

Accelerating Data Augmentation for Deep Learning with Cuda and OpenCL

Drew Afromsky
Email: daa2162@columbia.edu

Abstract—Two different data augmentation techniques were explored in this project. Specifically, image blurring and image transpose were implemented in Cuda and OpenCL. The data augmentation techniques were conducted on two different datasets. Dataset 1 is comprised 100 pathology images, each of size 3100 x 3100. This was considered a small dataset, but with large matrices. Dataset 2 is a large dataset consisting of 50,000 images, each of size 32 x 32.

Image blurring utilized both shared memory (optimized, with tiles) and global memory (naive methods). The optimized, tiled image blurring technique was performed in Cuda only. For image blurring, an average image blur was executed, with two different kernel sizes: 3 x 3 for dataset 1 and 11 x 11 for dataset 2 (this also avoids dealing with dynamic memory allocation). For the naive blur kernel, the input to the algorithm is a matrix of size 3100 x 3100 (an image from dataset 1) or a matrix of size 32 x 32 (an image from dataset 2). The output after blurring 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. A serial implementation of image average blur was performed using an openCV function (`cv2.blur`)

Similarly, image transpose was conducted solely using naive methods in both Cuda and OpenCL on both datasets. These methods were compared to a serial implementation using a numpy function (`numpy.transpose`).

Execution times for image blur in Cuda (naive), image blur Cuda (tiled), and image blur (serial) were recorded. Execution times for image blur in OpenCL (naive) were also recorded.

Four different plots were created:

Plot 1 demonstrates the four different execution times (Cuda Naive, Cuda Tiled, OpenCL, and Serial) for image blur on dataset 1 (100 pathology images of size 3100 x 3100). Plot 2 demonstrates the four different execution times (Cuda Naive, Cuda Tiled, OpenCL, and Serial) for image blur on dataset 2 (50,000 natural images of size 32 x 32). Plot 3 demonstrates the four different execution times (Cuda Naive, Cuda Tiled, OpenCL, and Serial) for image transpose on dataset 1. (100 pathology images of size 3100 x 3100). Plot 4 demonstrates the four different execution times (Cuda Naive, Cuda Tiled, OpenCL, and Serial) for image transpose on dataset 2. (50,000 natural images of size 32 x 32).

Index Terms—Cuda, OpenCL, blur, transpose, naive, tiled, kernel, halo/ghost cell

I. INTRODUCTION

Image augmentation is a popular technique to represent the same image in a different form. For example, a common image augmentation technique is flipping vertically or horizontally, where the image appears to be "mirrored" or "flipped" along an axis after this augmentation is completed. More specifically, image augmentation is regularly used to increase the size of a dataset that is later used for training algorithms utilizing deep learning. In computer vision applications using deep

learning, these algorithms tend to be data hungry, requiring many example images to learn from. Data or image augmentation is an effective way to increase the size of a dataset while maintaining the same kind or type of image examples the algorithm learns from. Unfortunately, many of these augmentation techniques are done on the CPU, which is proven as computationally intensive and time consuming. Consider an image dataset with 5000 images each of size 3100 x 3100. A single augmentation technique applied to all images would increase the size of this dataset by 5000, with a new dataset containing augmented images and a total of 10,000 images. Now consider applying 5 augmentation techniques. That brings the total dataset size to 30,000 images. Some augmentation techniques can be executed quickly, on the order of milliseconds per image and sometimes quicker. Now imagine, it takes 240 seconds (4 minutes) to conduct a single data augmentation technique on a single image. If we wish to apply that to every image in a dataset of size 5,000, that would take a minimum of 2,0000 minutes (33+ hours) not taking into account the read/write operations to save the images. This process can be grueling and time consuming, often leading to deep learning pipelines that are unnecessarily long. This presents a ripe opportunity to apply parallel processing to various image augmentation techniques. Output data elements can be calculated independently of each other, which is a desirable trait for parallel computing.

In this project, image blurring and image transpose were explored, where the blurring operation 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 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 $[2n+1 \times 2n+1]$, to ensure odd dimensions. Because image blur 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 pixel value.

