# Tiled and Naive 2-D Convolution and 1-D Histogram Calculation With Cuda and OpenCL

Drew Afromsky
Email: daa2162@columbia.edu

*Abstract*—**2D convolution and 1D histogram calculation was performed in both CUDA and OpenCL. 2D convolution was implemented, taking advantage of both shared memory/tiles and global memory (naive methods). Tiled 2D convolution was performed in CUDA only. For naive 2D convolution, the input to the algorithm is an [M X N] matrix and a [K X K] kernel of odd dimension sizes. The output after convolution remained the same size as the input; zero padding was performed to take into account halo/ghost cells, when the kernel was "acting" on "non-existent" pixels/matrix elements. For the tiled 2D convolution, the kernel size was fixed: [5 x 5] to avoid dealing with dynamic memory allocation. For both methods, a serial implementation of 2D convolution was performed using scipy function (signal.convolve2D). Execution times for 2D convolution CUDA naive, 2D convolution CUDA tiled, and 2D convolution serial were recorded and plotted for comparison. Execution times for 2D convolution in OpenCL were compared to 2D convolution in serial and plotted as well. Matrices for both CUDA and OpenCL were initialized and iteratively increased in size in both dimensions by the same factor. Similarly, for 1D histogram calculation, execution times were recorded in both OpenCL and CUDA, and serial code, and then plotted and compared. One plot was created for CUDA versus serial. A second plot was created for OpenCL versus serial.**

*Index Terms*—**CUDA, OpenCL, naive, tiled, convolution, histogram, halo/ghost cell**

## I. INTRODUCTION

Convolution is a popular array or matrix operation performed in various forms of signal processing, digital recording, image processing, video processing, and computer vision. In these application areas, convolution is often performed as a filter that transforms signals and pixels into more desirable values. Examples include convolution for image blurring and Gaussian filters for sharpening boundaries and edges of objects in images. In high-performance computing, convolution pattern is referred to as stencil computation. Output data elements can be calculated independently of each other, which is a desirable trait for parallel computing. There also exists a substantial level of input data sharing among output data elements with boundary conditions. More specifically, convolution is an array operation where each output data element is a weighted sum of a collection of neighboring input elements. The weights are defined by a convolutional kernel, also referred to as a mask array/matrix. The size of this kernel for 1D arrays is an array of length 2*n+1, and for 2D arrays [2n+1 x 2n+1], to ensure odd dimensions. Because convolution is defined in terms of neighboring elements, boundary conditions arise for output elements close to the ends of an array or matrix. For handling such boundary condition cases, zero padding is usually performed. Zero padding is essentially filling in missing cells or matrix elements, known as ghost or halo cells, with a zero value.

A histogram is a display of the frequency of data items in successive numerical intervals. The value intervals are plotted along the horizontal axis and the frequency of data items in each interval is represented as the height of a bar. Histograms provide useful summaries of data sets. Histograms of different types of object images, such as faces versus cars, as seen in computer vision tasks, tend to exhibit different shapes. By dividing an image into subareas and analyzing the histograms for these subareas, one can quickly identify interesting subareas of an image that may contain the objects of interest. This is the basis of feature extraction.

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. 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. We also introduce constant memory, where variables in this type of memory are also visible to all thread blocks. In contrast to global memory, constant memory cannot be changed by threads during kernel execution. Kernel functions also access constant memory variables as global variables, so their pointers don't need to be passed to the kernel as parameters. Constant memory variables are located in DRAM like global memory variables. CUDA-enabled hardware caches the constant memory variables during kernel execution. As a result, when all threads in a warp access the same constant memory variable, as in the case with convolutional kernels, the caches can provide tremendous amount of bandwidth to satisfy the data needs of threads.

