

The Basics of Caches

Sat Garcia (sat@cs)

1 Definitions

The following are some basic cache terms that you should be comfortable with.

cache block - The basic unit for cache storage. May contain multiple bytes/words of data.

cache line - Same as cache block. Note that this is not the same thing as a “row” of cache.

cache set - A “row” in the cache. The number of blocks per set is determined by the layout of the cache (e.g. direct mapped, set-associative, or fully associative).

tag - A unique identifier for a group of data. Because different regions of memory may be mapped into a block, the tag is used to differentiate between them.

valid bit - A bit of information that indicates whether the data in a block is valid (1) or not (0).

2 Locating data in the cache

Given an address, we can determine whether the data at that memory location is in the cache. To do so, we use the following procedure:

1. Use the set index to determine which cache set the address should reside in.
2. For each block in the corresponding cache set, compare the tag associated with that block to the tag from the memory address. If there is a match, proceed to the next step. Otherwise, the data is not in the cache.
3. For the block where the data was found, look at valid bit. If it is 1, the data is in the cache, otherwise it is not.

Figures 7.7 and 7.17 from the Patterson & Hennessy textbook give a graphical representation of this process.

If the data at that address is in the cache, then we use the block offset from that address to find the data within the cache block where the data was found.

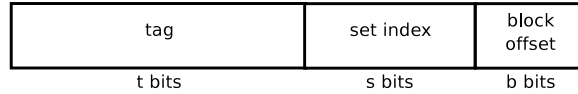


Figure 1: Divisions of the address for cache use.

All of the information needed to locate the data in the cache is given in the address. Fig. 1 shows which parts of the address are used for locating data in the cache. The least significant bits are used to determine the block offset. If the block size is B then $b = \log_2 B$ bits will be needed in the address to specify the block offset. The next highest group of bits is the set index and is used to determine which cache set we will look at. If S is the number of sets in our cache, then the set index has $s = \log_2 S$ bits. Note that in a fully-associative cache, there is only 1 set so the set index will not exist. The remaining bits are used for the tag. If ℓ is the length of the address (in bits), then the number of tag bits is $t = \ell - b - s$.

3 Loading data into the cache

If the requested address is not found in the cache, then it will be brought in from memory and placed there. Along with the requested address, we will also bring in other addresses that are near it to take advantage of spatial locality. To determine which addresses will also be brought in, we find the starting and ending address of the range that will be brought in. The starting address can be found by “zeroing out” the block offset part of the address. For the ending address, we replace the block offset with all 1’s. Note that the size of this range will always be the size of a cache block. The data in that range will be brought in and placed in one of the blocks in the cache.

Depending on the cache organization, there may be multiple places to put data. In a direct mapped cache, there is only one block in the cache where the data can go. In set-associative and fully-associative caches, we must choose among n different blocks. There are n locations in each set for an n -way set-associative cache while an incoming block of data can be

placed in any location in a fully-associative cache. If none of the possible locations are free (i.e. none of the blocks in the cache set have a valid bit set to 0) then we must evict some of the old data. There are various ways to decide which block of data will be replaced but the most popular is the least recently used (LRU) scheme. If using an LRU scheme, then each block in the cache must have LRU bits to track which block is the oldest.

4 Byte vs Word Addressing in Caches

A lot of confusion can arise when talking about byte vs word addressing and its impact on caches. Most modern architectures (MIPS included) use word addressing. Since most current ISAs are still 32-bit, registers have a size of 32 bits (which happens to be the size of a word). It doesn't make much sense to address specific bytes if you are going to have to put them in a word sized register. As we have seen in the Patterson & Hennessy textbook, addresses are shifted left by 2 bits to convert them from a word to a byte address. This translation from words to bytes is hidden from the ISA so if we specify a load from memory address X, we can't tell that it is really being treated as $4 \times X$ in the processor.

When dealing with caches, we have seen that the memory address is split into three parts: block offset, set index, and tag. These three parts are then used to determine whether the address requested is in the cache. The question is, if we are given a word address do we first have to "shift left" by 2 by adding "00" to the end of the address? Adding 2 bits to the end would affect which bits of the original address (i.e. the one without the 00 at the end) are used for the offset, index, and tag. As we will see, this leads to undesirable behavior in the cache, which will become apparent with a small example.

Suppose that your cache has a block size of 4 words. This means that the block offset is the 2 LSBs of your address. If we were to add "00" to the end of every address then the block offset would always be "00." This would mean that we could never access any word other than the one at position 00 in our block.

Always be aware of whether the architecture you are being asked about has word or byte addressing. We will, unless otherwise specified, be assuming a MIPS-like architecture. This means that you will generally be working with word addressing. Make sure that all your parameters (e.g. block size) have the same unit (either word or byte) as the type of addressing you are using.