

# Optimized and Naive Non-Square Matrix Multiplication With Cuda and OpenCL

Drew Afromsky  
Email: daa2162@columbia.edu

**Abstract**—In this assignment, matrix multiplication in both CUDA and OpenCL was implemented, taking advantage of both shared memory and global memory. A matrix of size  $[M \times N]$  and a second matrix of size  $[N \times M]$  are initialized. The matrix is then iteratively increased in both the x- and y-dimensions (i.e.  $[itr \times M \times N]$ ;  $[itr \times N \times M]$ ). First, a naive implementation of matrix multiplication, in both PyCUDA and PyOpenCL, is used, where both matrices are stored in global memory. Second, an optimized implementation of matrix multiplication is used, in both PyCUDA and PyOpenCL, where both matrices are stored in shared memory. Matrices in the naive and optimized PyCuda code were initialized to size  $16 \times 14$  and  $14 \times 16$ , with each dimension in each matrix was iteratively increased by a factor of 1 to 38. The dimension of the largest matrix was  $608 \times 608$ . Matrices in the naive and optimized PyOpenCL code were initialized to the same size as those in PyCuda for consistency. However, matrices in PyOpenCL naive and optimized implementations were able to compute matrices of size  $8000 \times 8000$  without encountering memory errors. Matrices in the optimized PyCuda implementation could also compute matrices of this size. Finally, run times were compared for PyCuda, Python (Serial/CPU), and PyOpenCL for both naive and optimized matrix multiplication implementations. PyOpenCL optimized code was clearly able to compute the matrix multiplication much faster than PyOpenCL naive and Python (serial/CPU). PyCuda optimized and naive had similar run-times, with the optimized code barely outperforming the naive implementation.

**Index Terms**—PyCuda, CUDA, PyOpenCL, OpenCL, naive, optimized, matrix multiplication

## I. INTRODUCTION

Matrix-matrix multiplication between two matrices of size  $[M \times N]$  and  $[N \times M]$  results in a new matrix of size  $[M \times M]$ . This operation presents opportunities for reduction of global memory accesses. The execution speed of matrix multiplication functions can vary by orders of magnitude, depending on the level of reduction of global memory accesses. A CUDA device contains several types of memory that can help programmers improve compute-to-global-memory-access-ratio, defined as the number of floating-point calculations performed for each access to the global memory within a region of a program. Specifically, we focus on global and shared memory. Global memory can be written and read by the device. Shared memory is on-chip memory. Variables that exist in this type of memory can be accessed at very high-speed in a highly parallel manner. Shared memory locations are allocated to thread blocks; all threads in a block can access shared memory variables allocated to that block. Shared memory is an efficient way for threads to share their input data and intermediate results. The global memory is large, but slow,

whereas shared memory is small, but fast. This causes a trade-off in the use of device memories. A common strategy is to partition the data into subsets called tiles, so that each tile fits into the shared memory.

## II. METHODS

### A. Naive Matrix Multiplication PyCuda and PyOpenCL

Unlike an optimized method for matrix multiplication, a naive method does not take advantage of shared memory and tiling. The process for computing matrix multiplication, of two rectangular matrices, using a naive approach can be broken down into the following steps: (1.) initialize the matrices of arbitrary size following the format  $m1 = [M \times N]$  and  $m2 = [N \times M]$ , where after multiplication the resulting matrix is of size  $[N \times N]$ , (2.) create the kernel code for matrix multiplication utilizing simple indexing for each thread and computing the dot product of the rows and columns of the corresponding matrices; global/device memory is used here only (3.) move input and expected output matrices to device memory, (4.) call kernel function to perform naive matrix multiplication, (5.) measure run-time and return the multiplied matrix. These steps were performed for matrices of various sizes created by iteratively increasing the dimension of the initial matrices in both dimensions.

```
self.kernel_code_naive = """
global void naive(float *M, float *N, float *P, const int Widthx, const int Widthy) {
    int Row = blockIdx.y * blockDim.y + threadIdx.y;
    int Col = blockIdx.x * blockDim.x + threadIdx.x;

    P[Row * Widthx + Col] = 0;
    if (Col < Widthx) {
        for (int k=0; k < Widthy; k++) {
            P[Row * Widthx + Col] += M[Row * Widthx + k] * N[k * Widthx + Col];
        }
    }
}
"""
```

