

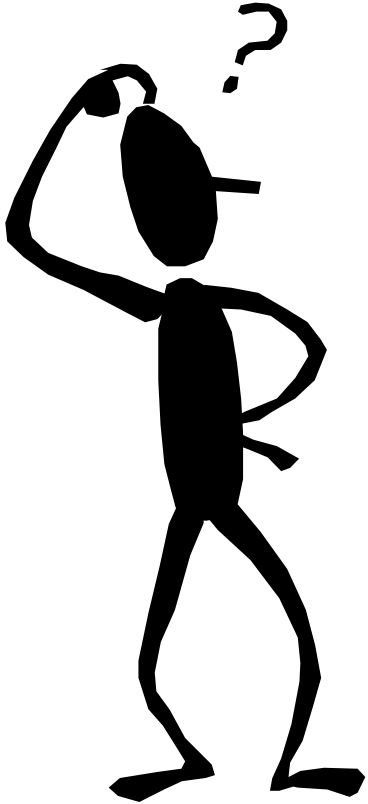
A simple cache design

- Caches are divided into **blocks**, which may be of various sizes.
 - The number of blocks in a cache is usually a power of 2.
 - For now we'll say that each block contains one byte. This won't take advantage of spatial locality, but we'll do that next time.
- Here is an example cache with eight blocks, each holding one byte.

index == row

index	8-bit data
000	
001	
010	
011	
100	
101	
110	
111	

Four important questions

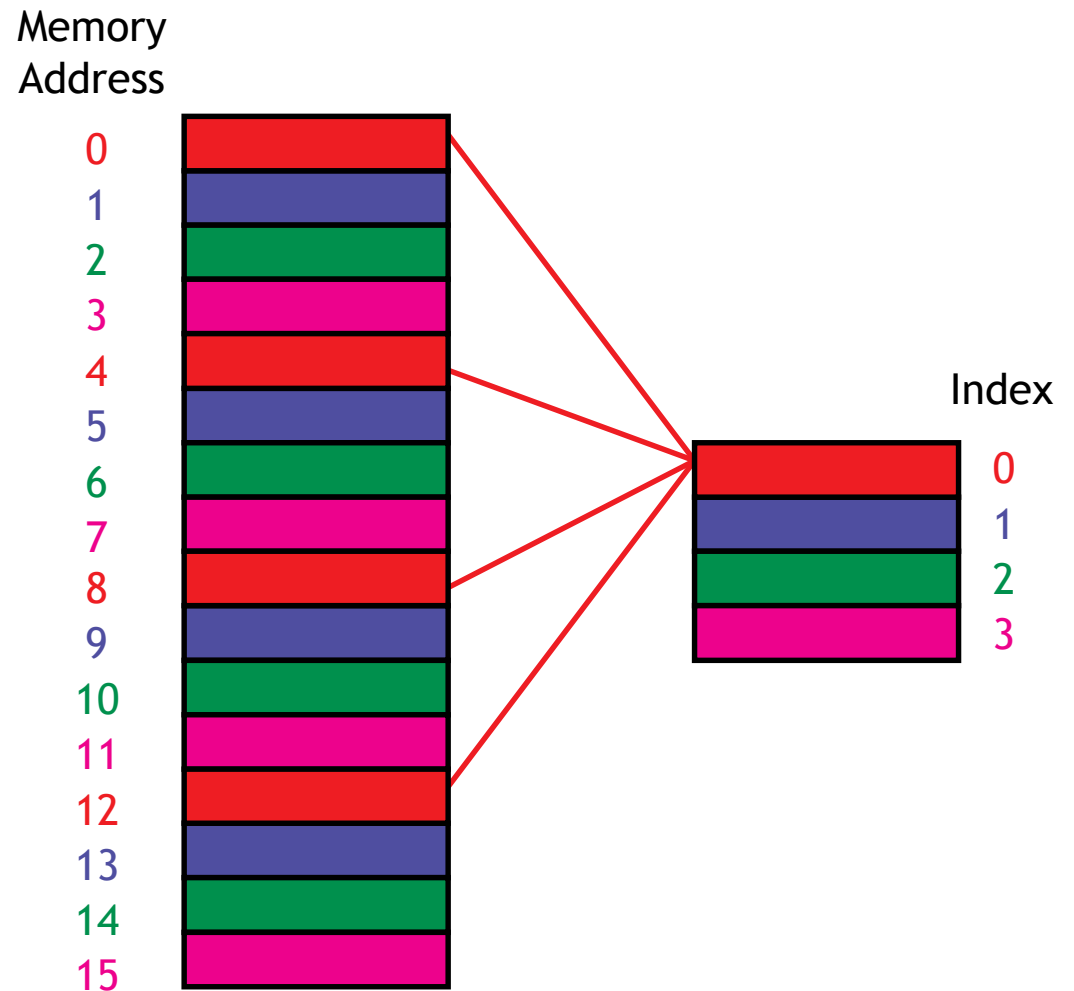


1. When we copy a block of data from main memory to the cache, where exactly should we put it?
2. How can we tell if a word is already in the cache, or if it has to be fetched from main memory first?
3. Eventually, the small cache memory might fill up. To load a new block from main RAM, we'd have to replace one of the existing blocks in the cache... which one?
4. How can *write* operations be handled by the memory system?

