

CUDA - K-means

Anand Kamble
Department of Scientific Computing
Florida State University

1 Introduction

In this homework, we will be implementing the k-means algorithm for image segmentation using CUDA.

2 Implementation

2.1 Host-side Initialization

The host program reads the input image file (in this case, a JPEG file) using the provided JpegFile library, obtaining the image width, height, and pixel data in RGB format.

Listing 1: main.cu

```
135  UINT width, height;  
136  uint8_t *hostDataBuf = JpegFile::JpegFileToRGB("test-large.jpg", &width,  
        &height);
```

2.2 Device Memory Allocation

The host program allocates device memory for the input RGB data and the output segmented data using `cudaMalloc`.

Listing 2: main.cu

```
148  cudaMalloc(&device_colors, width * height * sizeof(Pixel));
```

2.3 Data Transfer to Device

The input RGB data is copied from the host memory to the device memory using `cudaMemcpy`.

Listing 3: main.cu

```
148  cudaMemcpy(device_colors, hostDataBuf, width * height * sizeof(Pixel),  
        cudaMemcpyHostToDevice);
```

2.4 Kernel Launch

The k-means algorithm is implemented using three CUDA kernels.

2.4.1 Grouping Kernel

This kernel is responsible for grouping each pixel to the nearest generator (cluster center) based on the Euclidean distance

Listing 4: main.cu

```
42  __global__ void groupingKernel(UINT width, UINT height, int k, Pixel *colors,  
        Pixel *generators, int *groupColorSum, int *groupCount){  
43      int x = threadIdx.x + blockIdx.x * blockDim.x;  
44      int y = threadIdx.y + blockIdx.y * blockDim.y;  
45  }
```

```

46     if (x < width && y < height){
47         int index = y * width + x;
48         Pixel pixel = colors[index];
49         float minDist = INFINITY;
50         int minIndex = 0;
51
52         for (int i = 0; i < k; i++){
53             Pixel gen = generators[i];
54             float dist = sqrtf((pixel.r - gen.r) * (pixel.r - gen.r) +
55                               (pixel.g - gen.g) * (pixel.g - gen.g) + (pixel.b - gen.b) *
56                               (pixel.b - gen.b));
57             if (dist < minDist){
58                 minDist = dist;
59                 minIndex = i;
60             }
61         }
62         atomicAdd(&groupColorSum[minIndex * 3 + 0], (int)pixel.r);
63         atomicAdd(&groupColorSum[minIndex * 3 + 1], (int)pixel.g);
64         atomicAdd(&groupColorSum[minIndex * 3 + 2], (int)pixel.b);
65         atomicAdd(&groupCount[minIndex], 1);
66     }
67 }

```

2.4.2 Update Generators Kernel

This kernel is used to update the generators (cluster centers) based on the group color sums and counts.

Listing 5: main.cu

```

42 __global__ void updateGeneratorsKernel(int k, Pixel *generators, int
43   *groupColorSum, int *groupCount)
44 {
45     int i = threadIdx.x;
46     if (i < k)
47     {
48         generators[i].r = (int)groupColorSum[i * 3 + 0] / groupCount[i];
49         generators[i].g = (int)groupColorSum[i * 3 + 1] / groupCount[i];
50         generators[i].b = (int)groupColorSum[i * 3 + 2] / groupCount[i];
51     }
52 }

```

2.4.3 Replace Color Kernel

This kernel is used to replace the color of each pixel with the color of the group generator it belongs to, effectively segmenting the image.

Listing 6: main.cu

```

42 __global__ void replaceColorKernel(UINT width, UINT height, int k, Pixel
43   *colors, Pixel *generators, int *groupCount)
44 {
45     int x = threadIdx.x + blockIdx.x * blockDim.x;
46     int y = threadIdx.y + blockIdx.y * blockDim.y;
47
48     if (x < width && y < height)
49     {
50         int index = y * width + x;
51         Pixel pixel = colors[index];
52
53         float minDist = INFINITY;
54         int minIndex = 0;
55
56         for (int i = 0; i < k; i++)

```

```

56     {
57         Pixel gen = generators[i];
58         float dist = sqrtf((pixel.r - gen.r) * (pixel.r - gen.r) +
59                             (pixel.g - gen.g) * (pixel.g - gen.g) + (pixel.b - gen.b) *
60                             (pixel.b - gen.b));
61         if (dist < minDist)
62         {
63             minDist = dist;
64             minIndex = i;
65         }
66     }
67     colors[index] = generators[minIndex];
68 }

```

2.5 Data Transfer to Host

After the kernel executions, the resulting segmented image data is copied from the device memory back to the host memory using `cudaMemcpy`.

Listing 7: main.cu

```

198 cudaMemcpy(hostDataBuf, device_colors, N * sizeof(Pixel),
199             cudaMemcpyDeviceToHost);

```

2.6 Output File Writing

The host program writes the segmented image data to a new JPEG file using the `JpegFile` library.

