

# The assessment of parallelizing matrix multiplication using openCL

Cameron F Clark<sup>†</sup> and Kian S Frassek<sup>‡</sup>

EEE4120F Class of 2024

University of Cape Town

South Africa

<sup>†</sup>CLRCAM007 <sup>‡</sup>KIAFRS001

## I. INTRODUCTION

Parallel computing is a power tool which enables code to run at phenomenal rates compared to single thread. There are many different ways to implement parallel programs on many different problems. The implementation that this report focuses on is using OpenCL and C++ to write a parallel program to efficiently multiply square matrices. Importantly, it is futile to design parallel programs that do not actually increase the efficiency and/or speed of the non parallelized programs, thus, this report will additionally evaluate the performance of the parallel program versus the single thread program with different matrix sizes and number of matrix multiplications.

## II. METHODOLOGY

### A. Hardware

The hardware used to conduct these experiments is a MacBook Pro 13" 2019 model. All of these parallel programs were conducted on the GPU which is a Intel Iris Plus Graphics 1536 MB. All of the single thread programs were conducted on the CPU which is 2 GHz Quad-Core Intel Core i5.

### B. Implementation

1) *Single Thread Program:* The single thread program is a relative simple implementation in C code. It requires two for loops which iterate over the row and column of the result matrix and a third which calculates the dot product of the input matrices to find the value of the cell of the output matrix. There is one additional for loop which allows for multiplying a given number of matrices together

```
int matrix_output[matrix_size], matrixA[matrix_size], cell;
for (int c = 0; c < matrix_size; c++) matrixA[c] = matrices[0][c];

//Matrix multiplication
for (int m = 1; m < MATRIX_COUNT; m++) {
    for (int row = 0; row < Size; row++) {
        for (int col = 0; col < Size; col++) {
            cell = 0;
            for (int dot = 0; dot < Size; dot++) {
                cell += matrixA[row * Size + dot] * matrices[m][col + dot * Size];
            }
            matrix_output[row * Size + col] = cell;
        }
    }
}
//rewrite matrix A with output so loop can re-iterate
for (int c = 0; c < matrix_size; c++) matrixA[c] = matrix_output[c];
}
```

2) *Parallel Program:* The parallel program involves setting up OpenCL to run instances of a kernel program. The kernel program itself runs a very similar instance of the cell calculation for the single thread program above where the

row is equivalent to the workGroupNumber and the column is equivalent to the localGroupID. A for loop is used to calculate the dot product as in the single thread program. Lastly, there is an additional for loop in the OpenCL setup code which allows for successive matrix multiplication.

```
__kernel void matrixMultiplication(
    __global int* matrixA_buffer,
    __global int* matrixB_buffer,
    __global int* output_buffer
){
    int workItemNum = get_global_id(0); //Work item ID
    int workGroupNum = get_group_id(0); //Work group ID
    int localGroupID = get_local_id(0); //Work items ID within each work group
    int localSize = get_local_size(0); //Get the work number of items per group

    //Row <=> workGroupNum; Column <=> localGroupID
    //workItemNum <=> workGroup * localSize + localGroupID
    int cell = 0;
    for (int i = 0; i < localSize; i++) {
        cell = cell + *(matrixA_buffer + workGroupNum * localSize + i) *
            *(matrixB_buffer + localGroupID + i * localSize);
    }
    *(output_buffer + workItemNum) = cell;
}
```

### C. Experiment Procedure

To evaluate the performance of the single thread versus parallel program two experiments will be conducted and evaluate the difference in time taken and the speed up of the parallel program. The first will be multiplying two matrices with different sizes, and the second will be to multiply different numbers of successive matrices of the same size. Both methods of testing have different amounts of parallelization and different amounts of overhead and should give good understanding of the effectiveness of parallelizing matrix multiplication code.

## III. RESULTS

### A. Multiplying two matrices with different sizes

This section discusses the time taken both for parallelized and single thread multiplying two matrices of the same size. The test repeated the multiplication with the size of the matrix increasing each time. The results of the multiplication are shown in figure ??.