## II. METHODS

### A. Naive 2D Convolution PyCuda and PyOpenCL

Unlike an optimized/tiled method for 2D convolution, a naive method does not take advantage of shared memory and tiling. The process for computing convolution, of two matrices (input matrix and convolutional kernel), using a naive approach can be broken down into the following steps:

(1.) initialize the input matrix of arbitrary size following the format input_matrix = [M X N] and convolutional_kernel = [K X K], where K=2*n+1 and n is an integer. After convolution, the output matrix retains the same shape as the input matrix. (2.) create the kernel code for 2D convolution utilizing simple indexing for each thread and computing the weighted sum of pixels/matrix elements of the input array with those of the convolutional kernel; global/device memory is used here only (3.) move input and expected output matrices to device memory, (4.) compute block and grid size (5.) call kernel function to perform 2D naive convolution, (6.) measure runtime and return the output - a matrix resulting from the convolution of input matrix and convolutional kernel. These steps were performed for matrices of various sizes created by iteratively increasing the dimension of the initialized matrix in both dimensions. For naive 2D convolution, kernel size was also iteratively changing in size. Convolutional kernel size was kept constant for 2D tiled convolution.



Fig. 1. PyCuda Naive 2D Convolution Key Steps



Fig. 2. PyCuda Naive 2D Convolution Kernel Code



Fig. 3. PyOpenCL Naive 2D Convolution Key Steps



Fig. 4. PyOpenCL Naive 2D Convolution Kernel Code

### B. Optimized/Tiled 2D Convolution PyCuda

The optimized/tiled 2D convolution 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 the shared memory, all the threads in that tile access these memory variables and we avoid global memory accesses (larger, but slower access). Therefore, steps 1-6 above were also followed for the tiled 2D convolution 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.) Each thread of the kernel first calculates the y and x indices of its output element - col_o and row_o variables of the kernel, (2.) each thread then calculates the y and x indices of the input element it is to load into the shared memory by subtracting mask_width/2 from row_o and col_o and assigning the results to row_i and col_i, (3.) 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); loading input tiles into shared memory, (4.) check if the y and x indices of the input tiles are within range of the input; if not, the

input element it is attempting to load is a ghost element and a 0.0 value should be placed into shared memory, (5.) compute the output value using input elements in shared memory, (6.) ensure all threads have finished loading the tiles of input matrices into shared memory matrix variables before any can move forward (using __syncthreads), (7.) ensure that only the threads whose indices are smaller than tile_size should participate in the calculation of output pixels, (8.) 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, (9.) all threads whose output elements are in the valid range write their result values into their respective output elements

## C. 1D Histogram Calculation PyCuda and PyOpenCL

1D histogram calculation can be broken down into the following steps:

(1.) initialize the input array of length N and containing randomly sampled integers between 0 and 255. (2.) create the kernel code for 1D histogram calculation (3.) move input and expected output matrices to device memory, (4.) compute block and grid size (5.) call kernel function to perform 1D histogram calculation, (6.) measure run-time and return the output - a histogram of the input array. These steps were performed for arrays of various sizes created by iteratively increasing the dimension of the initialized array. Bin size of 26 was kept constant.



Fig. 7. PyOpenCL 1D Histogram Key Steps



Fig. 5. PyCuda Tiled 2D Convolution Key Steps



Fig. 8. PyOpenCL 1D Histogram Key Steps

## III. RESULTS

### A. Run-Time Comparisons, Code Equality, Profiling (PyCuda 2D Convolution)

The size of the input matrices in each dimension were iteratively increased to produce a new matrix during 2D convolution (both naive and tiled). For each iteration, the CPU, PyCuda kernel, and PyOpenCl kernel execution times for 2D convolution were recorded. PyCuda tiled and naive



Fig. 6. PyCuda Tiled 2D Convolution Kernel Code

```
self.kernel_hist = """
__kernel void hist(__global const int *buffer, __global int *hist, const int size)

{
int i = get_global_id(0);

if (i < size && i >= 0)
    {
    int bin = buffer[i];
    atomic_add(&hist[bin], 1);
    }
}
"""
```