- Questions 1 and 2 are related—we have to know where the data is placed if we ever hope to find it again later!

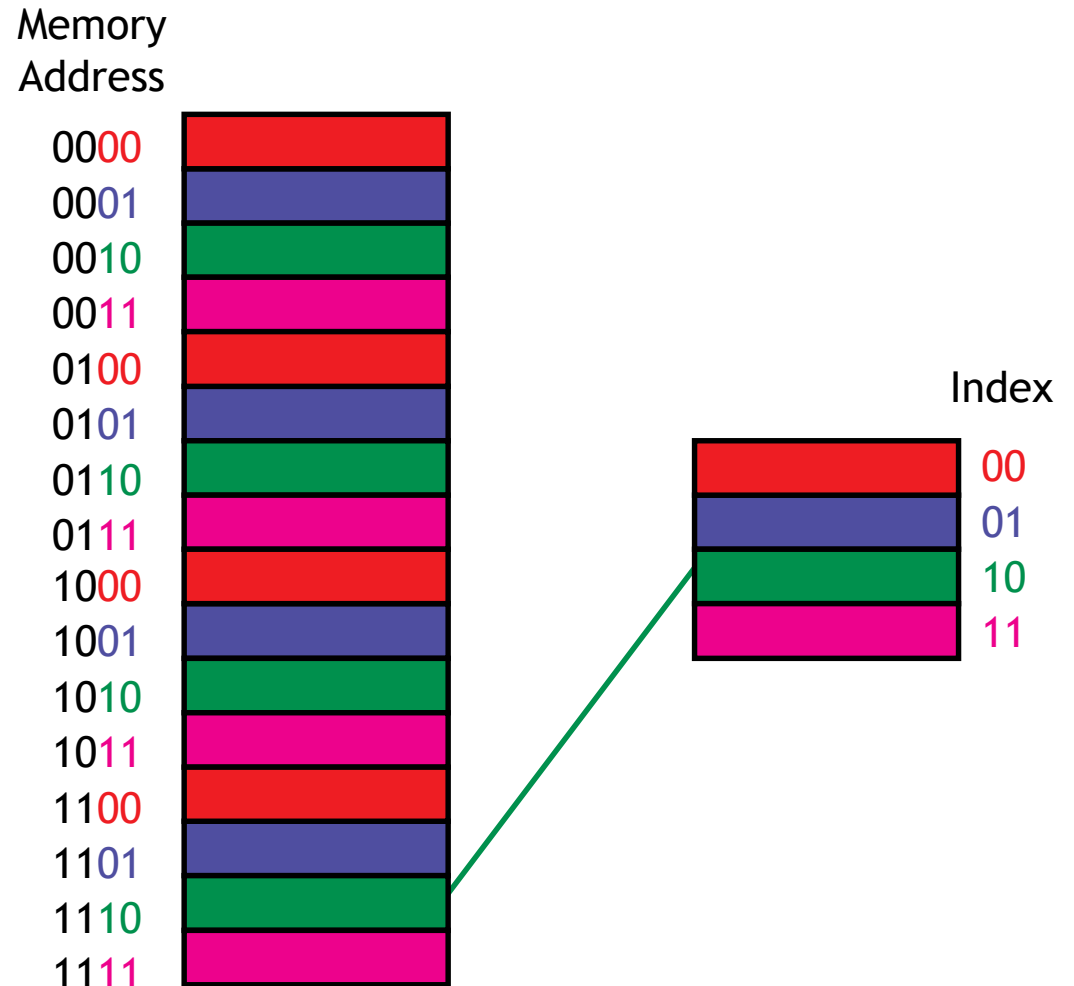
Where should we put data in the cache?

- A **direct-mapped** cache is the simplest approach: each main memory address maps to exactly one cache block.
- For example, on the right is a 16-byte main memory and a 4-byte cache (four 1-byte blocks).
- Memory locations **0**, **4**, **8** and **12** all map to cache block **0**.
- Addresses **1**, **5**, **9** and **13** map to cache block **1**, etc.
- How can we compute this mapping?