The significance of transposing a matrix is mainly to represent linear transformations such as rotation and scaling. By taking the transpose of a matrix that represents a linear transformation, properties of the transformation can be revealed. If the transpose of a matrix and the original matrix are equivalent, for example, then the transformation is a symmetric transformation and the matrix is symmetric.

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 image blur kernels, the caches can provide tremendous amount of bandwidth to satisfy the data needs of threads.

II. METHODS

A. Block Diagram

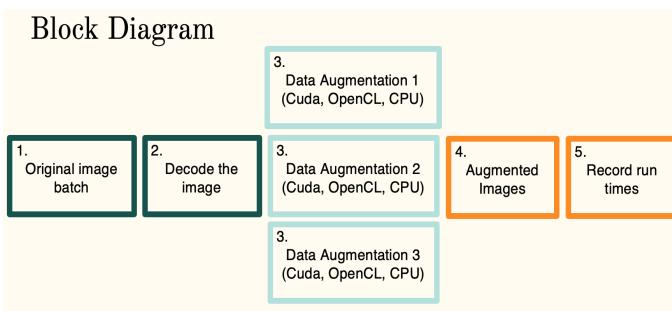


Fig. 1. Block Diagram for Image Augmentation with OpenCL and Cuda

B. Naive Image Blur Cuda and OpenCL

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

(1.) Initialize the input matrix with size of the dataset of interest (either 3100 x 3100, or 32 x 32) blur_kernel = [K X K], where K=2*n+1 and n 1 for dataset 2 (32 x 32 images) and 5 for dataset 1 (3100 x 3100 images). After blurring, the output matrix/image retains the same shape as the input

matrix. (2.) Create the kernel code for image blurring utilizing simple indexing for each thread and computing the weighted sum of pixels/matrix elements of the input array with those of the blur kernel; global/device memory is used here only. (3.) Move input and expected output matrices to device memory; matrices here refer to a single image matrix. (4.) Compute block and grid size based on the shape of the input matrices and the optimal block size (32 ,32, 1) for 2D. (5.) Call kernel function to perform image blurring (6.) Measure run-time and return the output - a matrix resulting from the convolution of input matrix and blur kernel.

These steps were performed for both datasets.

For naive image blurring and tiled blurring, convolutional kernel size was kept constant.

C. Optimized/Tiled Image Blurring Cuda

The optimized/tiled 2D image blur 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 image blurring 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

D. Naive Image Transpose Cuda and OpenCL

We take a 2D matrix of any size and transpose it, which means the rows become the columns and the columns become the rows.

The steps taken to produce a transposed image are similar to those described for image blurring.

(1.) Declare host variables a. Input matrix, shapes of this matrix, and a zeros matrix that will be the shape of the desired output and filled with the values calculated on the device (2.) Allocate memory on the device/GPU for the input matrix and the output matrix (3.) Create the kernel code/function b. We have 32x32 threads per block, and block size is 32x32, and grid size is calculated according to the size of the input array in the x-dimension divided by 32 and the same is done in the y-dimension. (4.) Similar as before, we launch and call the kernel code declaring the block and grid size, after transferring input data from the host (CPU) to the device (GPU). (5.) Call kernel function to perform image transpose (6.) To retrieve our desired output, the transposed matrix, we must then transfer the data on the GPU back to the CPU. (7.) Measure run-time and return the output - a matrix resulting from the transpose of the input matrix.