Fig. 9. PyOpenCL 1D Histogram Kernel Code

```
self.kernel_hist = """
__global__ void hist(const int *buffer, int *hist, const int size)

{
int i = threadIdx.x + blockIdx.x * blockDim.x;

if (i < size && i >= 0)
    {
    int bin = buffer[i];
    atomicAdd(&hist[bin], 1);
    }
}

"""
```

Fig. 10. PyCuda 1D Histogram Kernel Code

2D convoution kernel execution times were compared to CPU/serial code execution times. PyOpenCL and 2D convolution kernel execution times were compared to CPU/serial code execution times. For consistency between PyOpenCL and PyCuda, the figures below were produced using the same size input matrices beginning with dimension [20 x 25] iterating each dimension by a factor 1 to 39. The horizontal axis is the number of elements in the newly produced matrix during each iteration. The vertical axes are the run times in seconds. From the plots, we see that both the CUDA naive and tiled methods run faster than OpenCL, and naive and tiled methods for both OpenCL and CUDA run much faster than serial/CPU code. Similarly, for 1D histogram, the array was iteratively increased.

Figure 15 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

2D convolution and 1D histogram calculation was successfully implemented using three different methods: 1. CPU/Serial, 2. PyCuda (naive + tiled), and 3. PyOpenCL (naive). Tiled 2D convolution was only performed in PyCuda. From the results, it is clear that CPU histogram operation and 2D convolution are much slower compared to PyCuda and PyOpenCL. The naive PyCuda method execution time was slower than the tiled/optimized PyCuda execution time. This is to be expected, since reducing global memory accesses and utilizing tiled matrix multiplication with shared memory accesses should speed up kernel execution time. The execution time for PyOpenCL naive was much faster than CPU and slower than both the naive and tiled/optimized versions in Py-
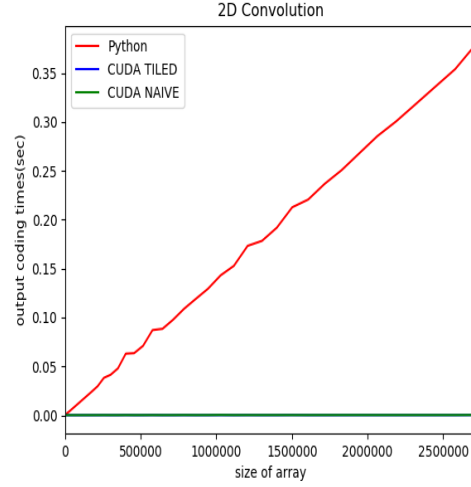


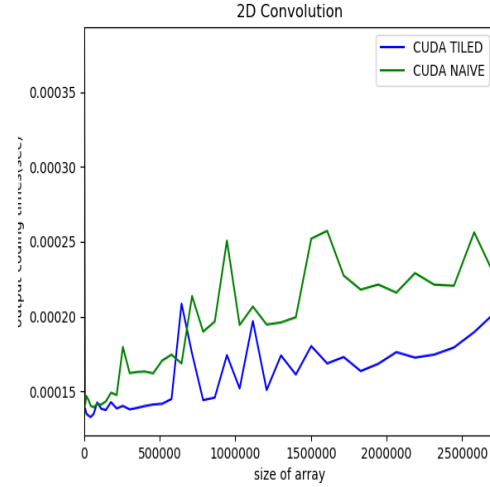Fig. 11. PyCuda Execution Times for CPU, Naive, and Tiled 2D Convolution



Fig. 12. Closer Look at PyCuda Execution Times for Naive and Tiled 2D Convolution ONLY

Cuda. Similarly for 1D histogram calculation, CUDA appeared to outperform OpenCL and both outperformed serial/CPU implementation.