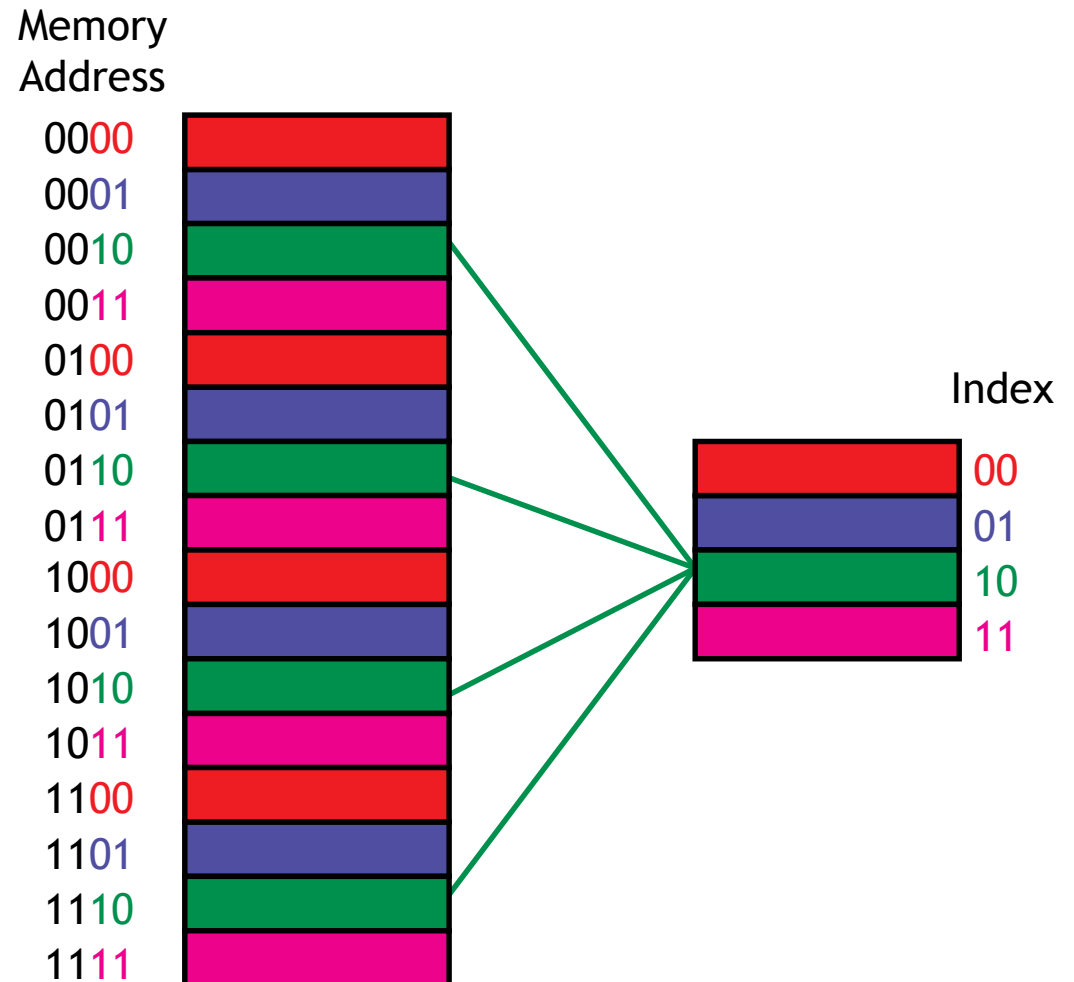
Least-significant bits

- If the cache contains 2^k blocks, then the data at memory address i would go to cache block index $i \bmod 2^k$
- An equivalent way to find the placement of a memory address in the cache is to look at the least significant k bits of the address.
- With our four-byte cache we would inspect the two least significant bits of our memory addresses.
- Again, you can see that address 14 (1110 in binary) maps to cache block 2 (10 in binary).
- Taking the least k bits of a binary value is the same as computing that value $\bmod 2^k$.



How can we find data in the cache?

- The second question was how to determine whether or not the data we're interested in is already stored in the cache.
- If we want to read memory address i , we can use the mod trick to determine which cache block would contain i .
- But other addresses might *also* map to the same cache block. How can we distinguish between them?
- For instance, cache block 2 could contain data from addresses 2, 6, 10 or 14.



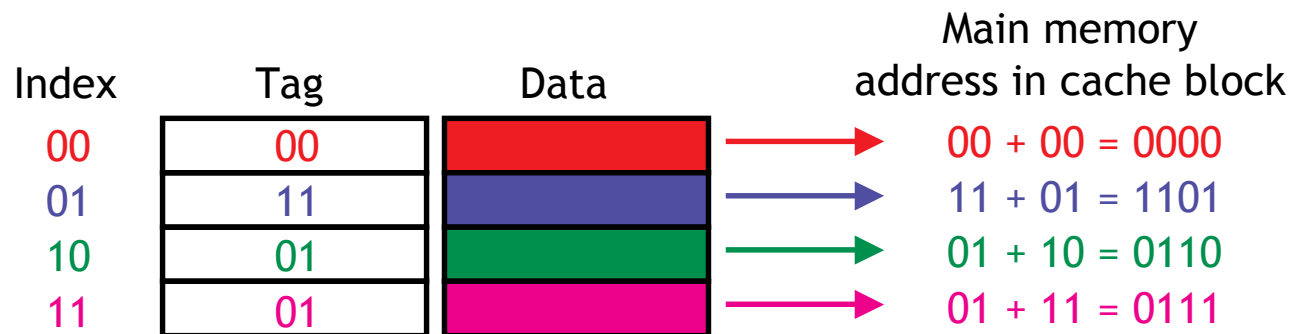
Adding tags

- We need to add **tags** to the cache, which supply the rest of the address bits to let us distinguish between different memory locations that map to the same cache block.



