

Image Histogram along A Direction with CUDA

CS677 Final Project Report Ying Liu

1. Problem description

When we have a 2D image of a measurement result, for example, pressure or temperature distribution on a surface or a CT image for medical purpose, other than the histogram of the whole image, sometimes we would like to see how the data is distributed horizontally or vertically, or even along a certain line we drew on the picture. Thus, getting histogram along a direction is useful in many applications. For histogram along a line, the line can be decided by the coordinates of two points (Fig 1).

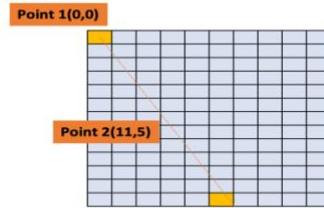


Fig. 1

To produce luminance histogram of an image along this line, we can rotate the image first with an angle, then calculate the histogram column wisely. The output result will be a 2D matrix with $\text{bin_Num} * \text{WIDTH}$ elements. This WIDTH will be the width after rotation. This 2D matrix contains the position information along the given direction. The pseudo code is as follows:

for pixel along the col:

Histogram[luminance(pixel)] [position] ++

After some calculation, we can get the histogram at this line after rotation, which will be just a column of the final result.

2. Suitability for GPU acceleration

- 1) Image rotation is a natural fit for GPU acceleration. Each thread rotates one pixel, there is no data dependency among threads.
- 2) Image histogram can also be paralleled processed by GPU either by each thread deals with one pixel or each thread gathers data for one bin.
- 3) Amdahl's Law: close to 100%: After image padding, no extra processing on CPU.

3. Content

There are 2 steps for the GPU code implementation: rotation and histogram.

1) Image Rotation

The new coordinates (x2, y2) of (x1, y1) when rotated by an angle θ around origin point (x0, y0) is¹:

$$x2 = \cos(\theta) * (x1 - x0) - \sin(\theta) * (y1 - y0) + x0$$

$$y2 = \sin(\theta) * (x1 - x0) + \cos(\theta) * (y1 - y0) + y0$$

For example, if the given points are (0, 0) and (1, 1), the rotate θ is 45 degree. I choose to rotate along the center of the image as given angle, like a swirl. In order to keep the original image non-cropped, I padded the image to the size of its diagonal size then rotate it. In this case, no matter rotating the image to any angle, it will keep all the original image safe and sound. Paddling part is done in CPU because there is not much GPU optimization can be done on it. Here is the code for paddling the image:

```
int deltaX = diaLen - xsize ;
int deltaY = diaLen - ysize ;
for(int i=0; i<diaLen; i++){
    for(int j = 0; j< diaLen; j++){
        if(i>=deltaY/2 && i< (ysize + deltaY/2) && j >= deltaX/2 && j < (xsize + deltaX/2)){
```

```

        h_Input[i*diaLen+j]=pic[(i-deltaY/2)*xsize + (j-deltaX/2)];
    }else{
        h_Input[i*diaLen+j] = 0;
    }
}
}
}

```

Three diagrams in Fig. 2 show the steps to rotate an image 45° clockwise.

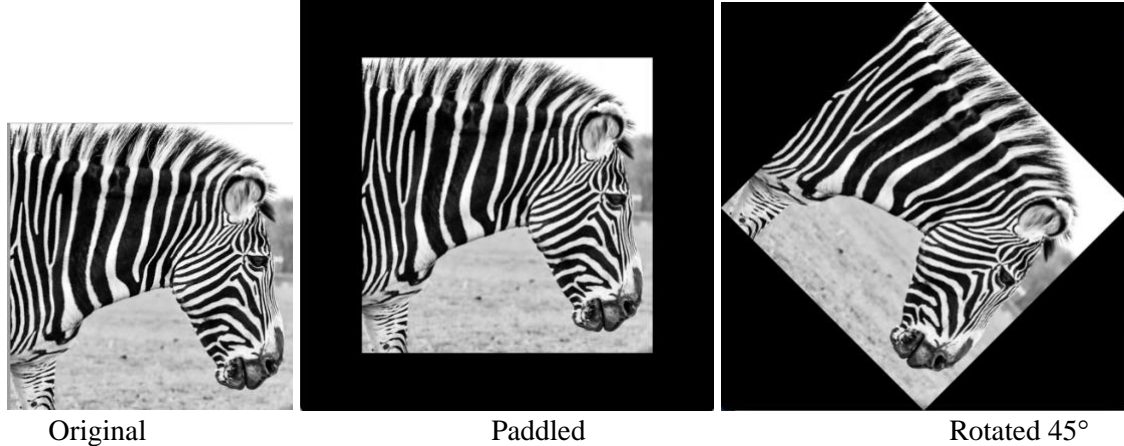


Fig. 2 (Original Image Source: Unsplash)