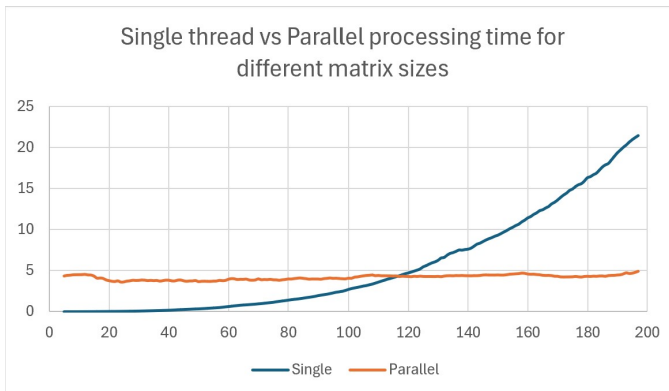


Fig. 1: Graph of single thread vs Parallel processing time for different matrix sizes

The results show that the time taken for single-thread matrix multiplication is initially faster than parallelized multiplication up to size of the matrix around 117. After size 117, parallelized multiplication becomes faster. This is indicated on the graph by time taken for the single thread processing increases with a larger exponential scale. These results occur because the time due to overhead of setting up parallel processing is 4 ms which is larger than the time taken when running the program with small sized matrices on a single thread. However, as the matrix size increases, the time taken for single-thread multiplication increases at a faster exponentially rate, resulting in longer processing time.

#### B. Multiplying different number of matrices with same sizes

This section discusses the time taken for parallelized verses single thread multiplying different numbers of matrices of the same size. Looking at the results in figure ??, ??, ??, ?? and ??, It can be seen that there are two different types of overheads the openCl setup overheads and the overhead of starting the kernel to run. The openCl setup overheads are the same for all the different matrix counts, this is because the openCl setup is done once and the same for all the different matrix counts. This is indicated in our graph because the time taken at matrix count of 5 remain around 5ms independent of matrix size.

The overhead of starting the kernel to run seem to result in a linear increase in time as the matrix count increases. This occurs because each time the a new test is run a the buffer need to be re initialized and the kernel need to be started. This can be seen in the graph because both the parallelized and single thread result in linear increase in time as the matrix count increases. however percentage of time spent parallelizing decreases as the matrix count increases. This is indicated in the in figure ?? where the linear graf of the time take vs matrix count for parallelized multiplication and single thread multiplication intersect. in other words there is an increase in efficiency of parallelized multiplication as the matrix count increases.

Looking at all the graph it can be seen that as the matrix size increases the speed up for parallelized multiplication vs the single thread multiplication improves. in other words

at smaller matrix sizes the time taken for parallelized multiplication is longer than the time taken for single thread multiplication and at larger matrix sizes the time taken for parallelized multiplication is shorter than the time taken for single thread multiplication.

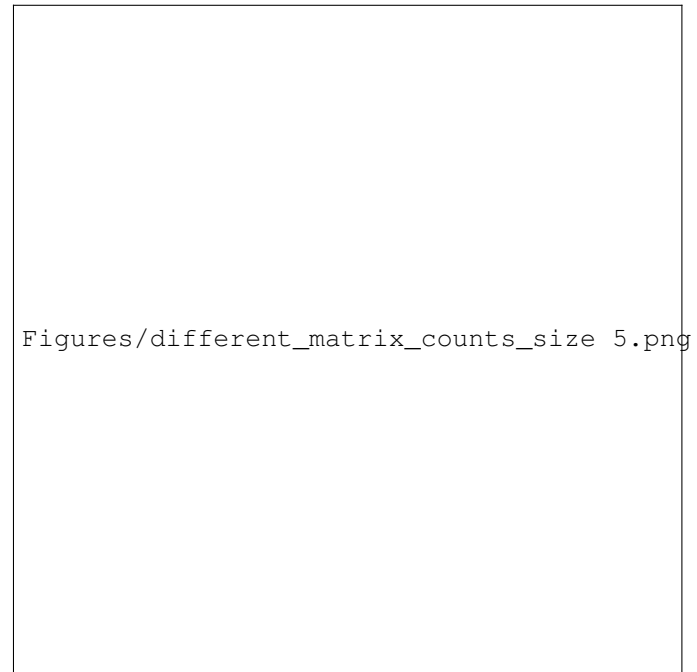


Fig. 2: Single thread vs Parallel processing time for different matrix counts all with size 10

