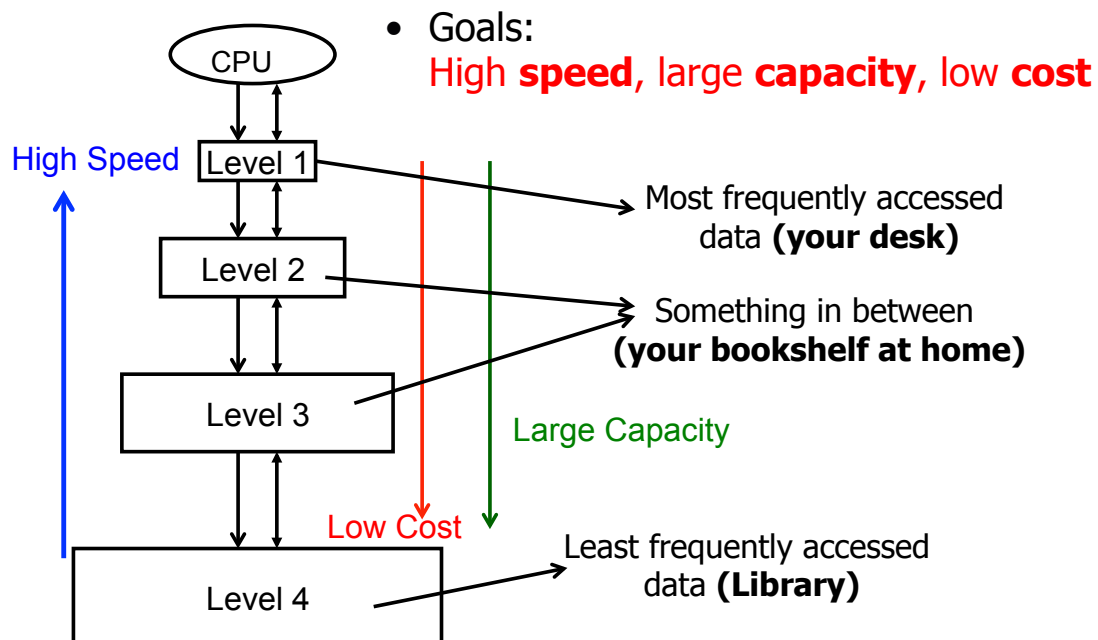
- Which memory accesses demonstrate spatial locality?
- Which memory accesses demonstrate temporal locality?

```
int sum = 0;
int x[1000];

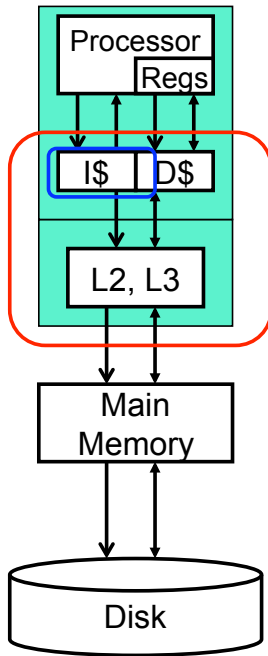
for(int c = 0; c < 1000; c++) {
    sum += c; // Temporal locality: A program tends to access
              the same memory location many times and all
              within a small window of time

    x[c] = 0; // Spatial locality: A program tends to reference a
              cluster of memory locations at a time
}
```

## Review: basic idea of memory hierarchy

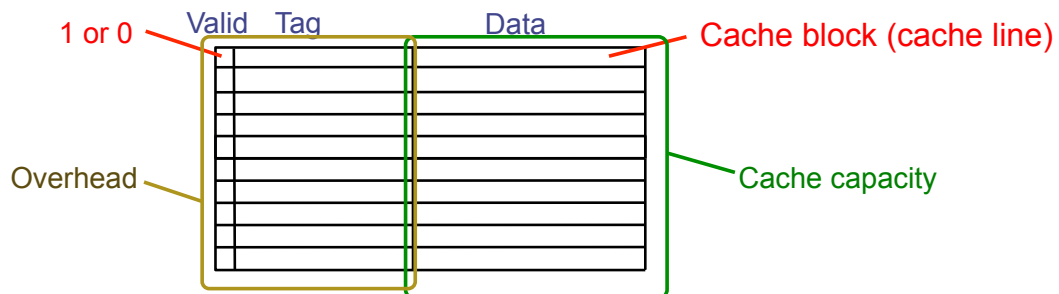


## Review: Cache basics



- I\$ (instruction cache): read only
- D\$ (data cache): read/write
- L2, L3 are lower-level caches for both instruction and data
- Inclusive vs. exclusive vs. non-inclusive cache

