# 50.002 COMPUTATIONAL STRUCTURES

## INFORMATION SYSTEMS TECHNOLOGY AND DESIGN

# Cache Issues

Natalie Agus (Fall 2018)

## 1   Basic Caching Algorithm

Figure 1 illustrates the basic caching algorithm. We can READ or WRITE to the cache, and it will only access the memory if theres a WRITE or when the READ command has a cache miss (HIT = 0).
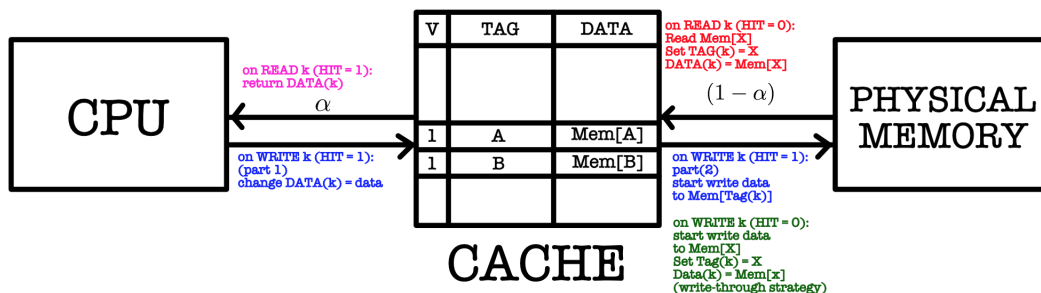


**Figure 1**

Note: $k$ is the row/cache line index of the cache, which is the full address for FA caches, or the lower $k$ bits for DM caches.

## 2   Cache Design Issues

**The access to the main memory is the bottleneck for computer performance**. Cache is used to reduce the frequency of access to the main memory whenever possible. There are several design issues for cache:

1. Associativity : how many different addresses can be stored in the cache

2. Replacement strategy

3. Block size

4. Write strategy : when to write from cache to memory

# 3    Pros and Cons of FA and DM Caches

|  | **FA Cache** | **DM Cache** |
|---|---|---|
| **Tag Content** | All address bits | Lower $t$ address bits |
| **Tag Index** | None | Higher $k$ address bits |
| **Data** | Mem[A] | Mem[A] |
| **Performance** | The gold standard on how well a cache should perform | Performs slower than FA cache on average |
| **Contention** | Does not have address contention. Address-data can fit on any cache line. | Has address contention. **The probability of contention is inversely proportional to cache size.** |
| **Cost** | Expensive, many boolean circuits for comparators | Cheap, only one comparator in a cache |
| **Replacement Strategy** | Yes: LRU, FIFO, Random | No |
| **Mapping** | No available cache line mapping, must check all cache lines and compare the entire address | Fast mapping, lower k-bits for cache line index, and then compare upper t-bits for tag content |
| **Application** | Good when cache size is small, less important when cache size is large | Bad when cache size is small (more contention), good when cache size is large |

**Table 1**

# 4    N-way set associative cache

## 4.1    Purpose

Finding the balance point:

1. Although DM cache is cheap, it suffers from contention problem.

2. Contention mostly occurs within a certain block of addresses (called independent hot-spots)

3. This is due to locality of reference in each of different address range.

4. Hence, some associativity (and not full associativity) is needed.

5. We should allow each location (lower k-bits directly mapped) to be stored in a restricted set of cache lines (see Fig 2 to know what is a 'set' and what is a 'cache line').

## 4.2  Method

One solution is to build N-way set associative cache. See Figure 2.

