

Parallel Computing: Image noise filtering project

Liam Pulles

September 30, 2016

1 INTRODUCTION

1.1 PROBLEM STATEMENT

1.2 PROPOSED SOLUTION OVERVIEW

1.3 BACKGROUND THEORY

2 SOLUTION TECHNICALITIES

2.1 SERIAL IMPLEMENTATION

The serial implementation works as follows:

1. An image is loaded into memory (a pointer to BMP).
2. The matrix structure is created (we've made a few standard 5x5 and 3x3 matrices).
3. The convolution function is called.
4. The convolution function makes a blank output BMP* structure with the same dimensions as the source image for storing the calculated pixels.
5. The convolution function calls the kernel function for each pixel, using two nested for loops.
6. The kernel function calculates the kernel for a pixel using the edge extrapolation function, the matrix structure, and the source image. The final calculation for the pixel is stored in the output BMP* structure directly for speed.

7. Once the whole convolution is complete, the program writes a new file with the filtered image.

2.2 PARALLEL IMPLEMENTATION

The parallel implementation differs from the serial implementation in that instead of the master thread (the serial program) iterating through each pixel, the parallel implementation has several threads going through different pixels at the same time.

This is achieved by use of an omp for loop structure. We use dynamic scheduling to accomodate for potential differences in complexity in regions of the image. Each iteration of the outer for loop calculates the kernel for a $(\text{image height}) \times (\text{matrix width} \times 2)$ sized chunk of the image (See figure ??) (Note that the final chunk is truncated to fit with the width of the image). There are a few reasons for this choice of task decomposition:

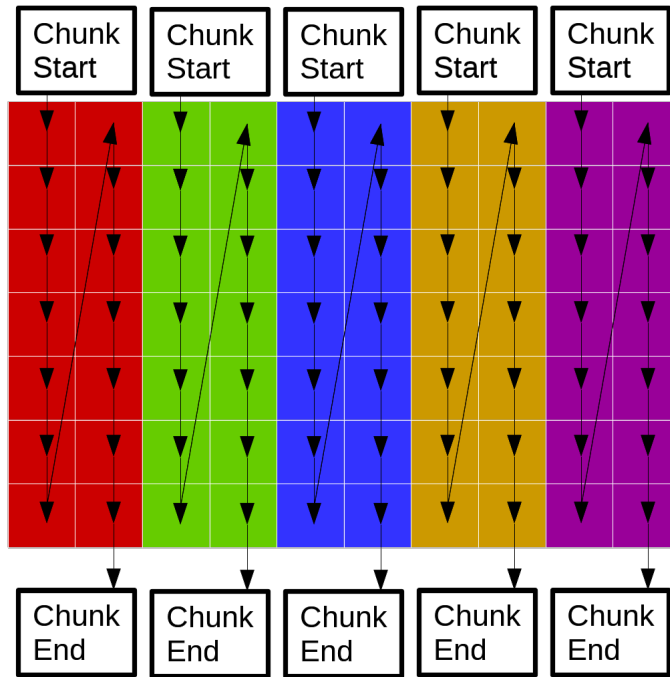


Figure 2.1: The decomposition of the workload of the image into chunks.

- We want to keep the number of chunks for threads to process low so as to reduce unproductive inter-thread communication time, but we also want there to be enough chunks so that threads that execute their chunks quickly can still have something useful to do while waiting for other threads.
- We want the dimensions of the chunks to minimize potential contention between threads, specifically when we are calculating the kernel for two separate pixels which are in kernel dimension range (this would result in both the kernel calculations

needing to contend for the same pixels for their own calculation, see figure ??). Assuming the image is reasonably sized, using $(\text{matrix width}) * 2$ width chunks will nearly always keep the kernel calculations far enough away from each other at each moment of execution, even when one thread is running faster than the others. For a 5x5 kernel, a thread would have to have done more than 5 columns of the image before one thread had done a single column for there to be contention issues in this regard.

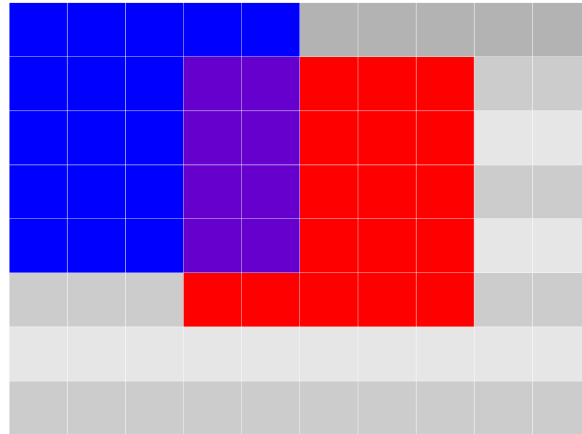


Figure 2.2: An example of contention with nearby simultaneous kernel calculations. The purple region represents potentially contested pixels.

Each thread works column by column. We do not create separate sub images for each thread to use for source and output because it would require splitting overhead at the beginning and joining overhead at the end. It is mentioned in the EasyBMP documentation (EasyBMP Manual 1.06, pg. 16) that the RGBApixel pointers are stored internally as a 2d array. We also know that this must be stored on the heap, since it has to do deal with varying sized images (and thus, varying sized 2d arrays). Finally, we can deduce that since the size of the array is width x column that these 2d arrays store arrays of rows (the x component of the image).

So what about the contention that arises when threads read from the same source image in memory at approximately the same time, and write to the same output image at approximately the same time. Will contention on the cache line arise from reading/writing the same row?

We compared two version of the parallel implementation - one which sequentially accessed by row and one which sequentially accessed by column. The difference in execution time was negligible.

Why is this? We believe it is because the 2d array stores pixel POINTERS.

Suppose there is contention between 2 threads trying to access far away pixels on the same row, and - we shall presume for argument's sake - the same cache line. Then one thread will have to wait while the other gets the pointer data (which is of size int). But

after that contention, modifying or reading to that RGBApixel - which forms the vast bulk of the process of modifying a pixel - will be independent of modifying or reading any other RGBApixel, since it is a pointer and can exist anywhere in the heap and thus on separate cache lines.

So since the delay caused by the contention is so small, it has an almost zero sum effect on the total execution time. The number of contentions also becomes less over time, as once a contention does arise, one of the threads will move forward while another has to wait. So ultimately, the threads will often be working on both different rows and different columns anyway.

3 RESULTS ANALYSIS

4 CONCLUSION