

1 Pre-Check

This section is designed as a conceptual check for you to determine if you conceptually understand and have any misconceptions about this topic. Please answer true/false to the following questions, and include an explanation:

- 1.1 We cannot use a 1KB cache in a 32-bit system because it's too small and cannot contain all possible addresses.

False. The purpose of the cache is not to hold every possible piece of memory at the same time, but rather to hold some parts of it only, so a 1KB cache is not "too small".

- 1.2 If a piece of data is both in the cache and in memory, reading it from cache is faster than reading from memory.

True. The cache is smaller and faster than memory.

- 1.3 Caches see an immediate improvement in memory access time at program execution.

False. A cache starts off 'cold', and required loading in values in blocks at first directly from memory, forcing compulsory misses. This can be somewhat alleviated by the use of a hardware prefetcher, that uses the current pattern of misses to predict and prefetch data that may be accessed later on.

- 1.4 Increasing cache size by adding more blocks always improves (increases) hit rate for all programs.

False. Whether this improves the hit rate for a given program depends on the characteristics of the program. As an example, it is possible for a program that only consists of a loop that runs through an array once to have each access be separated by more than one block (say, the block size is 8B, but we have an integer array and accessing every fourth element, so our access are separated by 16B). This makes every miss a compulsory miss, and there is no way for us to reduce the number of compulsory misses just by adding more blocks to our cache.

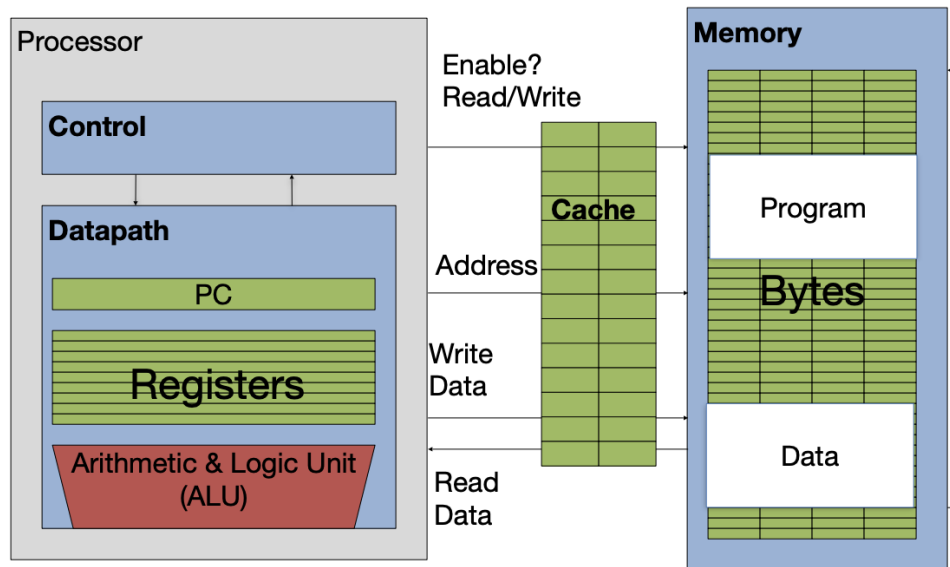
- 1.5 Decreasing block size to increase the **number** of blocks held by the cache improves the program speed for all programs.

False. This question is similar to the one above, in that the answer to it depends on the program that is running. If we have a program with a for loop that loops through continuous memory (like an array), having a bigger blocks size and fewer blocks might be helpful, as the single blocks will holds more continuous data. For example, lets say cache A has 10 lines and a block size of 8 bytes, while cache B has 20 lines with a a block size of 4 bytes and the array we loop through has 80

characters. Cache A in this case will have 10 cache misses and 70 hits, while Cache B will have 20 misses and 60 hits.

2 Understanding T/I/O

We use caches to make our access to data faster. When working with main memory (RAM), the main problem faced is the fact that access to the main memory is very slow. In fact, modern processors take about 100 instructions cycles or more to access the main memory, meaning memory accesses become the bottleneck of our programs. Caches help fix this problem for us - they hold a portion of the data in main memory, that we might access again later on. They are closer to the processor in the memory hierarchy, and thus accessing a cache is much faster than accessing the main memory.