Since there is no data sharing or write conflicts between each pixel, in the kernel, firstly I launches 1024 threads per block then I found there is a corner cropped in the image. By using “-Xptxas=“-v”” during the compiling, I found that the maximum register usage is the “rotation_naive_kernel”, which has 44 registers. Thus, I used less threads per block for the rotation kernel and got the right output. Here are the block dimension and grid dimension of rotation kernel. The “diaLen” as mentioned before, is the width/height of the padded image:

```

int blockSize = 8;
dim3 blockDim(blockSize,blockSize,1);
int gridSize = (diaLen + blockSize -1)/blockSize;
dim3 gridDim(gridSize,gridSize,1);

```

I implemented 3 versions of rotation kernel:

- “rotation_kernel_naive”: This kernel directly loads and stores data to the global memory.
- “rotation_kernel_2”: This kernel used more registers to save several precalculated values such as: xCenter, yCenter, sin(angle), cos(angle). These four values are repeatedly used by every thread.

```

347 //-----rotation kernel 2-----
348 _global_ void rotation_kernel_2(unsigned int *input,unsigned
349 //TO DO
350 double xCenter = (double)width / 2;
351 double yCenter = (double)height / 2;
352 int x = blockIdx.x * blockDim.x + threadIdx.x;
353 int y = blockIdx.y * blockDim.y + threadIdx.y;
354 int index = y * width + x;
355 double sinA = sin(angle);
356 double cosA = cos(angle);
357 double shiftX = (double)x - xCenter;
358 double shiftY = (double)y - yCenter;
359 int orgX = 0;
360 int orgY = 0;
361 //boundary check
362 if(x >=0 && x < width && y >=0 && y < height){
363     orgX = (int)(cosA * shiftX - sinA * shiftY + xCenter);
364     orgY = (int)(sinA * shiftX + cosA * shiftY + yCenter);
365 }
366
367 if(orgX>=0 && orgX < width && orgY>=0 && orgY<height){
368     output[index] = input[ orgY * width + orgX];
369 }
370
371 }
372

```

Fig. 3 Snapshot of “rotation_kernel_2”

- C. “rotation_kernel_3”: This kernel used constant memory to save those four values mentioned in “rotation_kernel_2”.

```
double *P;
P = (double*) malloc(4 * sizeof(double));
P[0] = (double)diaLen / 2; //xCenter
P[1] = (double)diaLen / 2; //yCenter
P[2] = sin(angle);
P[3] = cos(angle);

//load data to constant memory
cudaMemcpyToSymbol(PARAMS, P, 4 * sizeof(double));
//timer
```

Fig. 4 Snapshot of Constant Memory Kernel Code

Comparison between CPU rotation and GPU naïve kernel on different size of Images is showed in Table 1. From the time consumption comparison, it is obvious that the GPU rotation speed up a lot more than the CPU. For the XLarge image, the GPU is almost 1000x faster than the CPU.

Time Unit: milliseconds	Small size (640x651 pixels)	Middle size (1920x1953 pixels)	Large Size (2400x2442 pixels)	Xlarge Size (3932 x 4000 pixels)
CPU-Rotation	59	565	927	2845
GPU- Rotation_Naive	0.067	0.572	0.892	2.355

Table 1

The GPU optimization comparison by different versions of kernels based on different image size is showed in Table 2 and Fig. 5. We can see that the best performance among these three kernels is using constant memory.

Time unit: milliseconds	Naïve	Registers	Constant Memory
Small size (640x651 pixels)	0.067	0.03804	0.02076
Middle size (1920x1953 pixels)	0.572	0.2671	0.1395
Large Size (2400x2442 pixels)	0.892	0.4211	0.2398
Xlarge Size (3932 x 4000 pixels)	2.355	1.2034	0.5949