1. We have $N$ DM caches.

2. The cells in the same 'row' marked in red is called to belong in the same **set**.

3. The cells in the same 'column' of DM caches marked in blue is said to belong in the same **cache line**.

4. Given the same k-bits lower address, it can be stored in any of the N cache lines in the **same set**.

5. However a different combination of k-bits lower address will have to be direct mapped on different set.

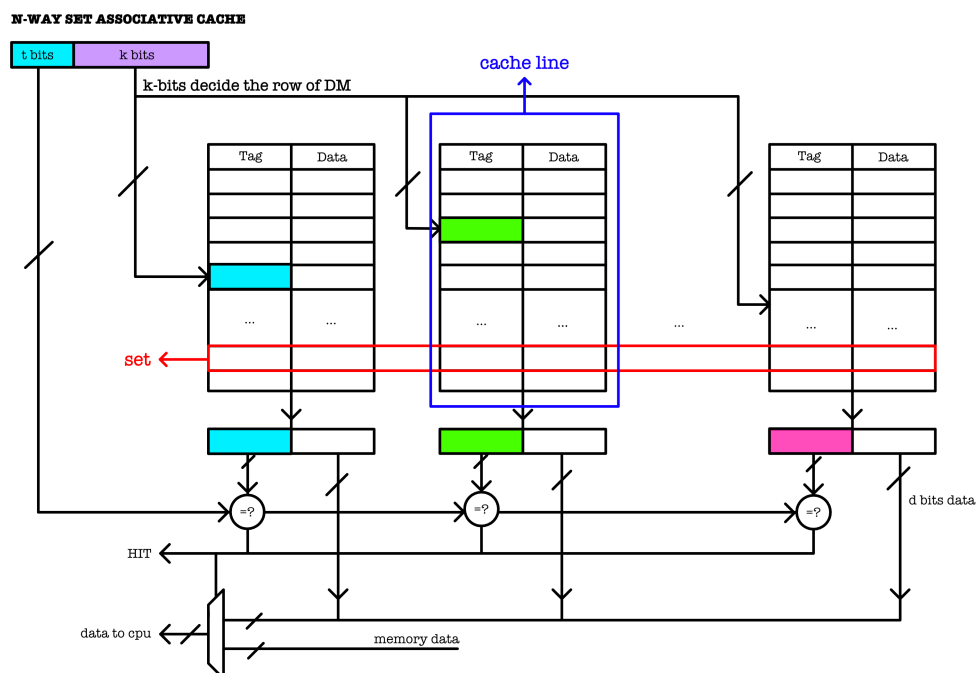6. Lookup operation: Parallel operation of N direct-mapped cache, each with $2^t$ lines in the cache line.

**Figure 2**

## 4.3  Replacement Strategy

The replacement strategy is required to find out which of the N-cache lines in the same set can be replaced when there's a cache miss to that particular address request. There are three common strategies:

1. **LRU : Least Recently Used**.

   - Replace least recently used item, good when N is small
   - Need to keep ordered list of N items ($N!$ orderings),
   - Overhead is $O(\log_2 N!) = O(N \log_2 N)$ "LRU bits" per set
   - Plus complex logic to re-order the list after every cache access

2. **FIFO / LRR : Least Recently Replaced**

   - Replace oldest item
   - Overhead is $O(\log_2 N)$ bits/set, because one needs just one pointer (indicator bits) that tells us the oldest item within each set.

3. **Random**

   - Use pseudo random generator to get reproducible behavior
   - Good when N is huge

There's no one best / winning replacement strategy. One replacement strategy can be better than the other depending on cases.

# 5   Increasing Block Size

Sometimes we would want to enlarge each line in cache (more data per tag), especially if there's high locality of reference:  Here we have:
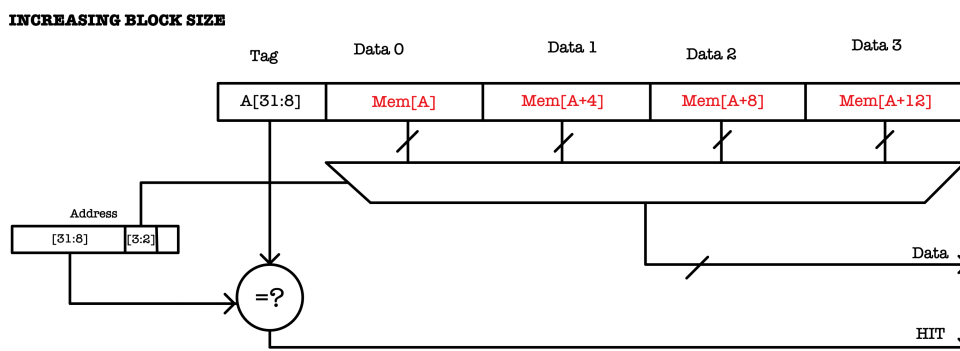


**Figure 3**

1. Blocks of $2^B$ words per row, in this case in Figure 3 $B = 2$.

2. **Pros**: locality of reference means there's a high likelihood that the words from the same block will be required together