## Review: Cache basis



- When data referenced
  - **HIT**: If in cache, use cached data instead of accessing memory
  - **MISS**: If not in cache, bring block into cache (**invalid** → **miss**)
    - Go to the next level of cache to bring this data up
    - Have to kick something else out to do it, if it is full

## Today: Cache (cont.)

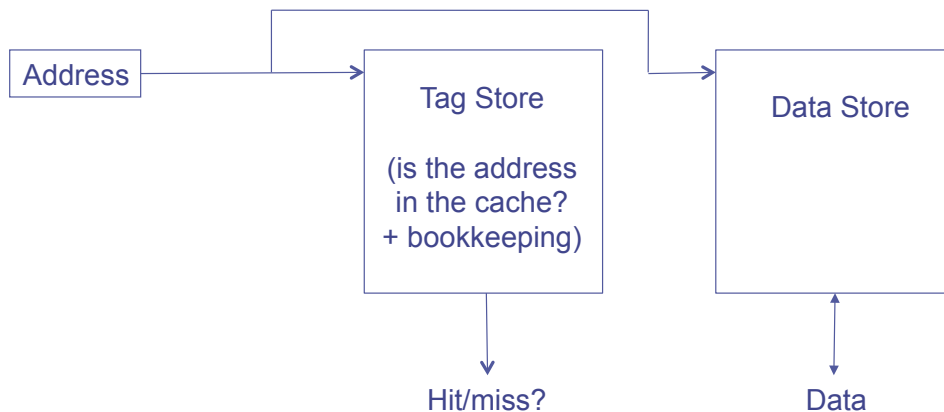
---

- How to access cache with a memory address?
- How to calculate tag overhead?
- Cache and the pipeline

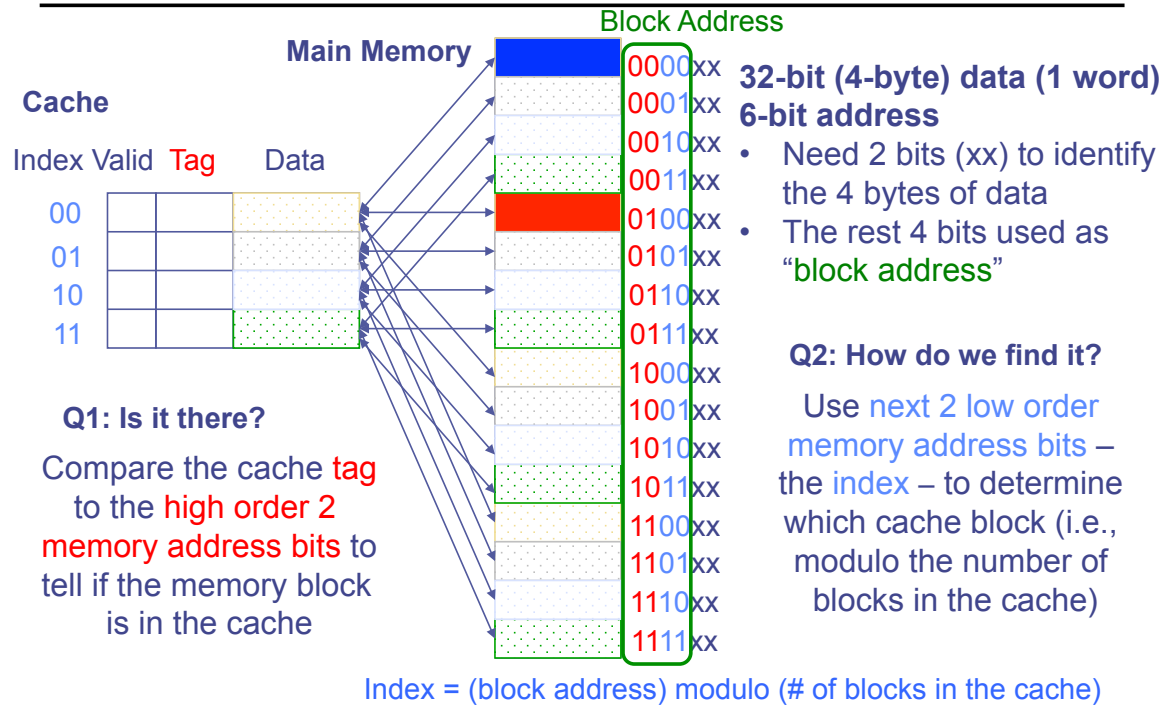
## How is cache accessed?