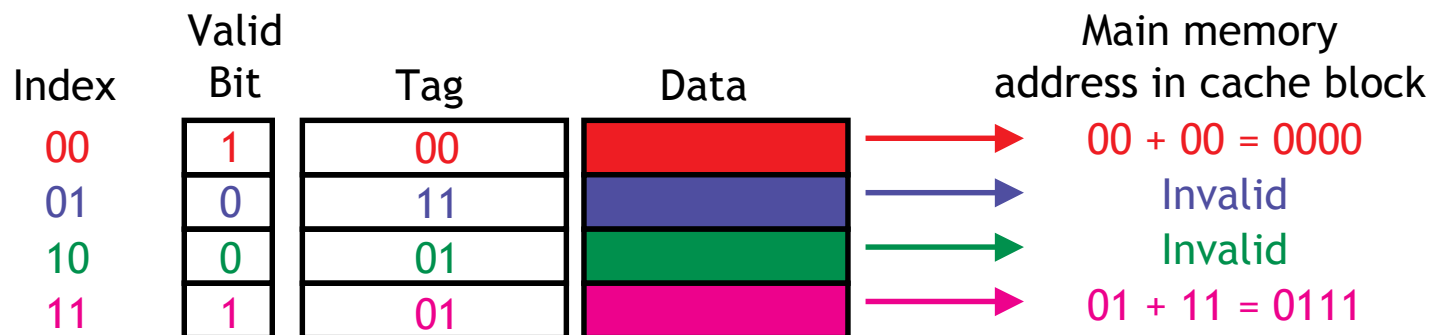
Figuring out what's in the cache

- Now we can tell exactly which addresses of main memory are stored in the cache, by concatenating the cache block tags with the block indices.



One more detail: the valid bit

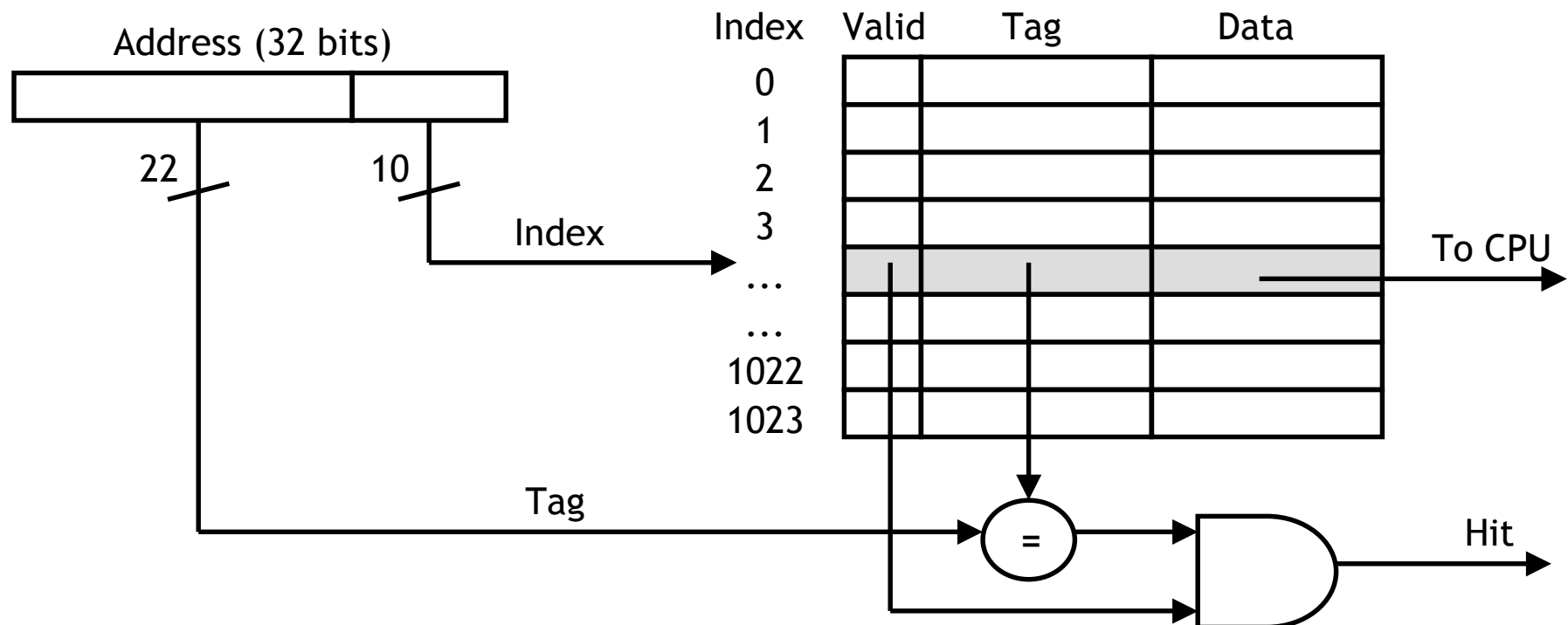
- When started, the cache is empty and does not contain valid data.
- We should account for this by adding a **valid bit** for each cache block.
 - When the system is initialized, all the valid bits are set to 0.
 - When data is loaded into a particular cache block, the corresponding valid bit is set to 1.



- So the cache contains more than just copies of the data in memory; it also has bits to help us find data within the cache and verify its validity.

What happens on a cache hit

- When the CPU tries to read from memory, the address will be sent to a **cache controller**.
 - The lowest k bits of the address will index a block in the cache.
 - If the block is valid and the tag matches the upper $(m - k)$ bits of the m -bit address, then that data will be sent to the CPU.
- Here is a diagram of a 32-bit memory address and a 2^{10} -byte cache.