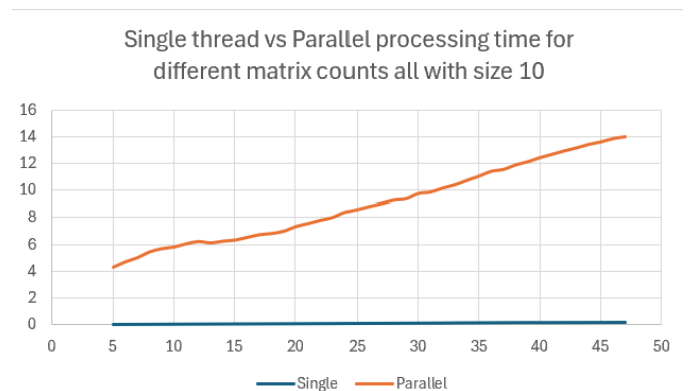


Fig. 3: Single thread vs Parallel processing time for different matrix counts all with size 10

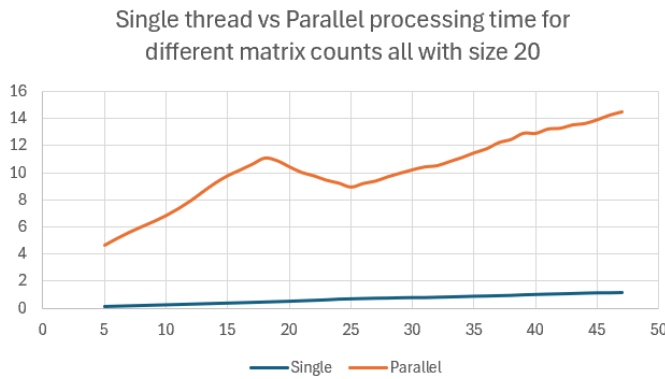


Fig. 4: Single thread vs Parallel processing time for different matrix counts all with size 20

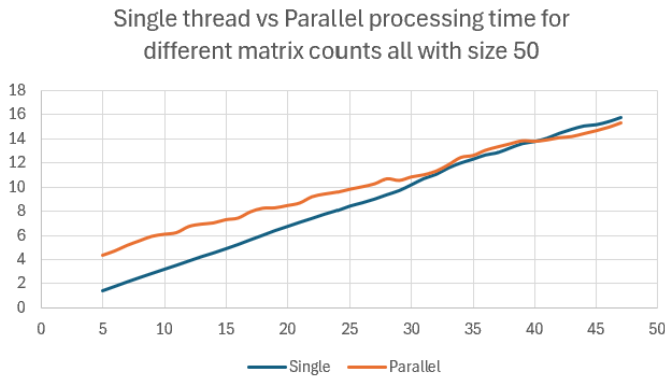


Fig. 5: Single thread vs Parallel processing time for different matrix counts all with size 50

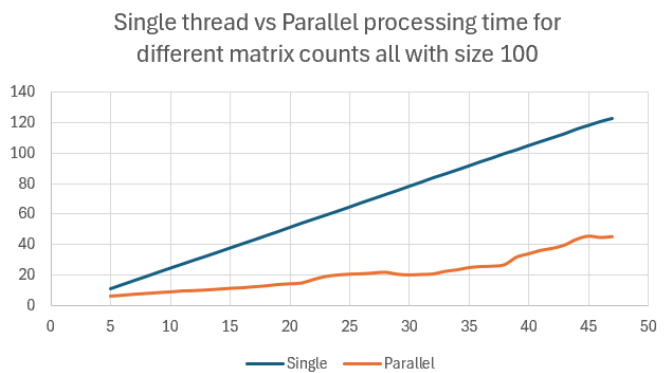


Fig. 6: Single thread vs Parallel processing time for different matrix counts all with size 100

