

# Template Matching with Cache Friendly Code

## COMP273 Assignment 4 - Fall 2018, Prof. Kry

School of Computer Science, McGill University

Available: 11:30 PM, 11 November 2018

Due date: 11:30 PM, 27 November 2018

Submit electronically via MyCourses



<http://pxlcon.jimmysomething.com>

## Template Matching

In this assignment, you will use template matching to find Waldo in a pixel art image. Waldo is easily recognizable because of his glasses and his red and white striped shirt and hat. Template matching is the simplest of a family of algorithms that are used for optical character recognition in scanned documents, and automatic face detection in images. It will be indeed face detection you will be doing in this assignment, although with a very small template of only 8-by-8 pixels. Using such a small template will allow your assignment to run at a reasonable speed in the simulator, in contrast to the larger template shown at right. This small fixed-size template matching process is also very closely related to the motion estimation step (optical flow computation) that video compression algorithms perform on pairs of successive frames.



There are many ways to compute the error of a match, but a common choice is the sum of absolute differences of pixel intensities,

$$e(x, y) = \sum_{i=0}^{w-1} \sum_{j=0}^{h-1} |I(x+i, y+j) - T(i, j)|. \quad (1)$$

The function  $e$  is the error for a given pixel location  $(x, y)$ , and  $I$  and  $T$  are functions that provide the pixel intensities in the image and template, respectively. The double sum loops over the width  $w$  and height  $h$  of the template (in our case,  $w = h = 8$ ). The better the match, the lower the error. If template and image intensities match exactly the error will be zero.

While we could consider doing a comparison of colours, it is simpler to only work with intensities. Thus the image and template provided in this assignment will be gray-scale images. That is, when you load the bitmap image, the red, green, and blue components of the word corresponding to a pixel will all be the same value, and we **only need to read one byte** (1bu) to know the pixel's brightness (a number between 0 and 255). This wastes memory, using four times what is required, but is convenient because it will allow use of the bitmap display for visualization and debugging.

## Provided Data and Code

Useful functions are provided for you in the `templatematch.asm` file. You will implement two different versions of the template matching function by modifying and submitting this file. Be sure to enter your name and student number at the top of the file.

In the provided file, the data section reserves space for the display buffer, an error buffer, and a template buffer. It also sets up some simple data structures to store information relating to each buffer, specifically, the address in memory where it is found, the width, the height, and the filename if there is one. It is this buffer information structure that is used as an argument to the provided functions, and to the `matchTemplate` and `matchTemplateFast` functions that you will write. The image and error buffer regions are placed at the beginning of the static data `0x10010000` such that you can easily visualize both in the memory mapped bitmap display.

```
.data
displayBuffer: .space 0x40000 # space for 512x256 bitmap display
errorBuffer:   .space 0x40000 # space to store match function
templateBuffer: .space 0x100   # space for 8x8 template
imageFileName: .asciiz "pxlcon512x256cropgs.raw"
templateFileName: .asciiz "template8x8gs.raw"
# struct bufferInfo { int *buffer, int width, int height, char* filename }
imageBufferInfo: .word displayBuffer 512 128 imageFileName
errorBufferInfo: .word errorBuffer 512 128 0
templateBufferInfo: .word templateBuffer 8 8 templateFileName
```

The following functions are provided.

```
void loadImage( bufferInfo* imageBufferInfo )
```

Loads an image or template from file.

```
(offset, score) = findBest( bufferInfo errorBufferInfo )
```

Finds the best match in the error buffer and returns the offset in bytes along with the score.

```
void highlight( bufferInfo imageBufferInfo, int offset )
```

Highlights in green the area corresponding to the best match in the image.

```
void processError( bufferInfo errorBufferInfo )
```

Processes the error into a score between 0 and 255 for viewing in the bitmap display. The error will generally be much larger than 255, and when viewed in the memory mapped display, the unprocessed error will be broken up into fields for the red, green, and blue intensities in a manner which is difficult to interpret. This function will allow you to better visualize the quality of the match across the entire image. A perfect match will show up as a bright green dot in the memory mapped display.

The provided functions make the assumption that **the template is always 8 by 8, and you should make the same assumption in your code too!** However, **do not assume that the image size is fixed at 512 by 128!** Having a variable image size allows it to be reduced when measuring performance in the last part of the assignment. You might notice that the `pxlcon512x256cropgs.raw` is actually 512 by 256, but we are only loading the first half of the file.