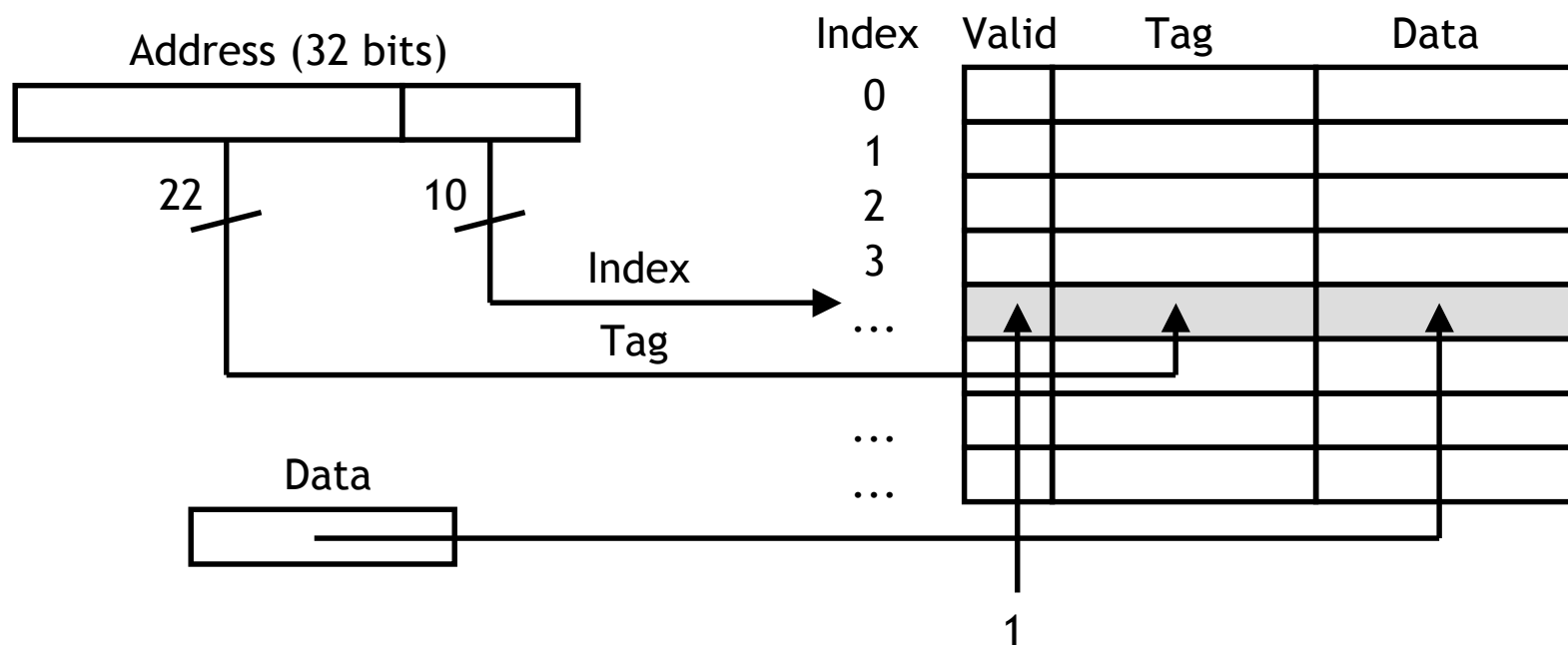


What happens on a cache miss

- The delays that we've been assuming for memories (e.g., 2ns) are really assuming cache hits.
 - If our CPU implementations accessed main memory directly, their cycle times would have to be much larger.
 - Instead we assume that most memory accesses will be cache hits, which allows us to use a shorter cycle time.
- However, a much slower main memory access is needed on a cache miss. The simplest thing to do is to stall the pipeline until the data from main memory can be fetched (and also copied into the cache).

Loading a block into the cache

- After data is read from main memory, putting a copy of that data into the cache is straightforward.
 - The lowest k bits of the address specify a cache block.
 - The upper $(m - k)$ address bits are stored in the block's tag field.
 - The data from main memory is stored in the block's data field.
 - The valid bit is set to 1.

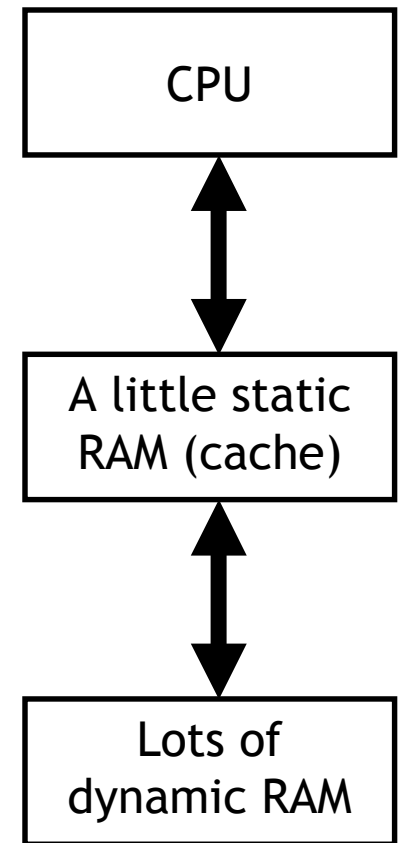


What if the cache fills up?

- Our third question was what to do if we run out of space in our cache, or if we need to reuse a block for a different memory address.
- We answered this question implicitly on the last page!
 - A miss causes a new block to be loaded into the cache, automatically overwriting any previously stored data.
 - This is a **least recently used** replacement policy, which assumes that older data is less likely to be requested than newer data.