Always remember to include axes text, units and a meaningful caption in your graphs. When typing units, a  $\mu$  sign has a tail! The letter “u” is not a valid unit prefix. When typing resistor values, use the  $\Omega$  symbol.

```
function FormatFig(X, Y, File);
set(gcf, 'PaperUnits', 'inches');
set(gcf, 'PaperOrientation', 'landscape');
set(gcf, 'PaperSize', [8, 4]);
set(gcf, 'PaperPosition', [0, 0, 8, 4]);

set(gca, 'FontName', 'Times New Roman');
set(gca, 'Position', [0.1 0.2 0.85 0.75]);

xlabel(['\n" X]);
ylabel(['Y "\n\n]);

setenv("GSC", "GSC"); # Eliminates stupid warning
print(...
[File '.pdf'],...
'-dpdf'...
);
end
```

Listing 1. Octave function to format a figure and save it to a high quality PDF graph

```
figure; # Create a new figure
# Some code to calculate the various variables to plot...
plot(N, r, 'k', 'linewidth', 4); grid on; # Plot the data
xlim([0 360]); # Limit the x range
ylim([-1 1]); # Limit the y range
set(gca, 'xtick', [0 90 180 270 360]); # Set the x labels

FormatFig(... # Call the function with:
'Phase shift \circ',... # The x title
'Correlation coefficient',... # The y title
['r_vs_N_f=' num2str(f) ';_P=' num2str(P)]... # Format the file name
);
close all; # Close all open figures
```

Listing 2. Example of how to use the FormatFig function

### C. Tables

Tables are often a convenient means by which to specify lists of parameters. An example table is presented in table ?? . You can use Tablesgenerator to make your  $\text{\LaTeX}$  tables.

### D. Pictures and Screen-shots

When you include screen-shots, pdf $\text{\LaTeX}$  supports JPG and PNG file formats. PNG is preferred for screen-shots, as it is a loss-less format. JPG is preferred for photos, as it results in a smaller file size. It's generally a good idea to resize photos (not screen-shots) to be no more that 300 dpi, in order to reduce file size. For 2-column article format papers, this translates to a maximum width of 1024. **Never change the aspect ratio of screen-shots and pictures!**

The source used to import a picture in an exact spot, with a caption and labels



Fig. 7: An example image

TABLE I  
MY INFORMATIVE TABLE

| Heading 1 | Heading 2 | Heading 3 |
|-----------|-----------|-----------|
| Data      | 123       | 321       |
| Data      | 456       | 654       |
| Data      | 789       | 987       |

#### E. Maths

L<sup>A</sup>T<sub>E</sub>X has a very sophisticated maths rendering engine, as illustrated by equation ?? . When talking about approximate answers, never use  $\pm 54$  V, as this implies “positive or negative 54 V”. Use  $\approx 54$  V or  $\sim 54$  V instead.

$$y = \int_0^{\infty} e^{x^2} dx \quad (1)$$

#### IV. CONCLUSION

The conclusion should provide a summary of your findings. Many people only read the introduction and conclusion of a paper. They sometimes scan the tables and figures. If the conclusion hints at interesting findings, only then will they bother to read the whole paper.

You can also include work that you intend to do in future, such as ideas for further improvements, or to make the solution more accessible to the general user-base, etc.

Publishers often charge “overlength article charges” [?], so keep within the page limit. In IEEE4084F we will simulate overlength fees by means of a mark reduction at 10% per page. Late submissions will be charged at 10% per day, or part thereof.

#### REFERENCES

- [1] “Voluntary Page and Overlength Article Charges,” <http://www.ieee.org/advertisement/2012vpcopec.pdf>, 2014.