III. RESULTS

A. Image Blurring and Image Transpose)

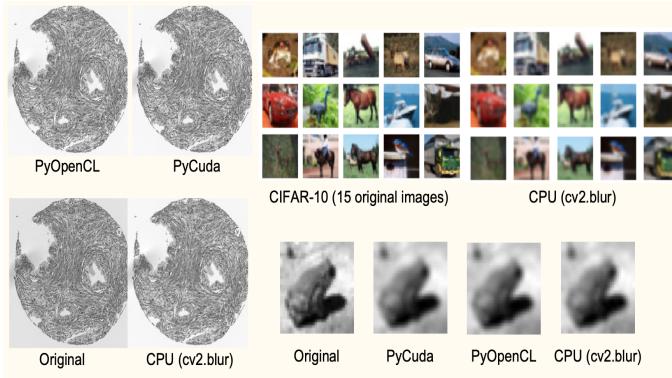


Fig. 2. Image Blurring with Cuda, OpenCL, and OpenCV

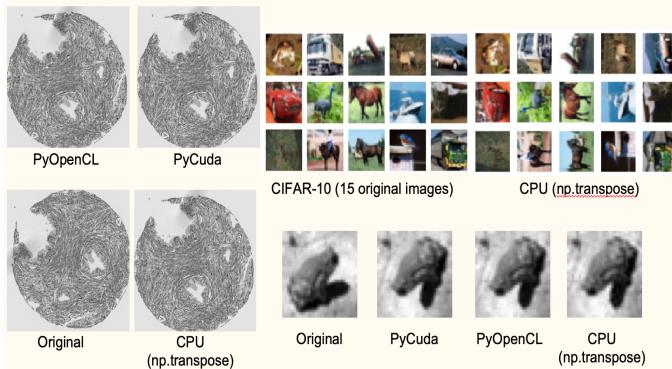


Fig. 3. Image Transpose with Cuda, OpenCL, and Numpy Transpose

Results: Image Blur Execution Times

PyOpenCL	PyCuda	CPU (cv2.blur)
Dataset 1 (Pathology Slides): N = 100 , Size = 3100 x 3100 Execution time: 1.067022144 seconds Dataset 2 (CIFAR-10): N= 50,000 images, Size = 32 x 32 Execution time: 3.4752e-05 sec/image*50,000 images=1.7376 seconds	Dataset 1 (Pathology Slides): N = 100 , Size = 3100 x 3100 Execution time: 0.02425217628479004 seconds Dataset 2 (CIFAR-10): N= 50,000 images, Size = 32 x 32 Execution time: 6.454583406448364 seconds 6.834920644760132 (Tiled) seconds	Dataset 1 (Pathology Slides): N = 100 , Size = 3100 x 3100 Execution time: 3.989729404449463 seconds Dataset 2 (CIFAR-10): N= 50,000 images, Size = 32 x 32 Execution time: 1.3029658794403076 seconds

Fig. 4. Image Blur Execution Times for Dataset 1 and Dataset 2

Results: Image Transpose Execution Times

PyOpenCL	PyCuda	CPU (np.transpose)
Dataset 1 (Pathology Slides): N = 100 , Size = 3100 x 3100 Execution time: 0.000852672 sec/image*100 images = 0.0852672 seconds Dataset 2 (CIFAR-10): N= 50,000 images, Size = 32 x 32 Execution time: 2.6592e-05 sec/image*50000 images = 1.3296 seconds	Dataset 1 (Pathology Slides): N = 100 , Size = 3100 x 3100 Execution time: 0.0421476364136 seconds Dataset 2 (CIFAR-10): N= 50,000 images, Size = 32 x 32 Execution time: 0.001200914383 sec/image *50,000 images = 60.0457191465 seconds	Dataset 1 (Pathology Slides): N = 100 , Size = 3100 x 3100 Execution time: 0.2578585147857666 seconds Dataset 2 (CIFAR-10): N= 50,000 images, Size = 32 x 32 Execution time: 0.17492103576660156 seconds