Memory System Performance

- To examine the performance of a memory system, we need to focus on a couple of important factors.
 - How long does it take to send data from the cache to the CPU?
 - How long does it take to copy data from memory into the cache?
 - How often do we have to access main memory?
- There are names for all of these variables.
 - The **hit time** is how long it takes data to be sent from the cache to the processor. This is usually fast, on the order of 1-3 clock cycles.
 - The **miss penalty** is the time to copy data from main memory to the cache. This often requires dozens of clock cycles (at least).
 - The **miss rate** is the percentage of misses.



Average memory access time

- The **average memory access time**, or **AMAT**, can then be computed.

$$\text{AMAT} = \text{Hit time} + (\text{Miss rate} \times \text{Miss penalty})$$

This is just averaging the amount of time for cache hits and the amount of time for cache misses.

- How can we improve the average memory access time of a system?
 - Obviously, a lower AMAT is better.
 - Miss penalties are usually much greater than hit times, so the best way to lower AMAT is to reduce the **miss penalty** or the **miss rate**.
- However, AMAT should only be used as a general guideline. Remember that **execution time** is still the best performance metric.

Performance example

- Assume that 33% of the instructions in a program are data accesses. The cache hit ratio is 97% and the hit time is one cycle, but the miss penalty is 20 cycles.

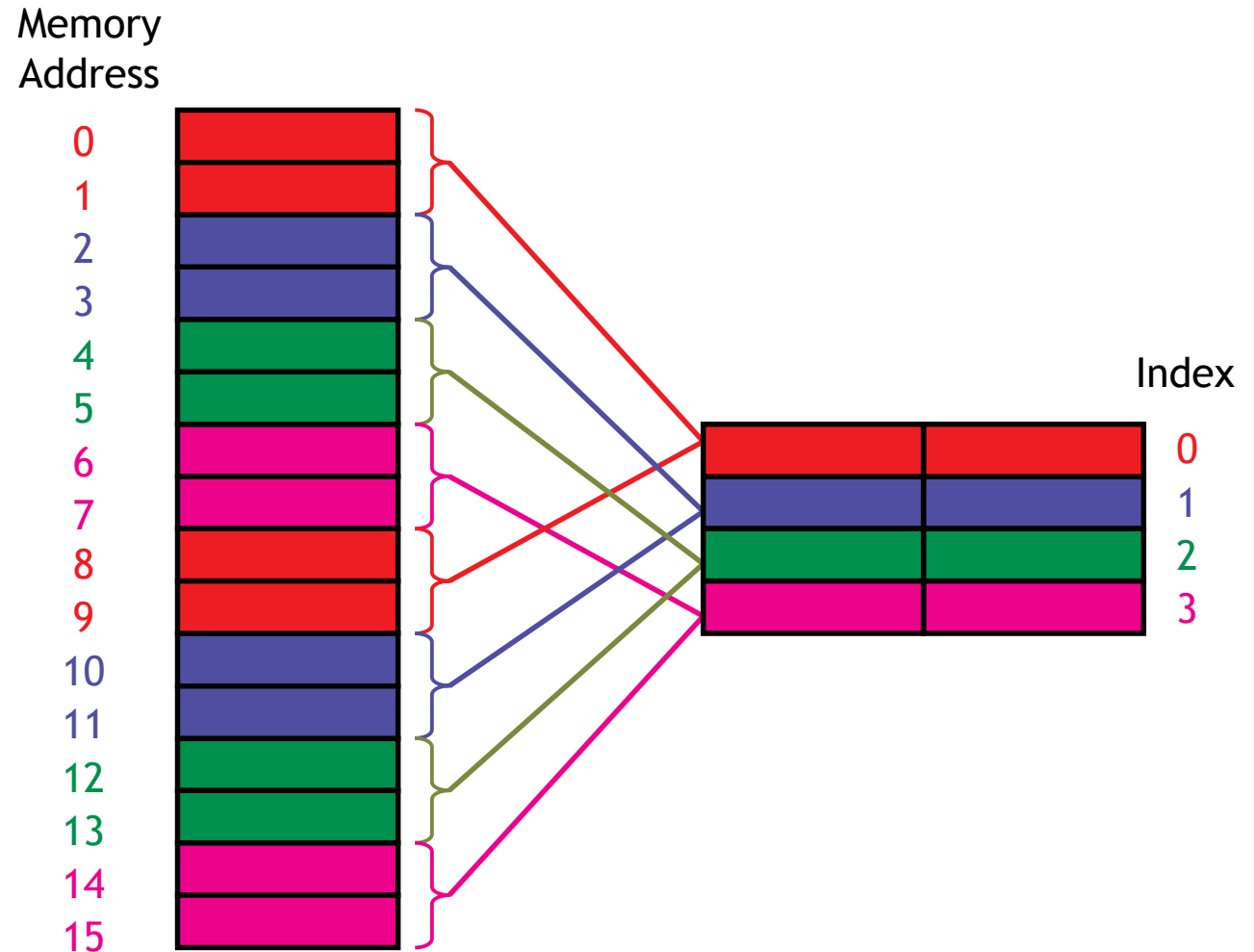
$$\begin{aligned} \text{AMAT} &= \text{Hit time} + (\text{Miss rate} \times \text{Miss penalty}) \\ &= \\ &= \end{aligned}$$

- How can we reduce miss rate?
 - One-byte cache blocks don't take advantage of **spatial locality**, which predicts that an access to one address will be followed by an access to a nearby address.
- What can we do?