---

- Two questions to answer (in hardware):
  - Q1: How do we know if a data item is in the cache?
  - Q2: If it is, how do we find it?

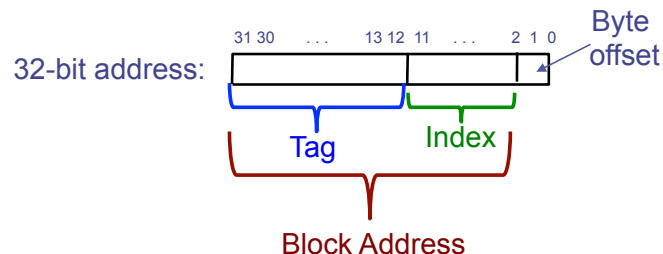


# Caching: A Simple First Example



## Three regions in memory address

- Address is a unique pointer to a data block in main memory



# How is cache accessed?

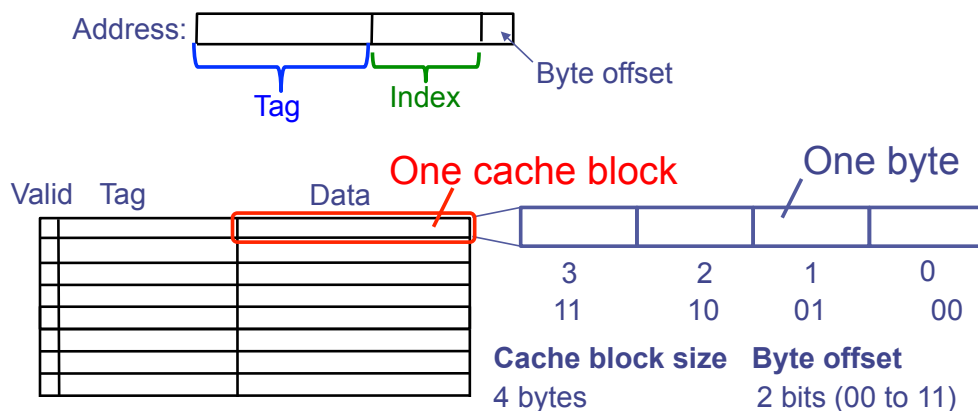
- **Input:** the **address** of data to be accessed
- **Output:**
  - Hit or miss?
  - If hit → we get the data from cache
  - If miss → we don't get the data from cache, have to access the next level in the memory hierarchy

## Cache access with three steps

- Step 1: Identify "byte offset" and "block address"
- Step 2: Calculate "cache index"
- Step 3: Compare "cache tag"

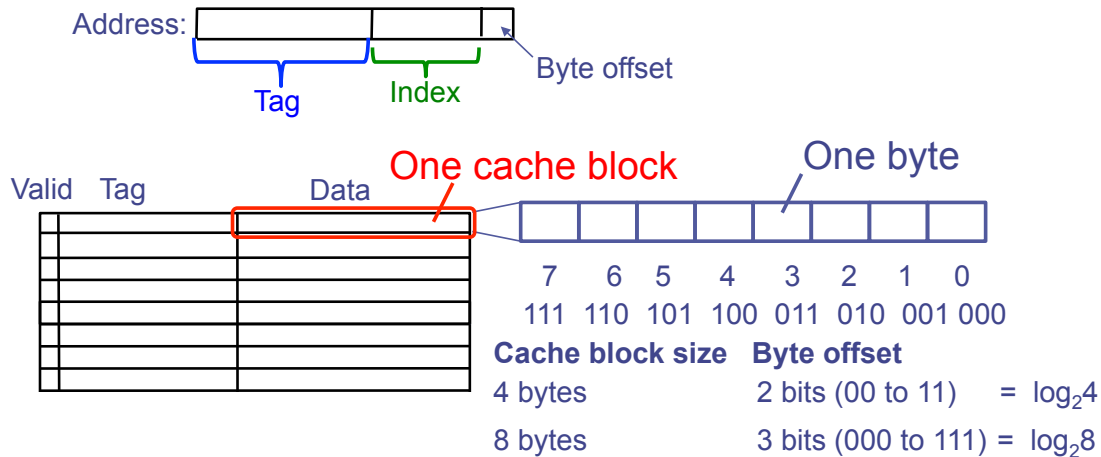
## Step 1: Byte offset and block address

- Byte offset: determined by the size of one cache block
- **Index:** determined by the total number of cache blocks the the cache
- **Tag:** the rest of bits in the address