3. **Cons**: risk of fetching unaccessed words

4. Figure 3 uses byte addressing (sometimes one can use word addressing too, especially for problem sets to make it easier for our computation), so recall that the lower two bits of address is always 0.
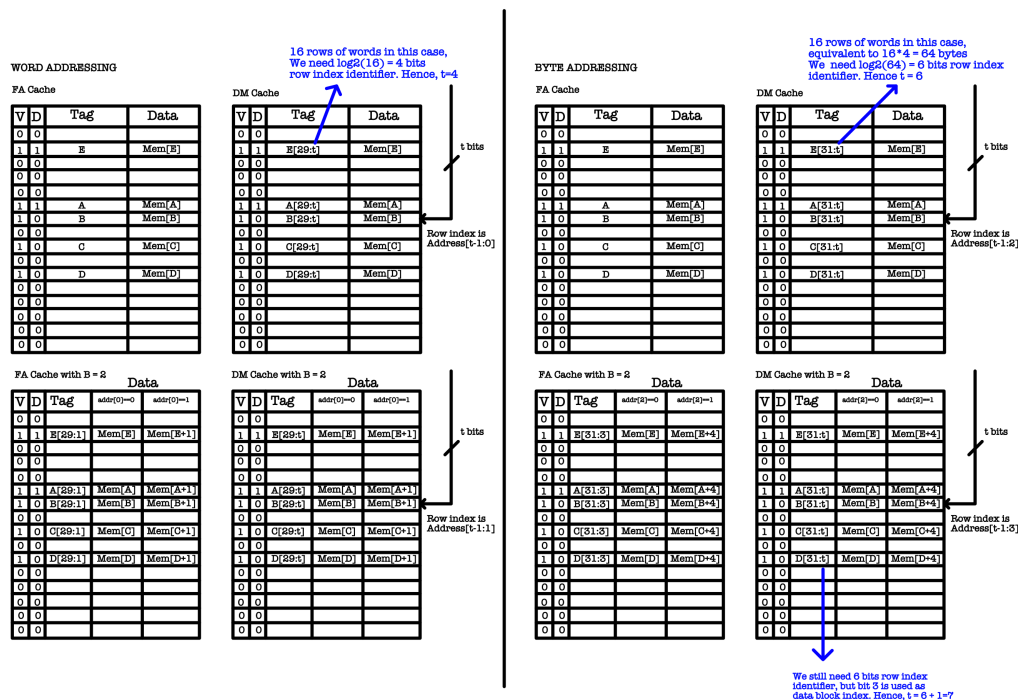
# 6 Byte VS Word Addressing



**Figure 4**

Byte addressing is common in practice. However for the ease of calculation in practice questions and problem sets, now we do word addressing. Figure 4 above illustrates the difference between byte addressing that we already know and word addressing.

# 7 Valid and Dirty Bit

In Figure 6, we add two extra bits of storage in the cache : $V$ and $D$.

## 7.1 V: Valid Bit

The valid bit indicates that the particular cache row (also called cache line, but it is a different graphical representation from 'cache line' in the N-way set associative cache) contains data from memory and not an empty or a redundant value. We only check the contents of cache lines with valid bit = 1. **The most obvious benefit of having a valid bit is that we can "flush" the cache real quick by setting V = 0.**

## 7.2 D: Dirty Bit

The dirty bit is set to 1 iff the CPU writes to cache and it hasn't been stored to the memory (memory is outdated).

# 8 Cache Writes

Most memory accesses are READs, but we can handle WRITEs in three different ways. All three methods require CPU to write data to cache first.

## 8.1 Write-through

CPU writes are done in the cache first by setting TAG = Addr, and Data = new Mem[Addr] in an available cache line, but also written to the main memory immediately. This **stalls** the CPU until write to memory is complete, but memory always holds the "truth", and is never oudated.

## 8.2 Write-behind

CPU writes are also cached, and write to the main memory is immediate but it is buffered or pipelined. CPU keeps executing next instructions while writes are completed (in order) in the background.

## 8.3 Write-back

CPU writes are also cached, but not immediately written to the main memory. Memory contents can be "stale". Typically CPU will write to the main memory only if the data in cache line needs to be replaced and that this data has been changed or is new. **This requires the dirty bit** in the cache.