Spatial locality

- What we can do is make the cache block size larger than one byte.

- Here we use two-byte blocks, so we can load the cache with two bytes at a time.
- If we read from address 12, the data in addresses 12 and 13 would both be copied to cache block 2.

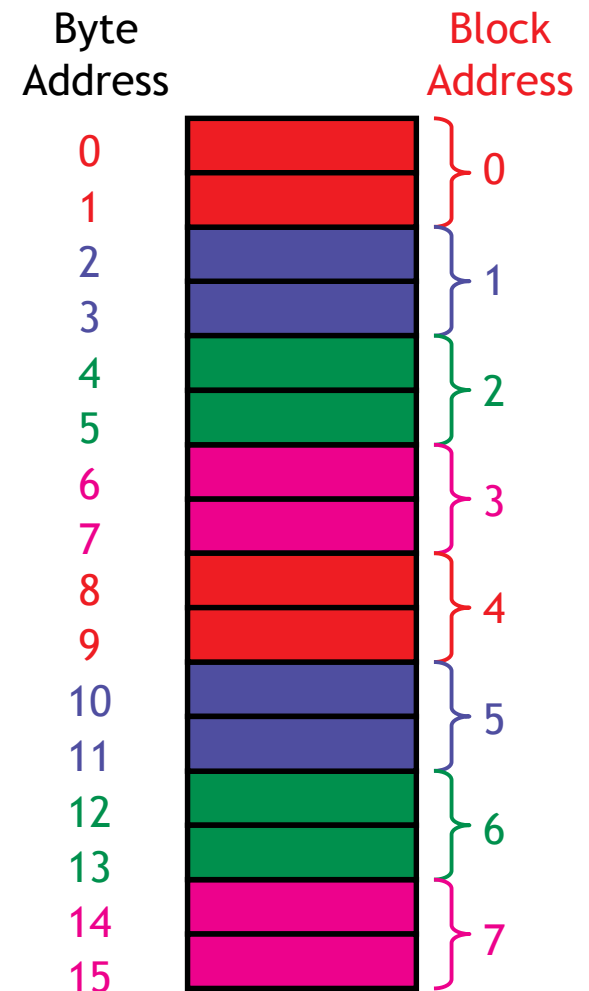


Block addresses

- Now how can we figure out where data should be placed in the cache?
- It's time for **block addresses**! If the cache block size is 2^n bytes, we can conceptually split the main memory into 2^n -byte chunks too.
- To determine the block address of a byte address i , you can do the integer division