## 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.
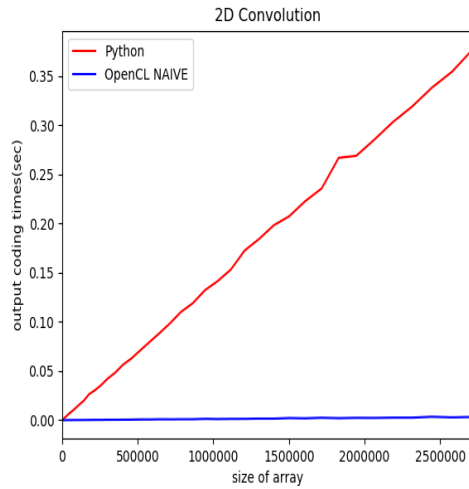
Fig. 13. PyOpenCL Execution Times for CPU and Naive 2D Convolution
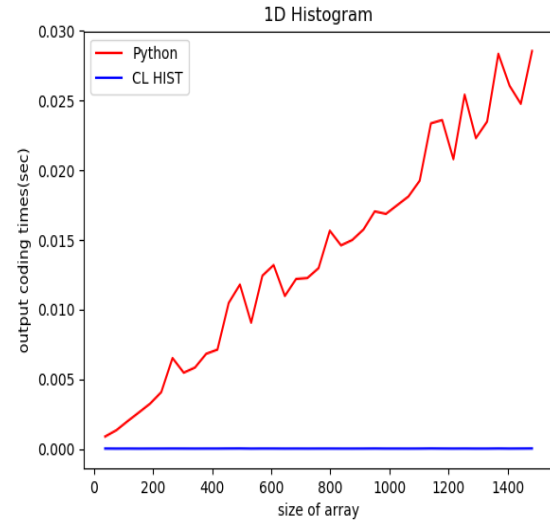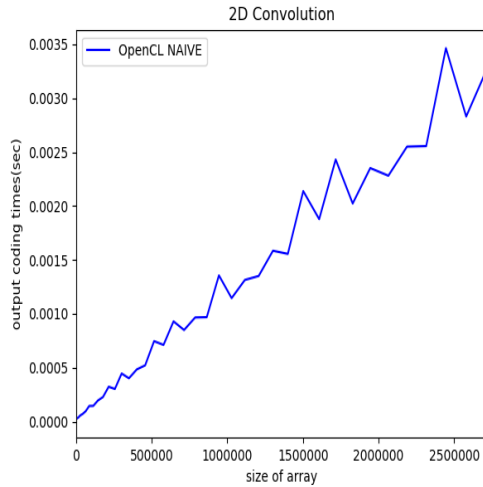


Fig. 14. Closer look at PyOpenCL Execution Time for Naive 2D Convolution



Fig. 15. Profiling for Tiled and Naive PyCuda 2D Convolution and 1D Histogram
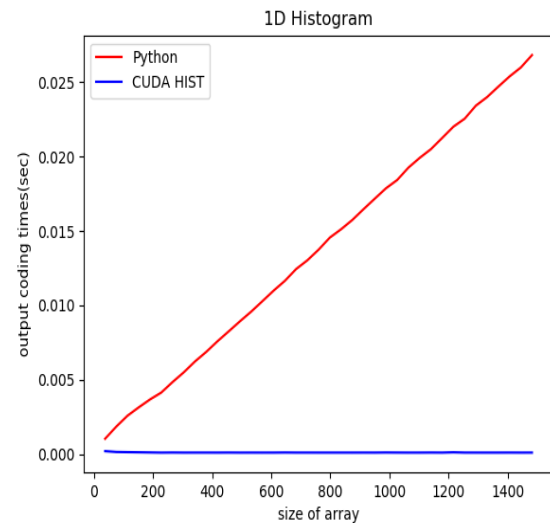


Fig. 16. PyOpenCL 1D Histogram Execution Time



Fig. 17. PyCuda 1D Histogram Execution Time