

## 1. Pre-check

This part is designed as a check to help you determine whether you understand the concepts covered in class. Answer "True" or "False" to the following questions and include an explanation:

- 1.1. For the same cache size and block size, a 4-associative cache will have fewer index bits than a direct-mapped cache.
  
- 1.2. When the cache is full, any cache miss that occurs is a "capacity miss".
  
- 1.3. Increasing the size of the cache by adding more lines always improves the success rate.

## 2. Understanding T/I/O

Manipulating data with caches relies on ingenious reading/decomposition of the memory address into three different bit fields:

- **Index** – The index of the cache line where the memory block will be placed  
$$\text{Width (in bits)} = \log_2(\text{number of cache lines})$$
- **Offset** – The offset position of the byte in the memory block  
$$\text{Width (in bits)} = \log_2(\text{size in bytes of a line of cache})$$
- **Tag** – Used to distinguish between different blocks of memory that may use the same cache line (i.e., same Index number).  
$$\begin{aligned} \text{Width (in bits)} &= \text{Width of memory address (in bits)} \\ &\quad - \text{Width(Index)} - \text{Width(Offset)} \end{aligned}$$

Thus, one can verify the following identity:

$$\begin{aligned} \log_2(\text{Memory size in bytes}) &= \text{Width of memory address in bit} \\ &= \text{Width(Tag)} + \text{Width(Index)} + \text{Width(Offset)} \end{aligned}$$

Another useful equality to remember is:

$$\text{Size of Cache} = \text{Size of Memory Block} \times \text{Number of Cache lines}$$

**Note:** In the lecture, the "Offset" field is further broken down into two sub-parts: The "word" field indicating the position of the word in a memory block and the "byte" field giving the index of the byte in the current word. In this tutorial, we limit ourselves to the "Tag / Index / Offset" decomposition of the memory address.

2.1. Assume a direct-mapped cache with a capacity of 32 bytes and a cache line of 8 bytes. For a 32-bit memory address, which bits correspond to the "Offset" field?

2.2. Which bits should we check to find the line to use in cache?

2.3. What about the "Tag" field?

2.4. Indicate for each of the following memory accesses whether it is a cache hit (H), a cache miss due to invalid entry (M), or a cache miss due to Replacement (R). **Hint:** Drawing a sketch of the cache can help you see the overrides more clearly.

Address	T	I	O	Hit, Miss, Replacement
0x00000004				
0x00000005				
0x00000068				
0x000000C8				
0x00000068				
0x000000DD				
0x00000045				
0x00000004				
0x000000C8				

### 3. Cache Associativity

To minimise cache misses due to “insufficient capacity” in a direct-mapped cache, we could simply increase the size of the cache line. This will ensure a better exploitation of the spatial locality principle, but the number of cache misses due to, for example, repetitive function calls will not be reduced – increasing the size of the cache line in a direct-mapped cache does not ensure a better exploitation of the temporal locality principle.

To take this principle into account in the cache design, we allow to associate (hence the name associativity) to a memory block several possible locations on a cache line. Thus, a cache is said to be N-associative when each memory block can go to N distinct locations on a cache line. For an N-associative cache, we have the rule:

$$N \times \text{Number of cache lines} = \text{Size of cache (in number of memory blocks)}$$

- 3.1 Assume a 2-associative cache with a Least Recently Used (LRU) replacement policy and an 8-bit addressable memory. The size of the cache is 32 bytes and the size of a cache line is 8 bytes. Indicate for each of the following memory accesses whether it is a cache hit (H), a cache miss (M) due to invalid entry, or a cache miss due to Replacement (R).

Address	T	I	O	Hit, Miss, Replacement
0x04				
0x05				
0x68				
0xC8				
0x68				
0xDD				
0x45				
0x04				
0xC8				

- 3.2 What is the hit rate for the previous question?

## 4. The three causes of cache misses

Review questions 2.4 and 3.1 and classify each cache miss as one of the 3 types of misses described below:

- **Compulsory misses:** A miss that must occur when a memory block is referenced for the first time. We can reduce first reference misses by having longer cache lines (i.e. larger blocks). We can also prefetch blocks in advance by using a special circuit to predict the next blocks that are likely to be required.
- **Capacity misses:** These misses are caused by the fact that the cache cannot hold all the blocks referenced during the execution of the program. The number of these misses can be reduced by increasing the size of the cache.
- **Conflict misses:** These misses occur in addition to the two previous types. A block already loaded into the cache is ejected from the cache because another block with the same index number is required by the program. The number of these misses can be reduced by increasing the associativity of the cache.

## 5. Source Code Analysis

Consider the C code below, running on a system with 1MB of memory and a 16KB direct cache organized in 1KB blocks.

```
1 #define NUM_INTS 8192    // pow(2,13)
2 int A[NUM_INTS];        // Assume A is allocated at address 0x10000
3 int i, total = 0;
4 for (i = 0; i < NUM_INTS; i += 128) {
5     A[i] = i;
6 }
7 for (i = 0; i < NUM_INTS; i += 128) {
8     total += A[i];
9 }
```

5.1 What is the width (in bits) of this system's memory address?

5.2 Give the "T/I/O" decomposition of the address

5.3 Calculate the hit rate for line 5 of the C code

5.4 Calculate the hit rate for line 8 of the C code

## 6. Average Memory Access Time (AMAT)

The “Average Memory Access Time” is defined by the following rule:

$$AMAT = hit\ access\ time + miss\ rate \times miss\ penalty$$

*Hit access time* = access time to data found in cache

*Miss rate* = number of cache misses ÷ number of cache accesses

For a hierarchical cache system, there are two miss rate metrics for each cache level:

- **Global miss rate:** The number of accesses to RAM divided by the total number of accesses to the **entire** cache system.
- **Local miss rate** for level L: The number of cache misses at level L divided by the number of total accesses to **this** level.

6.1 In a two-level cache system, there were 20 misses out of a total number of 100 accesses. What is the global miss rate?

6.2 If the L1 (i.e. level 1) cache has a 50% miss rate, what is the local miss rate of the L2 cache?

For the following questions, consider a system with the following characteristics:

- An L1 cache with a hit access time of 2 clock cycles and a local miss rate of 20%.
- An L2 cache with a hit access time of 15 clock cycles and a global miss rate of 5%.
- A main memory with an access time of 100 clock cycles.

6.3 What is the local miss rate for the L2 cache?

6.4 Give the Average Memory Access Time (AMAT) of the system

6.5 We would like to reduce the system's AMAT to 8 clock cycles or less by adding a level 3 cache. If the L3 cache has a local miss rate of 30%, what is the highest response time this cache should have?