Fig. 5. Image Transpose Execution Times for Dataset 1 and Dataset 2

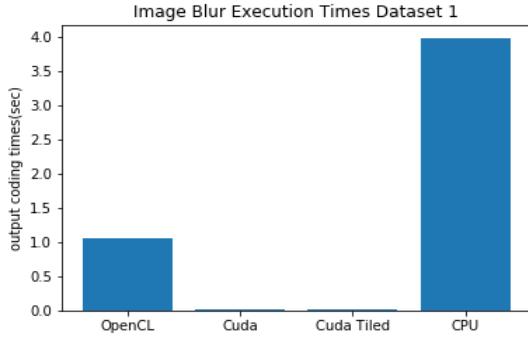


Fig. 6. Bar Plot for Image Blurring on Dataset 1

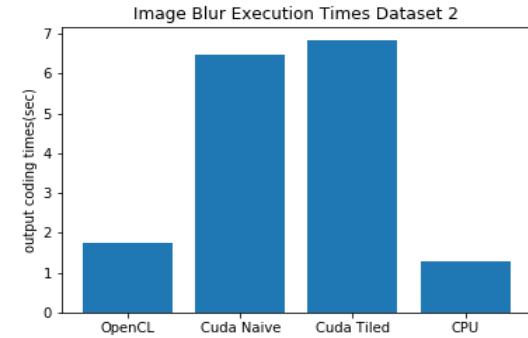


Fig. 7. Bar Plot for Image Blurring on Dataset 2

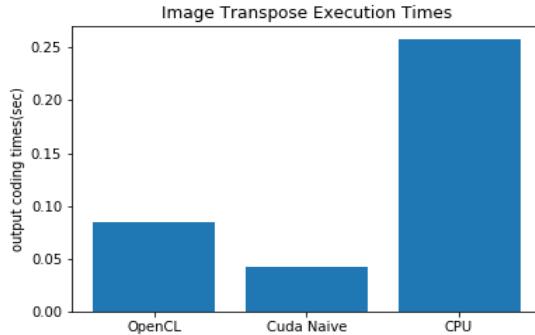


Fig. 8. Bar Plot for Image Transpose on Dataset 1

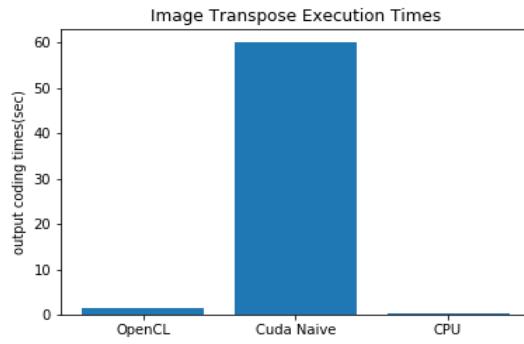


Fig. 9. Bar Plot for Image Transpose on Dataset 2

IV. DISCUSSION

In regards to figures 6 and 7, we see that tiled method for image blurring using PyCuda was the fastest at executing the task for dataset 1. However, for dataset 2, we see that cpu was the fastest. I believe that for larger images and smaller batches Cuda performs better, whereas for smaller images and larger batches, execution time increases with Cuda. OpenCL seems to remain steady regardless of batch size and image size.

In regards to figures 8 and 9, for dataset 1 we see that cuda outperforms both cpu and openCL on image transpose, where the images are large, but the batches are small. However, Cuda seems to take longer on larger batch sizes, and smaller matrices. Cpu outperforms both openCL and Cuda for these images/batches.

In general, we see that Cuda executes the quickest on larger matrices and smaller batches. OpenCL seems to work best with larger batch sizes and smaller image sizes. CPU performs the best on image transpose, when the batch size is large and images are small.

It would have been interesting to also compare these methods for large batch sizes with large matrices and small batch sizes and small matrices. Additionally, if given more time, it would have been interesting to compare to data augmentation tools publicly available such as albumentations, NVIDIA's DALI, and openCV's GPU-compatible framework.

I believe the tiled image blurring may have some issues in terms of ability to work with larger matrices, since some of the output images appeared to be partial images. Also, an attempt to generate the Cuda code for image rotation was done, but not completed and no image outputs were obtained.