Fig. 1. PyCuda Naive Matrix Multiplication Kernel Code

### B. Optimized Matrix Multiplication PyCuda and PyOpenCL

The optimized matrix multiplication was conducted in a similar manner to the naive method, however, the kernel code incorporates tiles and shared memory. The benefit for using shared memory and tiling is that through tiling we are able to partition the data that will be processed by threads and with utilizing shared memory (on-chip memory) we can access these variables quickly. With each tile fitting into

```

self.kernel_code_naive = """
__kernel void naive(__global float *M, __global float *N, __global float *P, const int Widthx, const int Widthy) {
    int index_y = get_global_id(0);
    int index_x = get_global_id(1);

    P[index_y * Widthx + index_x] = 0;
    if (index_x < Widthx) {
        for (int k=0; k < Widthy; k++) {
            P[index_y * Widthx + index_x] += M[index_y * Widthx + k] * N[k * Widthx + index_x];
        }
    }
}
"""

```

Fig. 2. PyOpenCL Naive Matrix Multiplication Kernel Code

the shared memory, all the threads in that tile access these memory variables and we avoid global memory accesses (slower access). Therefore, steps 1-5 above were also followed for the optimized matrix multiplication method. The kernel code differs, however, in that shared memory is declared and shared memory locations are allocated to tiles of thread blocks. The kernel code follows the following structure: (1.) declare shared memory variables, where all threads in a block have access to the elements loaded into these memory variables (one pair for each block), (2.) save threadIdx and blockIdx values into automatic variables, and thus, registers, for fast access, (3.) identify the row and column of the output matrix element to work on/produced by the thread - the horizontal position, or the column index of the P element to be produced by a thread is calculated using  $bx * BLOCKSIZE + tx$  because each block covers  $BLOCKSIZE$  elements in the horizontal (and this logic is also similar for the vertical position, using  $by * BLOCKSIZE + ty$ ) (4.) iterate through all the phases of calculating the P element, with each iteration corresponding to one calculation of a tiled matrix multiplication, (5.) load the appropriate input matrix element into shared memory, (6.) have each thread load an element 'tx' positions away from the beginning point, (7.) ensure all threads have finished loading the tiles of input matrices into shared memory matrix variables before any can move forward (using `__syncthreads()`), (8.) perform one phase of the dot product on the basis of tile elements, (9.) ensure all threads have finished using the input matrix elements in the shared memory before any move on to the next iteration and load elements from next set of tiles, avoiding loading elements too early and corrupting input values for other threads.

### III. RESULTS

#### A. Run-Time Comparisons, Code Equality, Profiling (Naive Matrix Multiplication)

The size of the input matrices in each dimension were iteratively increased to produce a new matrix. For each iteration the CPU, PyCuda kernel, and PyOpenCL kernel execution times for matrix multiplication were recorded. PyCuda optimized and naive matrix multiplication kernel execution times were compared to CPU/serial code. PyOpenCL optimized and naive matrix multiplication kernel execution times were compared to CPU/serial code. The plots for run-time using PyOpenCL and PyCuda can be seen in figures 5 and 6. For consistency

```

# Write the kernel code
self.kernel_code_optimized = """
__global__ void optimized(float *M, float *N, float *P, const int Widthx, const int Widthy)
{
    __shared__ float Mds[(BLOCK_SIZE)x](BLOCK_SIZE)s;
    __shared__ float Nds[(BLOCK_SIZE)s](BLOCK_SIZE)s;

    int tx=threadIdx.x;
    int ty=threadIdx.y;
    int bx=blockIdx.x;
    int by=blockIdx.y;

    int Row = ty * by * (BLOCK_SIZE)s;
    int Col = tx * bx * (BLOCK_SIZE)s;

    float Pvalue=0;

    for(int ph=0; ph< Widthy/(BLOCK_SIZE)s; ph++) {
        if (Row < Widthx && ph*(BLOCK_SIZE)s + tx < Widthy){
            Mds[ty][tx]=M[Row * Widthx + ph * (BLOCK_SIZE)s + tx];
        }
        else{
            Mds[ty][tx]=0;
        }
        if (Col < Widthx && ph*(BLOCK_SIZE)s + ty < Widthy){
            Nds[ty][tx]=N[(ph * (BLOCK_SIZE)s + ty)* Widthx + Col];
        }
        else{
            Nds[ty][tx]=0;
        }
        __syncthreads();

        for (int k = 0; k < (BLOCK_SIZE)s; ++k) {
            Pvalue += Mds[ty][k] * Nds[k][tx];
        }
        __syncthreads();
    }
    if (Row < Widthx && Col < Widthx) {
        P[Row * Widthx + Col] = Pvalue;
    }
}
"""