The image and template data files are stored in a `raw` binary format. This makes for a very simple `loadImage` function, which directly loads the contents of the file into memory. These files **must be placed in the directory where you launch the MARS**, otherwise they will not be found.

## Naive Implementation (8 marks)

Complete the function `matchTemplate` in the `templatematch.asm` file. Assuming the sum of absolute differences array is initialized to zero (which is how the memory is initialized in MARS), the following pseudo-code will naively compute the sum of absolute differences, i.e., the error.

```
for ( int y = 0; y <= height - 8; y++ ) {    // loop through the image
    for ( int x = 0; x <= width - 8; x++ ) {
        for ( int j = 0; j < 8; j++ ) {        // loop through the template image
            for ( int i = 0; i < 8; i++ ) {
                SAD[x,y] += abs( I[x+i][y+j] - T[i][j] );
            }
        }
    }
}
```

Note that the template is assumed to be 8 by 8, and that there will be a portion of the error buffer which is untouched because there is not a complete overlap of the template and image to test. This corresponds to the **last 7 words of each line, as well as 7 lines at the bottom of the error buffer** (notice the *less than equal* test in the for loops, instead of *less than*).

Despite the fact that you will need to make a faster implementation of the template match in the next objective, you must still complete this part of the assignment so that you have both versions of the function to compare in the last objective of the assignment.

## Faster Cache Friendly Implementation (8 marks)

The naive implementation will not work well with a small memory cache because of the inner loops which loads an entire 8 by 8 area of both the image and the template. The memory associated with an 8 by 8 area is 64 words, or 128 bytes, which means that 256 bytes of memory will be read for both template and image. The default settings of the *data cache simulation tool* is 128 bytes. **With this default cache, suppose we complete the two inner loops for the pixel at (0,0) and we advance  $x$  to the next pixel at (1,0). Our first loads to access the image and template will result in cache misses.** In fact, while we will still want most of the blocks that are currently in the cache, they will all be flushed out of the cache as the least recently used blocks before the inner loop code tries to load from them. Increasing the size of the cache can help, but there is also a very simple code modification that will help performance, even when the cache is small. Consider first (we will do better shortly) what happens when the order of the loops is changed.

```
for ( int j = 0; j < 8; j++ )
    for ( int y = 0; y <= height - 8; y++ ) {
        for ( int x = 0; x <= width - 8; x++ ) {
            for ( int i = 0; i < 8; i++ ) {
                SAD[x,y] += abs( I[x+i][y+j] - T[i][j] );
            }
        }
    }
```

Here, the inner loop processes only one row of template pixels. Thus the two inner loops only use a very small amount of memory, specifically 8 words of the image, and 8 words of the template, or 64 bytes. In the default 128 byte cache, all of these blocks *can still be* in the cache! The wording here is weak (*can still be* as opposed to *will still be*) because there could be conflict misses depending on the cache organization. In this assignment, we will mostly use the default settings of the cache. That is, we will always assume the cache holds a total of 8 blocks, with blocks being 4 words each, for a total of 128 bytes, and we will adjust the settings to consider direct, 2 way set associative, and fully associative caches.

As an improvement on the loop reordering, consider also unrolling the inner most loop. It only runs 8 times, and the 8 loads of the template can be stored in registers. This speeds up your code by avoiding the comparison and branching for the inner loop, but more importantly reduces the number of memory accesses needed in the inner loop, which can also help reduce conflict misses!

```
for ( int j = 0; j < 8; j++ )
    int t0 = T[0][j];
    int t1 = T[1][j];
    int t2 = T[2][j];
    int t3 = T[3][j];
    int t4 = T[4][j];
    int t5 = T[5][j];
    int t6 = T[6][j];
    int t7 = T[7][j];
    for ( int y = 0; y <= height - 8; y++ ) {
        for ( int x = 0; x <= width - 8; x++ ) {
            SAD[x,y] += abs( I[x+0][y+j] - t0 );
            SAD[x,y] += abs( I[x+1][y+j] - t1 );
            SAD[x,y] += abs( I[x+2][y+j] - t2 );
            SAD[x,y] += abs( I[x+3][y+j] - t3 );
            SAD[x,y] += abs( I[x+4][y+j] - t4 );
            SAD[x,y] += abs( I[x+5][y+j] - t5 );
            SAD[x,y] += abs( I[x+6][y+j] - t6 );
            SAD[x,y] += abs( I[x+7][y+j] - t7 );
        }
    }
}
```