Table 2

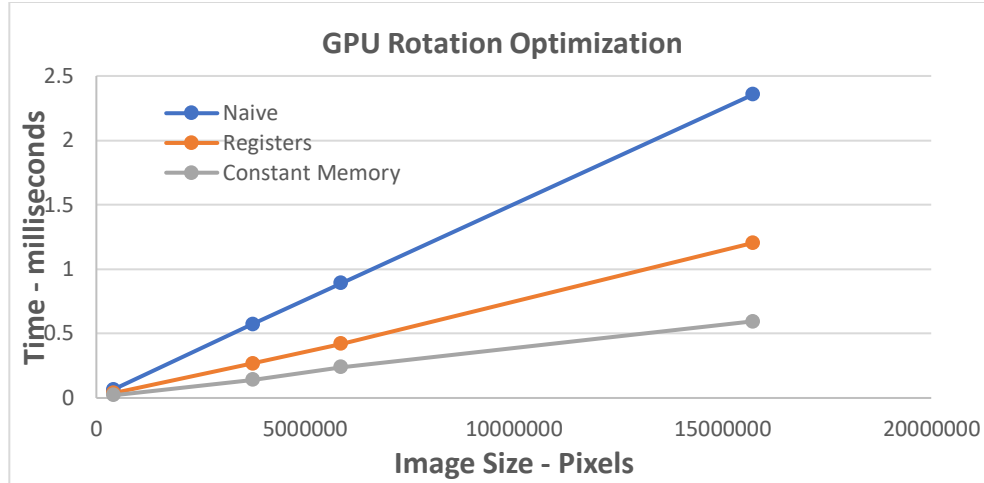


Fig. 5

2) Image Histogram

After padding and rotation, we get a square image with width and height both equal to the original image's diagonal length, noted as "diaLen". Now by calculating the histogram of every column, we get a 2D matrix with height equals the number of bins and width equals to "diaLen". I implemented 2 versions of histogram kernel: naïve one and the one using shared memory.

A. "histo_kernel_naive": Here are the block dimension and grid dimension of the naïve histogram kernel:

```
int blockSize = 32;
dim3 blockDim(blockSize,blockSize,1);
int gridSize = (diaLen + blockSize -1)/blockSize;
dim3 gridDim(gridSize,gridSize,1);
```

The naïve version directly loads and store using global memory. It uses "atomicAdd" method to deal with race condition. The row of the "histo" will be the gray scale value of the input and the column will be same as the image column. Fig. 6 is a snapshot of kernel:

```
__global__ void histo_kernel_naive(unsigned int *input,unsigned int *histo,int width)
{
    int x = blockDim.x * blockIdx.x + threadIdx.x;
    int y = blockDim.y * blockIdx.y + threadIdx.y;

    if(x >=0 && x < width && y >=0 && y < height){
        atomicAdd(&histo[input[y*width+x]*width+x],1);
    }
}
```

Fig. 6 Snapshot of histo_kernel_naive

B. "histo_kernel_2": In the naïve kernel, the "histo" global parameter has be accessed multiple times by each thread to do the accumulation, thus using shared memory can decrease the times of read and write global memory, which is a good choice to decrease the time consuming. So, in "histo_kernel_2", shared memory H[NUM_BINS] is used for storing one column of the final result. When the accumulation is done, each thread writes to global memory once. For this kernel the threads number launched equals to the histogram result element number:

```
histo_kernel_2<<<diaLen,256>>>(d_rotated,d_histo2,diaLen,diaLen);
```

Here is a snapshot of this kernel:

```

__global__ void histo_kernel_2(unsigned int *input,unsigned int *hist
__shared__ unsigned int H[NUM_BINS]; // one column of the output bin
int tx = threadIdx.x;
int col = blockIdx.x; // col number of input
int bd = blockDim.x;
H[tx] = 0;
__syncthreads();
for(int t =0;t<(height+bd-1)/bd;t++){ //when height is bigger than
    int row = t*bd+tx;
    if(row<height){
        int value = input[row*width+col];
        atomicAdd(&H[value],1);
    }
}
__syncthreads();
histo[tx*width+col] = H[tx];
}

```

Fig. 7 Snapshot of histo_kernel_2

Since threads number in one block could be smaller than the height of the image, there is a loop in the kernel to read all the sessions of one column of the image and does the “atomicAdd”. When it is done, each block writes one column of the final result.

Comparison between CPU histogram and GPU naïve kernel is showed in Table 3. We can see that for the XLarge size image, even using GPU naïve version can get almost 100x faster speed.

(milliseconds)	Small (640x651)	Middle (1920x1953)	Large (2400x2442)	Xlarge (3932 x 4000)
CPU-Histogram	6.8	39.96	60.8	156
GPU-Histogram_Naive	0.093	0.409	0.59	1.402

Table 3

GPU optimization comparison on different version of kernels based on different size of Images is showed in Table 4 and Fig. 8.