```

Fig. 3. PyCuda Optimized Matrix Multiplication Kernel Code

```

# Write the kernel code
self.kernel_code_optimized = """
__kernel void optimized(__global float *M, __global float *N, __global float *P, const int Widthx, const int Widthy)
{
    __local float Mds[(BLOCK_SIZE)s](BLOCK_SIZE)s;
    __local float Nds[(BLOCK_SIZE)s](BLOCK_SIZE)s;

    int tx=get_local_id(0);
    int ty=get_local_id(1);
    int bx=get_group_id(0);
    int by=get_group_id(1);

    int Row = ty * by * (BLOCK_SIZE)s;
    int Col = tx * bx * (BLOCK_SIZE)s;

    float Pvalue = 0;

    for(int ph = 0; ph < Widthy / (BLOCK_SIZE)s; ph++) {
        if (Row < Widthx && ph * (BLOCK_SIZE)s + tx < Widthy){
            Mds[ty][tx]=M[Row * Widthx + ph * (BLOCK_SIZE)s + tx];
        }
        else{
            Mds[ty][tx]=0;
        }
        if (Col < Widthx && ph * (BLOCK_SIZE)s + ty < Widthy){
            Nds[ty][tx]=N[(ph * (BLOCK_SIZE)s + ty)* Widthx + Col];
        }
        else{
            Nds[ty][tx]=0;
        }
        __syncthreads();

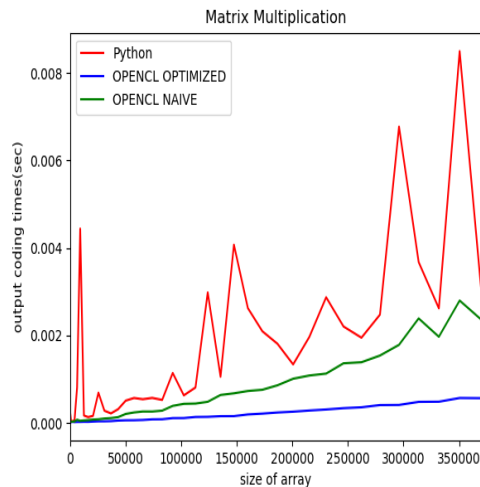
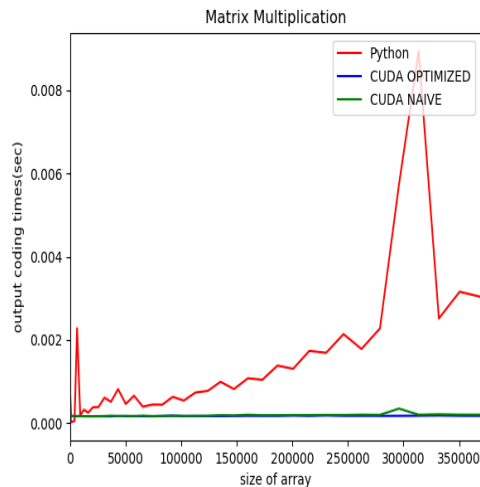
        for (int k = 0; k < (BLOCK_SIZE)s; ++k) {
            Pvalue += Mds[ty][k] * Nds[k][tx];
        }
        __syncthreads();
    }
    if (Row < Widthx && Col < Widthx) {
        P[Row * Widthx + Col] = Pvalue;
    }
}
"""

```

Fig. 4. PyOpenCL Naive Matrix Multiplication Kernel Code

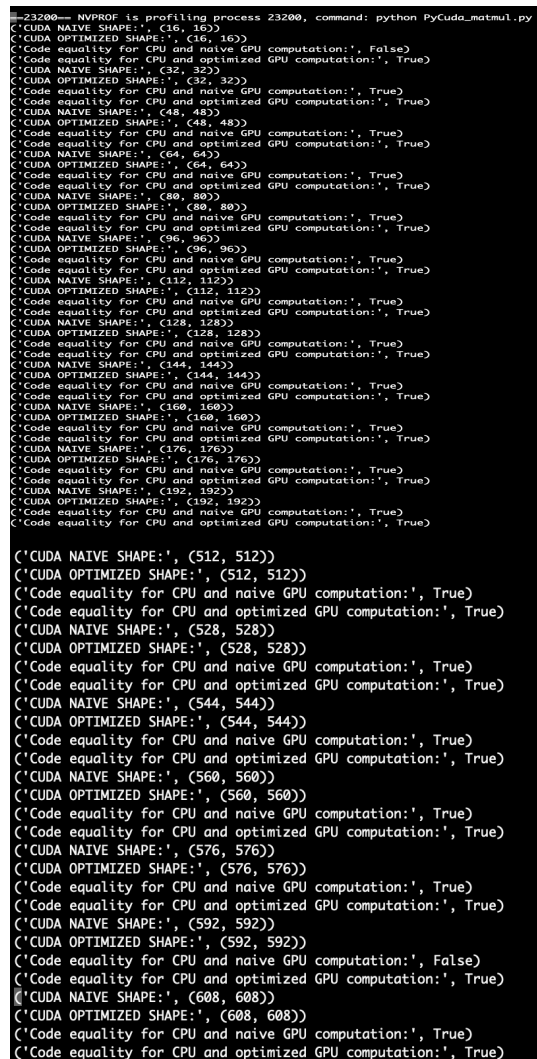
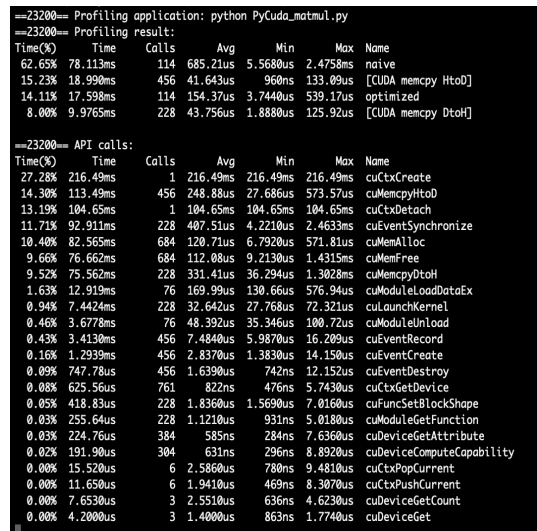
between PyOpenCL and PyCuda, the figures below were produced using output matrices beginning with dimension 16 x 16 iterating each dimension by a factor 1 to 38, with the largest matrix of size 608 x 608. The horizontal axis is the number of elements in the matrix during each iteration. The vertical axes are the run times in seconds. As stated earlier, the kernel execution time for matrices of much larger size, i.e. 8000x8000, was able to be recorded in PyOpenCL, without encountering memory errors. PyCuda optimized kernel code was able to compute matrices of this size as well without encountering memory issues, however, the naive method was limited to matrices of size 608 x 608.

