# REPORT.4.CUDA

Nguyen Ngoc Son 009, M2-ICT USTH                                         24/12/2022

## Summary

In this lab work, 1D and 2D GPU parallelism are compared in term of GPU processing time of grays cale. The framework used is numba in python under Google Colab GPU session (Tesla GPU).

Figure 1: Original RGB photo

As you can see in the figure 1, the image is RGB originally. After implementation of gray scale code with GPU, the resulting image is as below figure 2 is the same example.

Listing 1: Sample Python code – gray scale RGB image 1D and 2D implementations.

```
1
2  # 1D implementation
3  @cuda.jit
4  def grayscale(src, dst):
5    tidx = cuda.threadIdx.x + cuda.blockIdx.x * cuda.blockDim.x
6    g = np.uint8((src[tidx, 0] + src[tidx, 1] + src[tidx, 2]) / 3)
7    dst[tidx, 0] = dst[tidx, 1] = dst[tidx, 2] = g
```
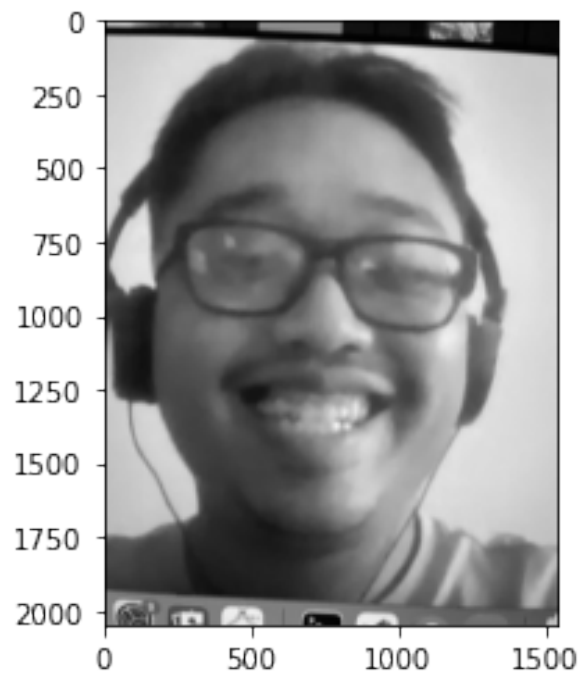
Figure 2: Gray scale photo

```
8
9  # 2D implementation
10 @cuda.jit
11 def grayscale(src, dst):
12   tidx = cuda.threadIdx.x + cuda.blockIdx.x * cuda.blockDim.x
13   tidy = cuda.threadIdx.y + cuda.blockIdx.y * cuda.blockDim.y
14   g = np.uint8((src[tidy, tidx, 0] +src[tidy, tidx, 1] + src[tidy, ↩
       tidx, 2])/3)
15   dst[tidy, tidx] = g
```

## Implementations

Listing 2: 1D block size calculation

```
1 pixelCount = imageWidth * imageHeight
2 blockSize = 64
3 gridSize = int(pixelCount / blockSize)
```

Listing 3: 2D block size calculation

```
1
2 blockSize_list = [[1,1], [1,1], [0.5, 2], [0.5,0.5], [0.25,0.25], ↩
       [0.125,0.125]]
```

```
3  blockSize_list_label = [str([1*32,1*32]), str([1*32,1*32]), str↩
       ([0.5*32, 2*32]), str([0.5*32,0.5*32]), str([0.25*32,0.25*32]), str↩
       ([0.125*32,0.125*32])]
4
5  for idx, warp_ratio in enumerate(blockSize_list):
6    warp_no_w = warp_ratio[0]
7    warp_no_h = warp_ratio[1]
8    ratio_w = int(imageWidth/(32*warp_no_w))
9    ratio_h = int(imageHeight/(32*warp_no_h))
10
11   gridSize = (ratio_w, ratio_h)
12   blockSize = (int(32*warp_no_w), int(32*warp_no_h)) #max thread/block↩
         = 1024
```

**warp multiplication block size sequence:**

blockSize list = [[1,1], [1,1], [0.5, 2], [0.5,0.5], [0.25,0.25], [0.125,0.125]]

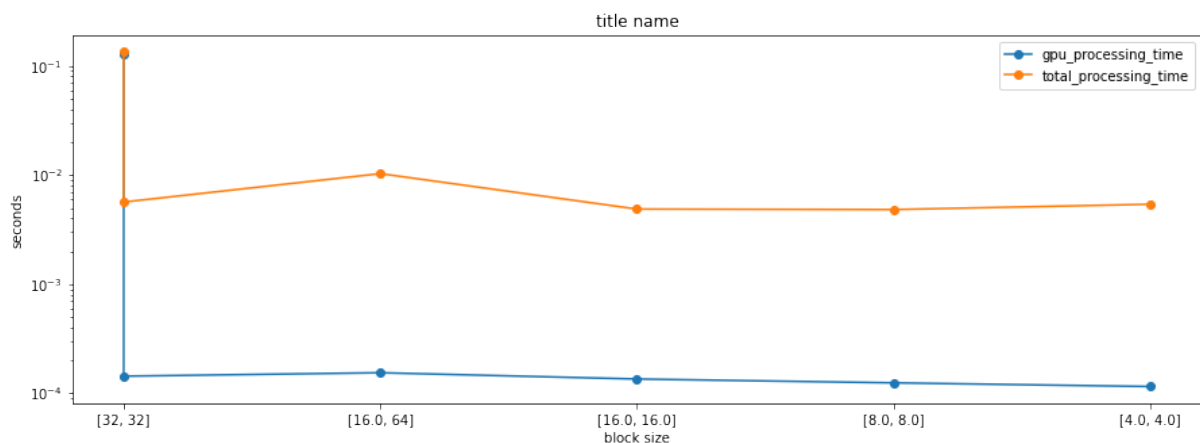*Multiply these ratio with 32 (max threads/ warp) to produce the block size.



