

# MIPS R2000 CACHE MEMORY

## DATA CACHE

### 1. Introduction

This lab session continues the previous one, where we focused on the code cache. We continue to use the same assembly program as in the last session, but our target is to study the data cache's operation and performance. We also continue to use the *PCSpim-Cache* simulator.

#### 1.1. Goals

This session shares the goals of the previous one but is centred on the data cache this time.

- To observe how cache memory helps reduce access time to data.
- To determine the address structure from the cache point of view
- Analyse how the cache organisation affects the hit rate.

#### 1.2. Material

The material is available in the Resources folder in PoliformaT:

- PCSpim-Cache Simulator MIPS R2000.
- Source file "program.s".

### 2. Test program: vector by constant product

```
#####  
# Data segment  
#####  
  
A:      .data 0x10000000  
        .word 0,1,2,3,4,5,6,7    # Vector A  
        .data 0x10001000  
B:      .space 32                 # Vector B (result)  
        .data 0x1000A030  
k:      .word 7                   # Scalar constant  
dim:    .word 8                   # Vector dimension
```

```
#####
# Code segment
#####

.text 0x00400000
.globl __start

__start:      la $a0, A           # $a0 = A address
              la $a1, B           # $a1 = B address
              la $a2, k           # $a1 = k address
              la $a3, dim         # $a2 = dimension address
              jal sax            # Subroutine call

#####
# Execution ending with a system call
#####

addi $v0, $zero, 10      # Exit code
syscall                # Execution end

#####
# Subroutine that computes Y <- k*X
# $a0 = Starting address of vector X
# $a1 = Starting address of vector Y
# $a2 = Address of scalar constant k
# $a3 = Address of dimension
#####

sax:          lw $a2, 0($a2)      # $a3 = constant k
              lw $a3, 0($a3)      # $a3 = dimension
loop:         lw $t0, 0($a0)      # Reading X[i] into $t0
              mult $a2, $t0       # Computes k*X[i]
              mflo $t0           # $t0 <- k*X[i] (HI value is 0)
              sw $t0, 0($a1)      # Writing Y[i]
              addi $a0, $a0, 4    # Address of X[i+1]
              addi $a1, $a1, 4    # Address of Y[i+1]
              addi $a3, $a3, -1   # Decrements the number of elements
              bgtz $a3, loop      # Jumps if elements remain
              jr $ra             # Subroutine return

.end
```

### 3. Data cache

We now consider a data cache with the features given in Table 1. We call this the *original* configuration because we will be changing it during this session.

Parameter	Value
Capacity	256 bytes
Mapping	Direct
Block or line size	16 bytes
Write policy	Write-through, write allocate

Table 1. *Original* data cache parameters.

In this section, we pay attention to data accesses from the program due to *load* and *store* instructions, that is, *lw*, *sw* and their type variations *lb*, *lbu*, *sh*, etc. The test program has two load instructions for constants *k* and *dim*, then *dim* reads of the source vector A and the same number of writes to the result vector B.

Take the following questions **before you use the simulator**.

1. Considering the given cache parameters, indicate how many lines are in the cache and how many lines are needed to store each vector.

$$256 / 16 = 16 \text{ cache lines}$$

$$8 \text{ words in a vector} * 4B = 32B \rightarrow 32B / 16B = 2 \text{ lines / vector}$$

2. Give the address structure from the point of view of this data cache (bits for the tag, line, and offset).

$$\log_2(16B \text{ block}) = 4 \text{ offset}$$

$$\log_2(16 \text{ lines}) = 4 \text{ line num}$$

$$32 - 4 - 4 = 24 \text{ tag}$$

If you apply this structure to the addresses of variables *dim* and *k*, you should verify that both of them map to line 3 in the data cache.

3. Give the cache lines where vectors *A* and *B* are mapped.

0x1000 0000 --> Last digit is offset, second last digit is starting line --> Vector A is in lines 0 and 1

0x1000 1000 --> Last digit is offset, second last digit is starting line --> Vector B is in lines 0 and 1

0x1000 A030 --> Last digit is offset, second last digit is starting line --> *k* and *dim* are in line 3

Given the mapping scheme, we will face conflict misses if both vectors are mapped to conflicting cache lines. We will check this **with the simulator**.

4. Load the original program in PCSpim-Cache. Configure the data cache with the parameters given in Table 1. Run the program step by step (F10) and observe the effects on the data cache. Verify that vectors are stored in the intended data cache lines and then fill the following table:

Nr. of data accesses	18
Nr. of hits	1
Nr. of misses	17
Hit rate achieved (H)	0.056

5. Which data access(es) have caused cache hits?

The only hit occurs when accessing *k* and *dim*, as they are in the same block. When using `lw $a2, 0($a2)` to access *k*, there's a miss, and the entire block is brought to cache. When using `lw $a3, 0($a3)` to access *dim*, the block is already in cache, so there's a hit.