As seen above, the access to cache is the middle step between the CPU asking for a memory bit, and the actual main memory access - if the data is not found in the cache, only then is main memory accessed. This way unnecessary trips to main memory are avoided. One important detail is that caches are much smaller in size than main memory - this is why we have to be efficient in what we hold in caches. When we are saving data in caches, we need to be as efficient as possible. In order to do this, we make use of locality. We have two different kinds of locality to consider.

Temporal Locality: If we have accessed a piece of information recently, it is possible that we will access it again. So, we hold this data in the cache.

Spatial Locality: If we have accessed a memory location recently, it is probable that we will access the neighbouring addresses as well. So, we also keep the neighbouring addresses within the cache. An example is array accesses - if we access the 0th element of an array, it is probable we will also access the 1st one.

Note that caches hold the data in blocks that have a size equal to the block size of the cache.

When working with caches, we have to be able to break down the memory addresses we work with to understand where they fit into our caches. There are three fields:

Tag - Used to distinguish different blocks that use the same index. Number of bits: ($\#$ of bits in memory address) - Index Bits - Offset Bits

Index - The set that this piece of memory will be placed in. Number of bits: $\log_2(\#$ of indices)

Offset - The location of the byte in the block. Number of bits: $\log_2(\text{size of block})$

Given these definitions, the following is true:

$$\log_2(\text{memory size}) = \text{address bit-width} = \# \text{ tag bits} + \# \text{ index bits} + \# \text{ offset bits}$$

Another useful equality to remember is:

$$\text{cache size} = \text{block size} * \text{num blocks}$$

- 2.1 Assume we have a direct-mapped byte-addressed cache with capacity 32B and block size of 8B. Of the 32 bits in each address, which bits do we use to find the index of the cache to use?

We can determine the number of index bits we need from the number of sets our cache has. Since our cache is direct-mapped, the number of sets is the same as the number of blocks, so we just need to figure out how many blocks our cache has. Using the equality from above, we see that $\text{num blocks} = \text{cache size}/\text{block size}$, so our cache has $32/8 = 4$ blocks. We need $\log_2(4) = 2$ bits to differentiate the 4 blocks, so we have 2 index bits.

In order to determine where exactly the index bits are, we need to calculate the number of offset bits and tag bits we have. The number of offset bits is just dependent on the block size, so since our blocks are size 8B, we need $\log_2(8) = 3$ bits to differentiate the 8 bytes in the block, so we have 3 offset bits.

our offset bits take up the least significant bits, with the index bits being the set of next most significant bits. Denoting the most significant bit (MSB, on the left) as 31 and the least significant bit (LSB, on the right) as 0, having 3 offset bits means our index bits start at bit 3, and thus we use bits 3 and 4 as the index bits.

- 2.2 Which bits are our tag bits? What about our offset?

The offset (in this case) is the 3 least significant bits, so reusing the convention from the previous question, the offset bits are bits 0, 1, and 2. Our tag is the remaining high-order bits, so our tag bits are bits 5-31.

- 2.3 Classify each of the following byte memory accesses as a cache hit (H), cache miss (M), or cache miss with replacement (R). Tip: Drawing out the cache can help you see the replacements more clearly.

Address	T/I/O	Hit, Miss, Replace
0x00000004		
0x00000005		
0x00000068		
0x000000C8		
0x00000068		
0x000000DD		
0x00000045		
0x000000CF		
0x000000F3		

Ignore miss types (compulsory/conflict/capacity) until Q4.

```

0x00000004    Tag 0, Index 0, Offset 4: M, Compulsory
0x00000005    Tag 0, Index 0, Offset 5: H
0x00000068    Tag 3, Index 1, Offset 0: M, Compulsory
0x000000C8    Tag 6, Index 1, Offset 0: R, Compulsory
0x00000068    Tag 3, Index 1, Offset 0: R, Conflict
0x000000DD    Tag 6, Index 3, Offset 5: M, Compulsory
0x00000045    Tag 2, Index 0, Offset 5: R, Compulsory
0x000000CF    Tag 6, Index 1, Offset 7: R, Conflict
0x000000F3    Tag 7, Index 2, Offset 3: M, Compulsory

```

Note that the M and R distinction here is for student understanding, and that the cache doesn't behave differently for these cases.