V. APPENDIX

```

... ssh -X daa2162@tesseract.ee.columbia.edu .d1a2162 ... bash
Drewfrosky@daa2162:~$ ssh -X daa2162@tesseract:~/train -ssh X daa2162@tesseract
File: slurm-126794.out

NUmernode 2,5,3
File: slurm-126794.out

#==23764-- NVPProf is profiling process 23764, command: python PyCuDNN_Blu.py
/home/daa2162/train
d2cnn
50000
C'NTILED Execution Times: : 6.454583406448364
C'TILED Execution Times: : 6.834926644768132
#==23764-- Profiling result:
Time(s)      Time       Calls    Avg     Min     Max      Name
 8.65        10.000000  13.316470 32.928000 13.000000  cuModuleAlloc
 8.17        247.25ms   50000   4.94400s  4.83200s  8.67200s  cuModuleBlur
 17.10        100000000  2.393600 2.24900s 15.80000s  [CUDA] memcpY DtoH
 15.32        100000000  1.390000 1.02400s 11.07000s  [CUDA] memcpY HtoD

#==23764-- API calls:
Time(s)      Time       Calls    Avg     Min     Max      Name
 3.83        26.700000  2500000  10.680us 11.980us  3.700000  cuMemcpyHtoD
 17.10        13.90075  100000000 13.190000 125.61us  7.3158ms  cuModuleLoadDataEx
 5.34        4.37185  150000000 28.961us 22.892us  45.217ms  cuMemcpyDtoD
 4.39%      3.570595  100000000 35.790000 31.653us  43.836ms  cuMemcpyDtoH
 1.92%      1.563245  200000000 7.2160us 5.7580us  664.99us  cuEventRecord
 0.67%      545.72ms   200000000 2.722000 7.880us  699.99us  cuEventCreate
 0.39%      319.95ms   200000000 1.590000 7.716ns  641.86us  cuEventDestroy
 0.39%      305.40ms   35000001  860ns   462ns   402.630us  cuEventynchronize
 0.17%      136.20ms   100000000 1.3620us 1.0440us  640.54us  cuFuncSetBlockShape
 0.17%      136.20ms   100000000 1.3620us 1.0440us  640.54us  cuFuncSetGridSize
 0.13%      186.44ms   100000000 1.0640us 836ns   602.34us  cuModuleGetFunction

#d Get Help      #d Write Out      #d Where Is      #d Cut Text      #d Read 39 Lines J
#d Exit          #d Read File      #d Replace      #d Uncut Text     #d To Spell      #d Cur Po
```

```
  ssh -X daa2162@tesseractee.columbia.edu ~daa2162@tesseractee:~/train -ssh -x -d
GNU nano 2.5.3
File: slurm-126792.out

==20261-- NVPROF is profiling process 20261, command: python PyCuda_Blu...dr
/home/daa2162/train
[0.00, 31000]  seconds
100
100
[NATIVE Execution Times: 0.02425217628479004]
[C-TILED Execution Times: 0.02221885147949727]
==20261-- Profiling application: python PyCuda_Blu.py
Time(s)      Time          Calls      Avg      Min      Max      Name
Time(s)      Time          Calls      Avg      Min      Max      Name
Time(s)      Time          Calls      Avg      Min      Max      Name
3.67 3.67031s 200 17.352ms 0.3-3.67ms 0.0-1.00ms 0.0-1.00ms cuMemcpyDtoH
33.79% 3.13653s 300 10.379ms 1.0-40.80ms 27.46ms 0.0-1.00ms [CUUDA memcpy HtoD]
13.62% 1.09280s 100 10.928ms 1.0-10.92ms 1.0-9.57ms image_blur
3.52% 2.82.85ms 100 2.8285ms 2.457msns 2.8294ms 0.0-1.00ms blur_tiled
PyCuda_Blu.py
Time(s)      Time          Calls      Avg      Min      Max      Name
Time(s)      Time          Calls      Avg      Min      Max      Name
Time(s)      Time          Calls      Avg      Min      Max      Name