6. Explain whether the low hit rate is due to the absence of locality in the program or due to one or more of the cache parameters used.

It's because we are using direct mapping. Vectors *A* and *B* are mapped to the same line, so there are conflict misses. If we were using 2-way s set-associative cache, we could have both vectors loaded in cache, greatly reducing conflict misses.

In the following experiments, we will explore the effect of three different approaches to improve this low rate:

- Change the write-miss policy to *no-allocate*.
- Modify the vector addresses in the data segment using `.data` directives.
- Increase associativity using 2-way set associative mapping.

**NOTE:**

Each proposed change must be independently applied to the original configuration in Table 1. Be careful not to combine them when you move to the next exercise.

### 3.1. First alternative: change the write-miss policy

In the original configuration, a write miss to vector B implies the allocation of the missed block. Since vector B is only accessed for writing, we could avoid collision with vector A by not allocating B, thus changing the write miss policy to write no-allocate. Let's see if that works. See Table 2 for the cache configuration to apply (change highlighted *in red*).

Parameter	Value
Capacity	256 bytes
Mapping	Direct
Block or line size	16 bytes
Write policy	Write-through, <b>no-write-allocate</b>

Table 2. Data cache parameters for the first alternative.

7. Configure the data cache in PCSpim-Cache as indicated in Table 2. Run the original program step by step (F10) to observe its behaviour. Then fill in the table:

Nr. of data accesses	18
Nr. of hits	7
Nr. of misses	11
Hit rate achieved (H)	0.389

Double-check that you have counted one miss for every access to vector B, two misses accessing vector A (read elements A [0] and A [4]) and a miss reading variable *k*. Note that we have avoided collision between A and B since B is only accessed for writing, and there is no allocation in case of a write miss. However, all accesses to B are misses. Let's try another approach.

### 3.2. Second alternative: change the vectors' allocation

A different way to avoid collisions is to modify the allocation of vectors in the data segment. We can do that using the `.data` directive. In the original program and with the original cache configuration, the addresses where A and B are allocated are causing collisions. To avoid collisions, we may change the starting address of either A or B so that they are mapped to non-conflicting lines in the cache.

**8.** Change the directive `.data 0x10001000` so that vector B maps to lines 4 and 5 of the cache. This would avoid collision with lines 0 and 1 (vector A) and 3 (k and dim). Write the new address for vector B here.

0x1000 1040

**9.** Assemble and run the program with this change. Remember to keep the original write policy (write allocate) and fill the following table:

Nr. of data accesses	18
Nr. of hits	13
Nr. of misses	5
Hit rate achieved (H)	0.72

Compare the hit rate achieved after taking this second alternative.

### 3.3. Third alternative: increase the cache associativity

The third alternative is to use a more flexible mapping scheme. Restore the original location of vector B (`.data 0x10001000` before its declaration) and configure the data cache as shown in Table 3.

Parameter	Value
Capacity	256 bytes
Mapping	2-way set associative
Block or line size	16 bytes
Write policy	Write-through, write-allocate
Replacement policy	LRU

Table 3. Data cache configuration for the third alternative

Using 2-way set associative mapping, a block from main memory may now be brought to any of the two lines of the *set* it maps to.

**Before you use the simulator again**, take the following questions

*10. Give the address structure from the point of view of this data cache (tag, line, and offset bits).*

24 + 1 = 25 bits for tag

4 - 1 = 3 bits for set (sets of 2 lines)

4 bits for offset

*11. Give the cache sets where vectors A and B are mapped.*

After the first reading, both vectors will use the set 0 (one will take one of the ways of the set, and the other the other way). When elements 4, 5, 6, and 7 are accessed, the set 1 will be used. One of the vectors will take one of the ways of the set, and the other the other way.

Note that the blocks of vectors *A* and *B* map to coincident cache sets. Direct mapping caused continuous replacements at run time. Are we in the same situation now? If not, what has changed?

No, now both vectors can be in cache at the same time, so continuous replacements don't happen at runtime.

It's time to **use the simulator** again.

*12. Load the original program and run it with the new cache configuration (Table 3). Then fill in the following table:*

Nr. of data accesses	18
Nr. of hits	13
Nr. of misses	5
Hit rate achieved (H)	0.72

Note that there are no replacements in this case. All misses are compulsory, *i.e.*, caused by the first access to the block. You should only find misses of this kind.

Compare alternatives 2 and 3. Which one do you think is preferable in general?