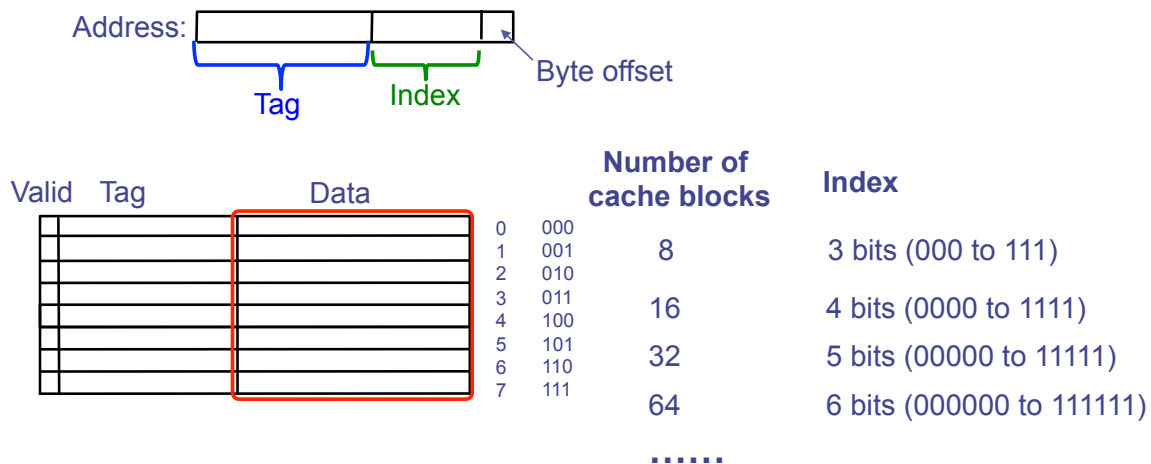
## Step 1: Byte offset and block address

- Byte offset =  $\log_2(\text{Bytes in one cache block})$
- Index: determined by the total number of cache blocks the the cache
- Tag: the rest of bits in the address



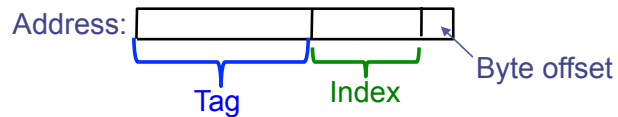
## Step 2: Index region

- Byte offset =  $\log_2(\text{Bytes in one cache block})$
- Index: determined by the total number of cache blocks the the cache
- Tag: the rest of bits in the address



## Step 2: Index region

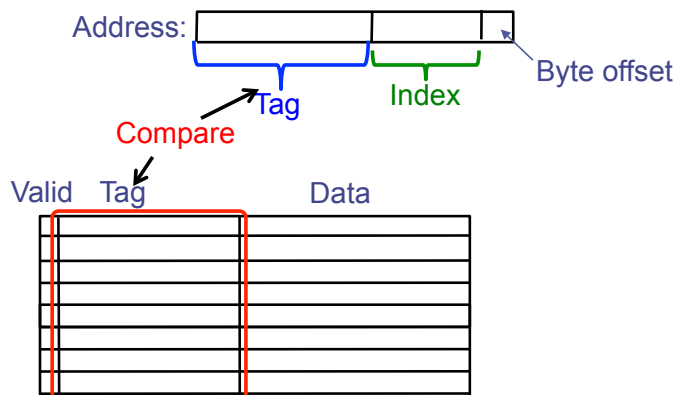
- Byte offset =  $\log_2(\text{Bytes in one cache block})$
- Index =  $\log_2(\text{Number of cache blocks in the cache})$
- Tag: the rest of bits in the address



Valid	Tag	Data		Number of cache blocks	Index
<input type="checkbox"/>			0	000	
<input type="checkbox"/>			1	001	
<input type="checkbox"/>			2	010	8 3 bits (000 to 111)
<input type="checkbox"/>			3	011	
<input type="checkbox"/>			4	100	16 4 bits (0000 to 1111)
<input type="checkbox"/>			5	101	
<input type="checkbox"/>			6	110	32 5 bits (00000 to 11111)
<input type="checkbox"/>			7	111	64 6 bits (000000 to 111111)
				.....	

## Step 3: tag region

- Byte offset =  $\log_2(\text{Bytes in one cache block})$
- Index =  $\log_2(\text{Number of cache blocks in the cache})$
- Tag: the rest of bits in the address





## Put it all together

- Cache block = 1 word (4 bytes),
- cache capacity = 1K words (4K bytes)