3 The 3 C's of Cache Misses

In order to evaluate cache performance and hit rate, especially with determining how effective our current cache structure is, it is useful to analyze the misses that do occur, and adjust accordingly. Below, we categorize cache misses into three types:

- (I) Compulsory: A miss that must occur when you bring in a certain block for a first time, hence "compulsory". Reduce compulsory misses by having longer cache lines (bigger blocks), which bring in the surrounding addresses along with our requested data. Can also pre-fetch blocks beforehand using a hardware prefetcher (a special circuit that tries to guess the next few blocks that you will want).
- (II) Conflict: Occurs if the block was fetched before, but evicted while the cache was not full. Increasing the associativity or improving the replacement policy would remove the miss. Note that Conflict misses tend to occur with inefficient cache usage, compared to Capacity misses.
- (III) Capacity: Occurs if the block was fetched before, but evicted while the cache was full. Capacity misses are independent of the associativity of your cache.

The only way to remove the miss is to increase the cache capacity.

3.1 Go back to question 2 and classify each miss (M) and replacement (R) as one of the 3 types of misses described above.

See solutions for Q3.

4 Cache Associativity

Previously, we focused on Direct-Mapped caches, in which blocks map to specifically one slot in our cache. This is good for quick replacement and finding out block, but not good for spatial efficiency!

This is where we bring associativity into the matter. We define associativity as the number of slots a block can potentially map to in our cache. Thus, a Fully-Associative cache has the most associativity, meaning every block can go anywhere in the cache. Our Direct-Mapped cache, on the other hand, has the least, being only 1-way set associative.

For an N -way associative cache, the following are true:

$$N * \# \text{ sets} = \# \text{ blocks}, \quad \text{Index bits} = \log_2(\# \text{ sets})$$

- 4.1 Here's some practice involving a 2-way set associative cache. This time we have an 8-bit address space, 8 B blocks, and a cache size of 32 B. Classify each of the following accesses as a cache hit (H), cache miss (M) or cache miss with replacement (R). For any misses, list out which type of miss it is (Compulsory, Conflict, or Capacity). Assume that we have an LRU replacement policy (in general, this is not always the case).

Address	T/I/O	Hit, Miss, Replace
0b0000 0100		
0b0000 0101		
0b0110 1000		
0b1100 1000		
0b0110 1000		
0b1101 1101		
0b0100 0101		
0b0000 0100		
0b0011 0000		
0b1100 1011		
0b0100 0010		

Since our cache is 2-way set associative, there are 2 blocks in a set. Given the cache size and the block size, we have $32 / 8 = 4$ blocks. Thus, there are $4 / 2 = 2$ sets in our cache. We need $\log_2(2) = 1$ bit to differentiate the 2 sets, so we have 1 index bit. Our block size of 8 B means we have $\log_2(8) = 3$ offset bits, and that the rest of our bits are our tag bits. Therefore, our TIO breakdown means bits 0, 1, and 2 are our offset bits, the only index bit is bit 3, and bits 4-7 being the tag bits.

```

0b0000 0100    Tag 0000, Index 0, Offset 100 - M, Compulsory
0b0000 0101    Tag 0000, Index 0, Offset 101 - H
0b0110 1000    Tag 0110, Index 1, Offset 000 - M, Compulsory
0b1100 1000    Tag 1100, Index 1, Offset 000 - M, Compulsory
0b0110 1000    Tag 0110, Index 1, Offset 000 - H
0b1101 1101    Tag 1101, Index 1, Offset 101 - R, Compulsory
0b0100 0101    Tag 0100, Index 0, Offset 101 - M, Compulsory
0b0000 0100    Tag 0000, Index 0, Offset 100 - H
0b0011 0000    Tag 0011, Index 0, Offset 000 - R, Compulsory
0b1100 1011    Tag 1100, Index 1, Offset 011 - R, Conflict
0b0100 0010    Tag 0100, Index 0, Offset 010 - R, Capacity

```

4.2 What is the hit rate of our above accesses?

$$\frac{3 \text{ hits}}{11 \text{ accesses}} \approx 27.3\% \text{ hit rate}$$

5 Code Analysis

Given the follow chunk of code, analyze the hit rate given that we have a byte-addressed computer with a total memory of **1 MiB**. It also features a **16 KiB** Direct-Mapped cache with **1 KiB** blocks. Assume that your cache begins cold.

```

#define NUM_INTS 8192    // 2^13
int A[NUM_INTS];        // A lives at 0x10000
int i, total = 0;
for (i = 0; i < NUM_INTS; i += 128) {
    A[i] = i;           // Line 1
}
for (i = 0; i < NUM_INTS; i += 128) {
    total += A[i];      // Line 2
}

```

5.1 How many bits make up a memory address on this computer?

We take $\log_2(1 \text{ MiB}) = \log_2(2^{20}) = 20$.