Implement this cache friendly version of the template matching loop in the `matchTemplateFast` function of the `templatematch.asm` file. This can be the algorithm described here, but you are also free to explore other ideas in implementing the algorithm in any way that improves cache performance. Keeping your approach simple is good advice, and please provide clear comments to document your code. Finally, be sure to **read the rest of this assignment before starting this objective!**

## Measuring Cache Performance and Questions (4 marks)

Once you have implemented both regular and fast versions of the template matching function, you will test your implementation and document its behaviour by tabulating performance numbers in the file `stats.csv` (see details below). Because the *data cache simulator* and the *instruction counter* tools slow down execution of your program, you will reduce the image and error buffer sizes to only 16 lines for measurements.

```
imageBufferInfo:    .word displayBuffer  512 16   imageFileName
errorBufferInfo:    .word errorBuffer    512 16   0
```

Replace these lines **only during measurements!** Be sure to submit your code with a value of 128 in the height of the `imageBufferInfo` and the `errorBufferInfo`!

**Measure only the instructions for `matchTemplate` to run.** Set one breakpoint at `jal matchTemplate` and another at the next line. Run your code up to the breakpoint, press the *connect to MIPS* button on the instruction counter, press the run button to continue execution, then make note of the *instructions so far* when you reach the next breakpoint. Measure the instruction count for both your naive and cache friendly implementations.

Use a similar process for measuring the cache performance of `matchTemplateFast`. Record the memory accesses and cache misses for your naive and cache friendly implementations. Test with fully asso-

ciative, a 2-way set associative, and direct mapped caches. Be sure to reset the cache simulator between measurements!

To make your cache friendly version work well with a direct mapped cache, you will want to consider if there are conflict misses that can be avoided.

1. Do the base addresses of image and error buffers fall into the same block of the direct mapped cache?
2. For the templateMatchFast, does it matter if the template buffer base address falls into the same block as the image or template buffer base address in the direct mapped cache?

Consider altering the memory layout in your data segment by adding `.space` directives as necessary to make your fast implementation work well with the direct mapped cache. In an `asm` file comment, **provide a brief explanation answering the questions above** to justify your choice of spacing.

Your `csv` file **must exactly match the required filename and format**. You may include comments in the file by starting a line with `#`, but otherwise there are six lines in this file to complete with six *comma separated values*, or fields, on each line. These fields consist of your student number, the test name (Naive or Fast implementation, with fully-associative (FA), 2-way set associative (2way), or Direct mapped cache), the instruction count, the number of memory accesses, the number of cache misses, and an execution time in microseconds. Here follows an example, but please note that you should not consider the values shown below to be those of a correct solution!

```
# StudentID, Case, InstCnt, MemAccess, Misses, MicroSeconds
260123456, NaiveFA,      2308937, 627243, 177209, 20029
260123456, Naive2Way,   2308937, 627243, 190466, 19277
260123456, NaiveDirect, 2308937, 627243, 246806, 26989
260123456, FastFA,      1963978, 363685, 18385, 3802
260123456, Fast2Way,    1963978, 363685, 18385, 3802
260123456, FastDirect,  1963978, 363685, 18385, 3802
```

You must compute the time in microseconds assuming a processor that runs at 1 GHz and executes one instruction every cycle, and assuming the cache miss penalty to be 100 cycles. In the example above, the fast implementation uses fewer instructions and would be 88  $\mu$ s faster if we were only counting instructions, but there are also far fewer cache misses, which can make the fast implementation almost 6.3 times faster than the naive implementation!

If you are unable to complete either of the first two objectives, you will not be able to receive full marks on this objective. To assist in marking, please either omit the file from your submission, or omit the appropriate lines in the `stats.csv` file, or include a comment to signal that either one or both of your functions is incomplete or has bugs.

## Bonus (4 marks)

Bonus marks will be awarded to students with the best direct mapped cache performance using the 100 cycle cache miss penalty described above. These marks will be given to the top 10% of solutions that **respect register conventions** and are *submitted on time*.

## Submission Instructions

All work must be your own, and must be submitted by MyCourses. **Include your name and student number in your source file**. You must **submit exactly two files**: `templatematch.asm` and `stats.csv`. Use comments to explain any optimizations you make beyond those suggested in the assignment. **Do not use an archive**, and do not submit any of the other provided data files. Be sure to **check your submission**

by downloading your submission from the server and checking that it was correctly submitted. You will not receive marks for work that is incorrectly submitted. *Note that you are encouraged to discuss assignments with your classmates, but not to the point of sharing code and answers. All code and written answers must be your own.*