

Project5 Write Up

Kernal Implementation

In the project we were given a CPU example of the convolve kernel and were asked to implement GPU version. Writing the GPU kernel was a minor tweak from the CPU version. The two most outer for loop does point-wise loop over the image pixels. I just had to update iterator variable by adding $x0 = (\text{blockIdx.x} * \text{blockDim.x} + \text{threadIdx.x})$ to outer loop that loops thorough x-axis and $y0 = (\text{blockIdx.y} * \text{blockDim.y} + \text{threadIdx.y})$ to inner loop that loops thorough y-axis. $x0$ and $y0$ gives start of each block thus each block will loop starting from its location.

Striding was implemented by updating the iteration variable of outer for loop by $\text{blockDim.x} * \text{gridDim.x}$ and updating the iteration variable of inner for loop by adding $\text{blockDim.y} * \text{gridDim.y}$ which is the total number of thread in the grid. In the case where input cell number is larger than thread available, after $\text{blockDim.x} * \text{gridDim.x}$, $\text{blockDim.y} * \text{gridDim.y}$, there will still be input- $\text{blockDim.x} * \text{gridDim.x}$, and input- $\text{blockDim.y} * \text{gridDim.y}$ cell to be processed so by striding by $\text{blockDim.x} * \text{gridDim.x}$, loop would not terminate and process the left over cells after iteration untill all have been computed

Data distribution

With block size and grid sizes, we distribute data to thread in each block by starting loop at $\text{blockIdx.x} * \text{blockDim.x} + \text{threadIdx.x}$. and $\text{blockIdx.y} * \text{blockDim.y} + \text{threadIdx.y}$ in the inner loop.

Since each block have different blockIdx and each thread different threadIdx , each thread in each block will start processing cells at different location.

Performance result

Execution times to run the GPU version on the input file video.mp4 with block sizes 2,4,8,16,32

Grid size was adjusted accordingly:

$\text{gridSizeX} = \text{height} / \text{blockDimSize} + (1 \text{ if } X_Limit \% \text{blockDimSize} \neq 0, \text{ else } 0)$

$\text{gridSizeY} = \text{width} / \text{blockDimSize} + (1 \text{ if } Y_Limit \% \text{blockDimSize} \neq 0, \text{ else } 0)$

where width and height give the dimensions of the video frame.

Video used was 1920 x 1080 width x height

Grid Sizes:

Block dim size 2 => 540 x 960

Block dim size 4 => 270 x 480

Block dim size 8 => 135 x 240

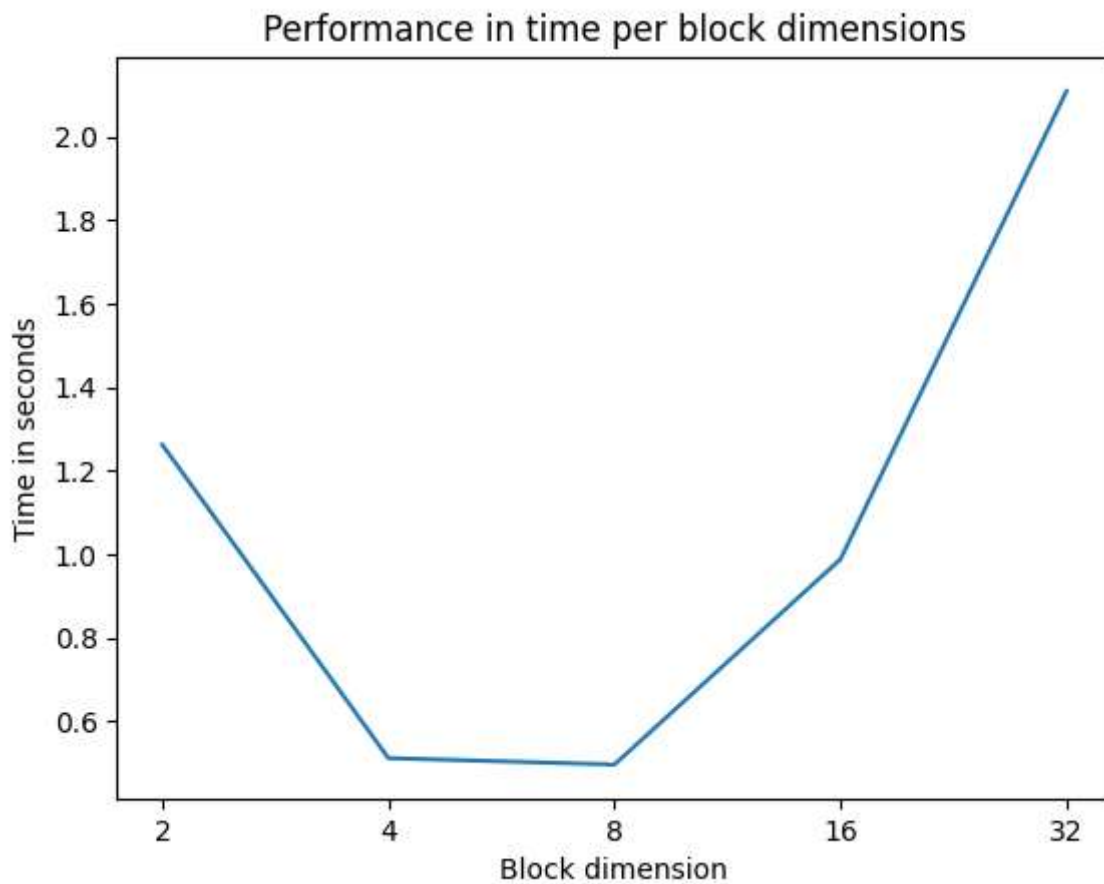
Block dim size 16 => 68 x 120

Block dim size 32 => 34 x 60

```
In [ ]: import numpy as np
import matplotlib.pyplot as plt

x = np.array(['2', '4', '8', '16', '32'])
y = np.array([1.26197, 0.511747, 0.496408, 0.987956, 2.10882])

plt.plot(x,y)
plt.title("Performance in time per block dimensions")
plt.xlabel("Block dimension")
plt.ylabel("Time in seconds")
plt.show()
```



Performance increased until block dim size of 8 then gets worse. I did expect that at some point when increasing block dim, the performance would get worse due to each SM fitting less threadblock resulting in actually less total threads per SM but I did not expect it to happen so soon.