5.2 What is the T:I:O breakdown?

Offset = $\log_2(1 \text{ KiB}) = \log_2(2^{10}) = 10$
Index = $\log_2(\frac{16 \text{ KiB}}{1 \text{ KiB}}) = \log_2(16) = 4$
Tag = $20 - 4 - 10 = 6$

5.3 Calculate the cache hit rate for the line marked Line 1:

The integer accesses are $4 * 128 = 512$ bytes apart, which means there are 2 accesses per block. The first accesses in each block is a compulsory cache miss, but the second is a hit because $A[i]$ and $A[i+128]$ are in the same cache block. Thus, we end up with a hit rate of **50%**.

5.4 Calculate the cache hit rate for the line marked Line 2:

The size of A is $8192 * 4 = 2^{15}$ bytes. This is exactly twice the size of our cache. At the end of Line 1, we have the second half of A inside our cache, but Line 2 starts with the first half of A. Thus, we cannot reuse any of the cache data brought in from Line 1 and must start from the beginning. Thus our hit rate is the same as Line 1 since we access memory in the same exact way as Line 1. We don't have to consider cache hits for total, as the compiler will most likely store it in a register. Thus, we end up with a hit rate of **50%**.

6 Cache Performance

Recall that AMAT stands for Average Memory Access Time. The main formula for it is:

$$\text{AMAT} = \text{Hit Time} + \text{Miss Rate} * \text{Miss Penalty}$$

In a multi-level cache structure, we can separate miss rates into two types that we consider for each level.

- **Global:** Calculated as the number of accesses that missed at that level divided by the total number of accesses *to the cache system*.
- **Local:** Calculated as the number of accesses that missed at that level divided by the total number of accesses *to that cache level*.

6.1 In a 2-level cache system, after 100 total accesses to the cache system, we find that the L2\$ (L2 cache) ended up missing 20 times. What is the global miss rate of L2\$?

$$\frac{20}{100} = 20\%$$

6.2 Given the system from the previous subpart, if L1\$ had a local miss rate of 50%, what is the local miss rate of L2\$?

$\frac{20}{50\% * 100} = \frac{20}{50} = 40\%$. We know that L2\$ is accessed when L1\$ misses, so if L1\$ misses 50% of the time, that means we access L2\$ 50 times, of which we ended up having 20 misses in L2\$.

Suppose your system consists of:

1. An L1\$ that has a hit time of 2 cycles and has a local miss rate of 20%
2. An L2\$ that has a hit time of 15 cycles and has a global miss rate of 5%
3. Main memory where accesses take 100 cycles

6.3 What is the local miss rate of L2\$?

The number of accesses to the L2\$ is the number of misses in L1\$, so we divide the global miss rate of L2\$ with the miss rate of L1\$.

$$\text{L2\$ Local miss rate} = \frac{\text{Misses In L2\$}}{\text{Accesses in L2\$}} = \frac{\text{Misses in L2\$}}{\text{Total Accesses}} / \frac{\text{Misses in L1\$}}{\text{Total Accesses}} =$$

$$\frac{\text{Global Miss Rate}}{\text{L1\$ Miss Rate}} = \frac{5\%}{20\%} = 0.25 = 25\%$$

6.4 What is the AMAT of the system?

$\text{AMAT} = 2 + 20\% \times (15 + 25\% \times 100) = 10$ cycles, as the Miss Penalty of the L1\$ is the 'local' AMAT of the L2\$.

Using global rates of each level, alternatively, $AMAT = 2 + 20\% \times 15 + 5\% \times 100 = 10$ cycles (using global miss rates)

- 6.5 Suppose we want to reduce the AMAT of the system to 8 cycles or lower by adding in a L3\$. If the L3\$ has a local miss rate of 30%, what is the largest hit time that the L3\$ can have?

Let H = hit time of the cache. Extending the AMAT equation so that the Miss Penalty of the L2\$ is the 'local' AMAT of the L3\$, we can write:

$$2 + 20\% * (15 + 25\% * (H + 30\% * 100)) \leq 8$$

Solving for H, we find that $H \leq 30$. So the largest hit time is 30 cycles.