Figure 3: Result of different block size over execution time of 2D implementation

**Output:**

blockSize list = [[1,1], [1,1], [0.5, 2], [0.5,0.5], [0.25,0.25], [0.125,0.125]]

GPU Processing time:

"[0.128113, 0.00014, 0.00015, 0.00013, 0.00012, 0.00011]"

**This is the implementation on multiple different block size with 2D array. From here we can tell block size of (4 cores, 4 cores) is fastest block size among these.**

Comparing with the previous lab work running only 1D implementation, the GPU processing is less than about 1 thousands times, figure 4 !
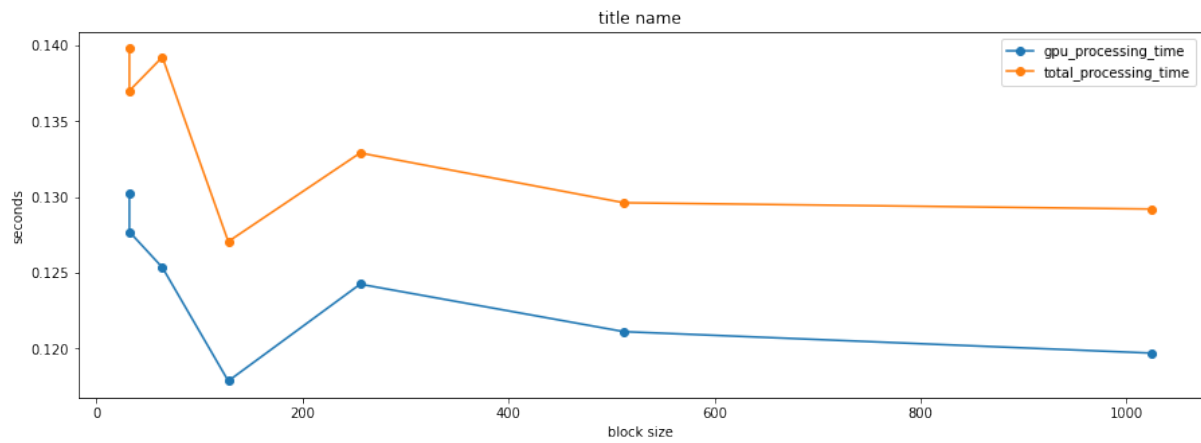
## Reference: 2D implementation code

Figure 4: Result of different block size over execution time of 1D implementation

Listing 4: Snip implemented on GPU

```
1
2  @cuda.jit
3  def grayscale(src, dst):
4    tidx = cuda.threadIdx.x + cuda.blockIdx.x * cuda.blockDim.x
5    tidy = cuda.threadIdx.y + cuda.blockIdx.y * cuda.blockDim.y
6    g = np.uint8((src[tidy, tidx, 0] +src[tidy, tidx, 1] + src[tidy, ↩
         tidx, 2])/3)
7    dst[tidy, tidx] = g
8
9
10 total_processing_time =[]
11 gpu_processing_time = []
12
13
14 blockSize_list = [[1,1], [1,1], [0.5, 2], [0.5,0.5], [0.25,0.25], ↩
       [0.125,0.125]]
15 blockSize_list_label = [str([1*32,1*32]), str([1*32,1*32]), str↩
       ([0.5*32, 2*32]), str([0.5*32,0.5*32]), str([0.25*32,0.25*32]), str↩
       ([0.125*32,0.125*32])]
16
17 for idx, warp_ratio in enumerate(blockSize_list):
18   warp_no_w = warp_ratio[0]
19   warp_no_h = warp_ratio[1]
20   ratio_w = int(imageWidth/(32*warp_no_w))
21   ratio_h = int(imageHeight/(32*warp_no_h))
22
23   gridSize = (ratio_w, ratio_h)
24   blockSize = (int(32*warp_no_w), int(32*warp_no_h)) #max thread/block↩
         = 1024
25
26   image = image.copy()
27   start = time.time()
```

```
28    devData = cuda.to_device(image)
29
30    devOutput = cuda.device_array(
31                              (imageHeight, imageWidth, 3),
32                              np.uint8)
33
34    gpu_start = time.time()
35    grayscale[gridSize, blockSize](devData, devOutput)
36    gpu_end = time.time()
37
38    hostOutput = devOutput.copy_to_host()
39    # hostOutput = hostOutput.reshape((imageHeight, imageWidth,3)).↩
          astype('uint8')
40    end = time.time()
41
42    # if idx > 0:
43    total_processing_time.append(end - start)
44    gpu_processing_time.append(gpu_end - gpu_start)
```