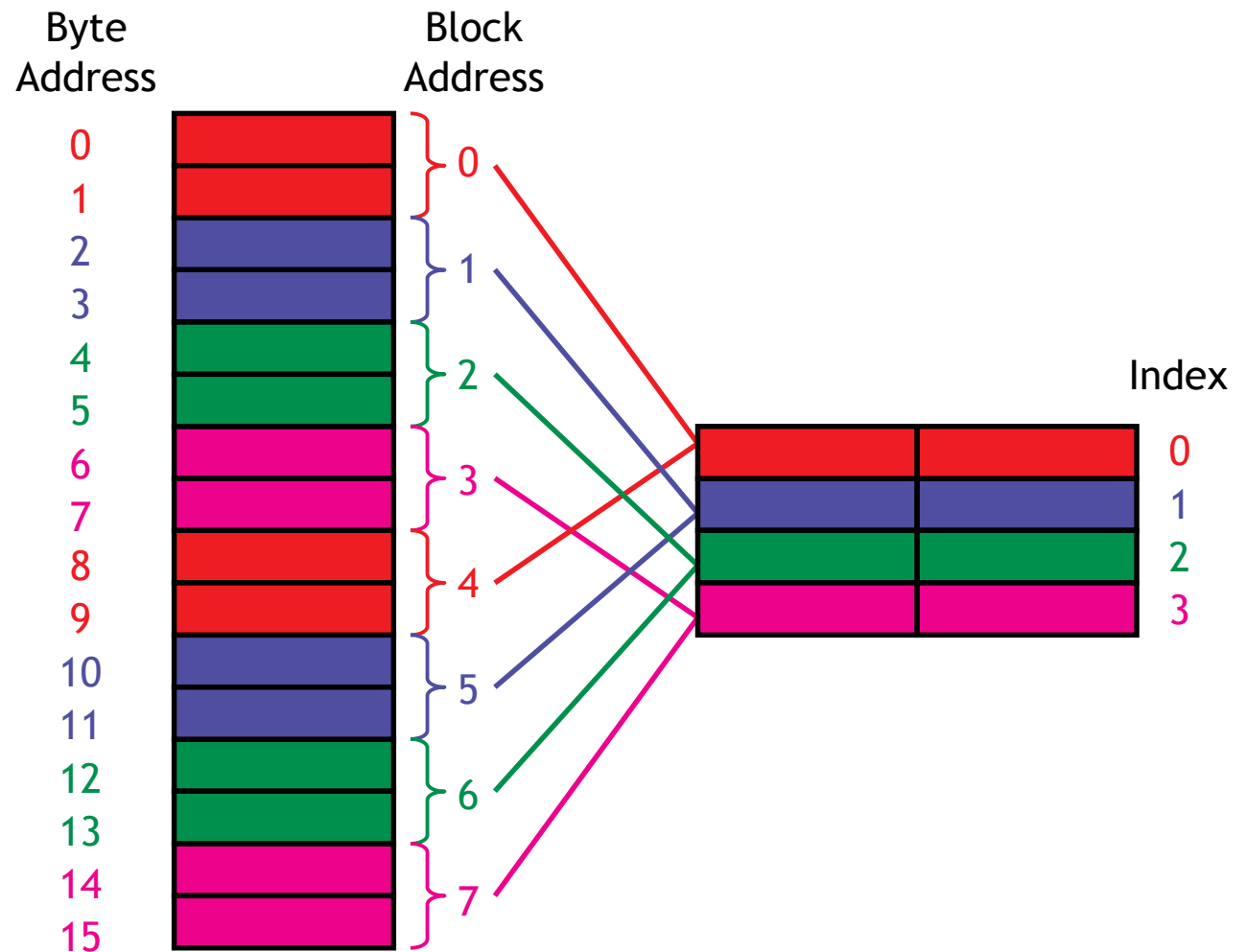
$$i / 2^n$$

- Our example has two-byte cache blocks, so we can think of a 16-byte main memory as an “8-block” main memory instead.
- For instance, memory addresses 12 and 13 both correspond to block address 6, since $12 / 2 = 6$ and $13 / 2 = 6$.



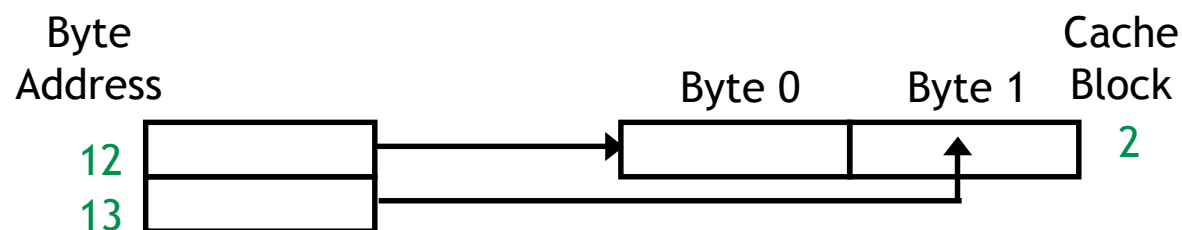
Cache mapping

- Once you know the block address, you can map it to the cache as before: find the remainder when the block address is divided by the number of cache blocks.
- In our example, memory block 6 belongs in cache block 2, since $6 \bmod 4 = 2$.
- This corresponds to placing data from memory *byte* addresses 12 and 13 into cache block 2.



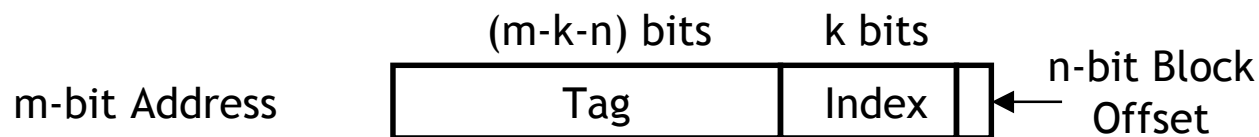
Data placement within a block

- When we access one byte of data in memory, we'll copy its entire *block* into the cache, to hopefully take advantage of spatial locality.
- In our example, if a program reads from byte address 12 we'll load all of memory block 6 (both addresses 12 and 13) into cache block 2.
- Note byte address 13 corresponds to the *same* memory block address! So a read from address 13 will also cause memory block 6 (addresses 12 and 13) to be loaded into cache block 2.
- To make things simpler, byte i of a memory block is always stored in byte i of the corresponding cache block.

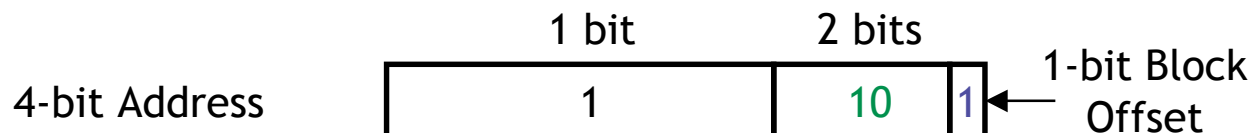


Locating data in the cache

- Let's say we have a cache with 2^k blocks, each containing 2^n bytes.
- We can determine where a byte of data belongs in this cache by looking at its address in main memory.
 - k bits of the address will select one of the 2^k cache blocks.
 - The lowest n bits are now a **block offset** that decides which of the 2^n bytes in the cache block will store the data.



- Our example used a 2^2 -block cache with 2^1 bytes per block. Thus, memory address 13 (1101) would be stored in byte **1** of cache block **2**.



Summary

- Today we studied the basic ideas of **caches**.
 - By taking advantage of **spatial and temporal locality**, we can use a small amount of fast but expensive memory to dramatically speed up the average memory access time.
 - A cache is divided into many **blocks**, each of which contains a **valid bit**, a **tag** for matching memory addresses to cache contents, and the data itself.
- Next, we'll look at some more advanced cache organizations.