37.61% 3.67031s 200 15.352ms 14.435msns 34.383msns 0.0-1.00ms cuMemcpyDtoH
35.00% 3.13653s 300 15.352ms 14.435msns 34.383msns 0.0-1.00ms cuMemcpyDtoH
14.06% 3.137234s 6 5.8615msns 2.7893msns 16.950msns 0.0-1.00ms cuEventSynchronize
7.01% 6.83..64ms 500 1.3673msns 333.36msns 12.2-825msns 0.0-1.00ms cuMemAlloc
2.13% 207.61msns 1 207.61msns 207.61msns 207.61msns 0.0-1.00ms cuCtxCreate
1.09% 106.46msns 1 106.46msns 106.46msns 106.46msns 0.0-1.00ms cuModuleLoadDataEx
0.98% 1.09.40msns 200 5.4505msns 0.0-5.45msns 0.0-5.45msns 0.0-1.00ms cuModuleLoadDataEx
0.12% 11.430msns 200 57.147usns 0.0-54.93usns 16.44usns 0.0-1.00ms cuLaunchKernel
0.08% 2.2388msns 400 5.5890usns 0.0-5.5890usns 0.0-5.5890usns 0.0-1.00ms cuEventRecord
0.04% 3.7260msns 400 9.3010usns 0.0-9.3010usns 31.628usns 0.0-1.00ms cuEventCreate
0.02% 2.2388msns 400 1.4436usns 0.0-1.4436usns 284.91usns 0.0-1.00ms cuEventDestroy
0.01% 0.00.00msns 1 0.0000usns 0.0-0.0000usns 0.0-0.0000usns 0.0-1.00ms cuModuleGetFunction
0.00% 4.08.90msns 200 2.0430usns 0.0-1.6189usns 3.3700usns 0.0-1.00ms cuModuleGetFunction
[Read 39 lines]
```

```
● DrewAfrormsky — daa2162@tesseract.ee.columbia.edu: ~ /train — ssh -x daa2162
File: slurm-126843.out

GNU nano 2.5.3

--14414-- NVPROF is profiling process 14414, command: python PyCuda_Image_Trans.py
[32, 32]
0.00120891438293
--14414-- Profiling application: python PyCuda_Image_Trans.py
--14414-- Profiling result:
Time(s)      Avg Time    Calls   Avg Cycles  Min   Max  Name
0.000000  2.6800us  1  2.6800us  2.6800us  2.6800us  [CUDA memcpym HtoD]
35.49%  2.2720us  1  2.2720us  2.2720us  2.2720us  [CUDA memcpym DtoH]
22.25%  1.4400us  1  1.4400us  1.4400us  1.4400us  [CUDA memcpym DtoH]

--14414-- API calls:
Time(s)      Avg Time    Calls   Avg Cycles  Min   Max  Name
0.000000  214.21ms  1  214.21ms  214.21ms  214.21ms  cuCtxCreate
35.60%  119.49ms  1  119.49ms  119.49ms  119.49ms  cuCtxDetach
0.18%  109.99us  1  109.99us  109.99us  109.99us  cuMemAlloc
0.12%  109.99us  2  274.99us  19.32us  53.70us  cuMemFree
0.12%  109.99us  2  202.59us  49.327us  35.86us  cuMemFree
0.02%  58.59us  1  58.59us  58.59us  58.59us  cuLaunchKernel
0.01%  41.66us  1  41.66us  41.66us  41.66us  cuMemcpyHtoB
0.01%  37.72us  1  37.72us  37.72us  37.72us  cuEventRecord
0.01%  19.360us  2  9.6800us  9.6800us  9.7510us  cuEventRecord
0.00%  11.46us  2  5.7230us  2.3370us  9.1090us  cuEventCreate
0.00%  10.46us  4  5.2100us  2.5050us  9.1090us  cuEventCreate
0.00%  7.1270us  4  1.7810us  0.4640us  5.4360us  cuCtxGetCurrent
0.00%  6.2120us  9  690ns  378ns  1.4330us  cuDeviceGetAttribute
0.00%  6.0000us  3  2.0000us  1.0000us  3.0000us  cuEventSynchronize
0.00%  1.1130us  4  1.5280us  592.9600us  2.7070us  cuCtxPopCurrent
0.00%  5.3180us  1  5.3180us  5.3180us  5.3180us  cuEventSyncronize
0.00%  2.9330us  1  1.9330us  1.9330us  1.9330us  cuEventSynchronize
0.00%  3.9330us  3  1.3110us  793ns  1.6970us  cuDeviceGet
0.00%  3.6300us  4  907ns  392ns  1.8090us  cuDeviceComputeCapability
```

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.