Figure 7 is a screenshot of the profiling for PyCuda. The Visual Profiler is a graphical profiling tool that displays a timeline of CPU and GPU activity, and that includes an automated analysis engine to identify optimization opportunities.



## IV. DISCUSSION

In this assignment, matrix multiplication was successfully implemented using three different methods: 1. CPU/Serial, 2. PyCuda (naive + optimized), and 3. PyOpenCL (naive + optimized). From the results, it is clear that CPU matrix multiplication operation is much slower compared to PyCuda and PyOpenCL. The naive PyCuda method and the optimized PyCuda method had similar execution times for matrices of dimension 16 x 16 as large as 608 x 608. It would be interesting to see if the execution times would begin to diverge with much larger matrices. However, the execution time for PyOpenCL naive was longer than PyCuda. This is to be expected, since reducing global memory accesses and utilizing tiled matrix multiplication with shared memory accesses should speed up



```

('OPENCL NAIVE SHAPE:', (16, 16))
('OPENCL OPTIMIZED SHAPE:', (16, 16))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (32, 32))
('OPENCL OPTIMIZED SHAPE:', (32, 32))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (48, 48))
('OPENCL OPTIMIZED SHAPE:', (48, 48))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (64, 64))
('OPENCL OPTIMIZED SHAPE:', (64, 64))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (80, 80))
('OPENCL OPTIMIZED SHAPE:', (80, 80))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (96, 96))
('OPENCL OPTIMIZED SHAPE:', (96, 96))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (112, 112))
('OPENCL OPTIMIZED SHAPE:', (112, 112))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (128, 128))
('OPENCL OPTIMIZED SHAPE:', (128, 128))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (144, 144))
('OPENCL OPTIMIZED SHAPE:', (144, 144))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (160, 160))
('OPENCL OPTIMIZED SHAPE:', (160, 160))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (176, 176))
('OPENCL OPTIMIZED SHAPE:', (176, 176))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (192, 192))
('OPENCL OPTIMIZED SHAPE:', (192, 192))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (208, 208))

('OPENCL OPTIMIZED SHAPE:', (512, 512))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (528, 528))
('OPENCL OPTIMIZED SHAPE:', (528, 528))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (544, 544))
('OPENCL OPTIMIZED SHAPE:', (544, 544))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (560, 560))
('OPENCL OPTIMIZED SHAPE:', (560, 560))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (576, 576))
('OPENCL OPTIMIZED SHAPE:', (576, 576))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (592, 592))
('OPENCL OPTIMIZED SHAPE:', (592, 592))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)
('OPENCL NAIVE SHAPE:', (608, 608))
('OPENCL OPTIMIZED SHAPE:', (608, 608))
('Code equality for CPU and naive GPU computation:', True)
('Code equality for CPU and naive GPU computation:', True)

```

Fig. 9. PyOpenCL Code Equality Between CPU and Naive and Optimized Matrix Multiplication Kernel Code

kernel execution time.

## REFERENCES

- [1] David B. Kirk and Wen-mei W. Hwu. 2016. Programming Massively Parallel Processors, Third Edition: A Hands-On Approach (3rd ed.). Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.