Time unit: milliseconds	Naïve	shared memory
Small size (640x651 pixels)	0.093	0.048
Middle size (1920x1953 pixels)	0.409	0.362
Large Size (2400x2442 pixels)	0.59	0.532
Xlarge Size (3932 x 4000 pixels)	1.402	1.356

Table 4

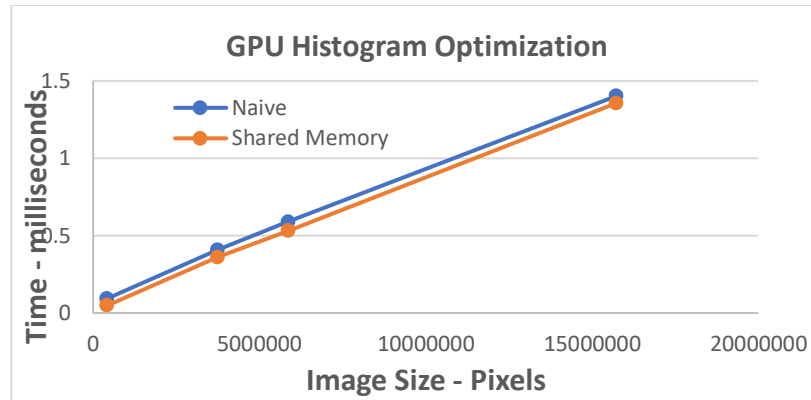


Fig. 8

4. Instructions on running the project

All the kernels and main function are in one file. When running the project without arguments:

```
!./histo
```

The main function will do default 45 degree rotation clockwise and the default image “small-zebra-unsplash.ppm” will be used. If one argument – the rotation angle be supplied, it will rotate the default image using the given angle, for example, rotating the image to 30 degree.

```
!./histo 30
```

For 2 args, the first one will be angle and the second one will be the image file. For example, the below command line will rotate the image “mid-zebra-unsplash.ppm” to 30 degree and then do the histogram. The histogram result will be in a “.csv” file.

```
!./histo 30 mid-zebra-unsplash.ppm
```

5. Track the line

Remember the arbitrary line we’d like to know the histogram along it? Now we have the histogram result, by using the method “getNewXY” in the project we can have the answer. For example, we have the start point (0,0) and end point (4,4) and the angle we get is 45 degree clockwise. After rotation we use this “getNewXY” function we get the (newX, newY) is (461,1). The 462nd column of the result is the histogram along this line. For other coordinates, just modify the 96th and 97th line of code to modify the coordinates. Fig. 9 shows the 45° line start from (0,0) after rotation. Fig. 10 shows the 462nd column of the histogram result in the “.csv” file, which is just the histogram of this line.

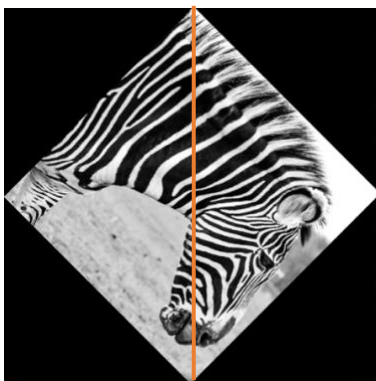


Fig. 9

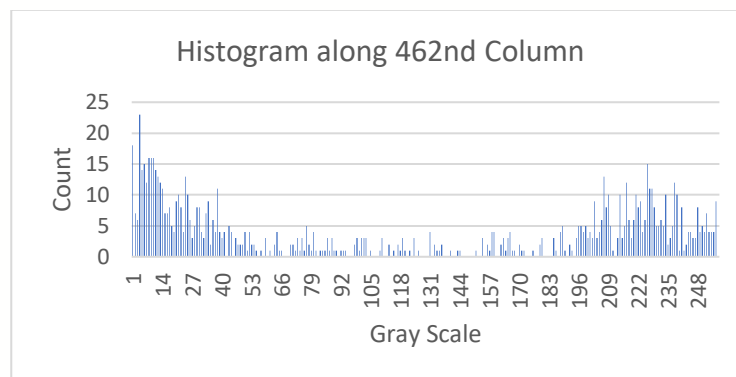


Fig. 10

6. Conclusion

This project produces histogram on an image along a direction by firstly rotating the image to desired angle and then do histogram column wisely on the rotated image. By applying parallel programming on CUDA, the program gained very satisfying speed up. Optimization methods used during the GPU implementation include using registers, using constant memory and using shared memory.

Code and results of this project can be seen at:

<https://github.com/yingliu928/CS677FinalProject.git>

References:

1. <http://cuda-programming.blogspot.com/2013/03/optimization-in-histogram-cuda-code.html>