Listing 8: main.cu

```

198 JpegFile::RGBToJpegFile("output.jpg", hostDataBuf, width, height, 100, false);

```

2.7 Memory Deallocation

The host program frees the allocated device and host memory.

Listing 9: main.cu

```

211 delete[] generators;
212 free(hostDataBuf);
213 cudaFree(device_colors);
214 cudaFree(device_generators);
215 cudaFree(device_groupColorSum);
216 cudaFree(device_groupCount);

```

2.8 Benchmarking

We are timing the program using the `cudaEventRecord` function provided by CUDA.

Listing 10: main.cu

```

211 cudaEvent_t start, stop;
212 cudaEventCreate(&start);
213 cudaEventCreate(&stop);
214 cudaEventRecord(start); // Start recording
215 // ... processing part
216 cudaEventRecord(stop);
217 cudaEventSynchronize(stop);
218 float elapsedTime; // Variable to store the elapsed time
219 cudaEventElapsedTime(&elapsedTime, start, stop);

```

3 Results

After successful execution, we get the following results:

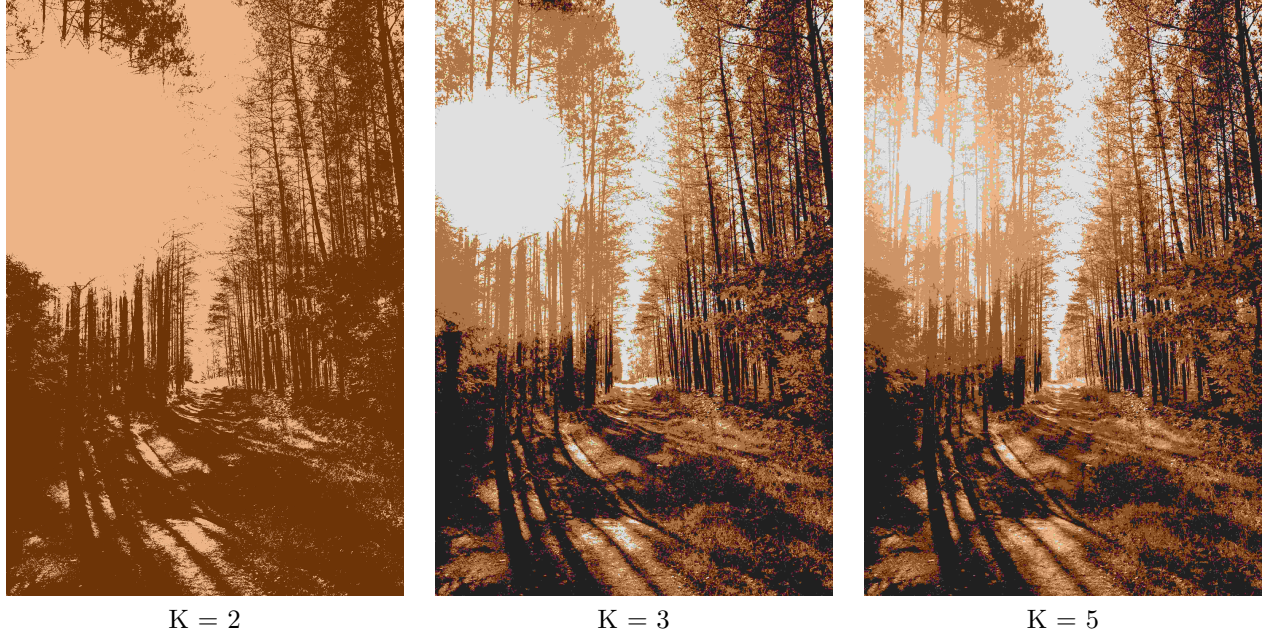


Table 1: K-means

The following chart shows the time taken by CPU compared to the GPU.

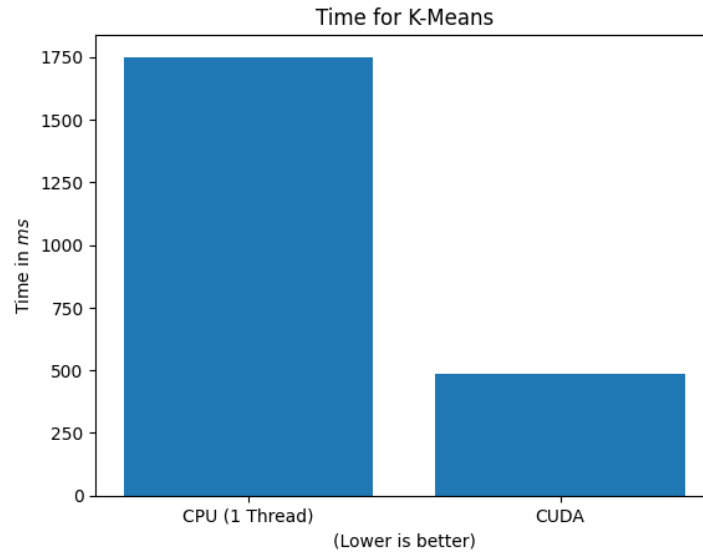


Figure 1: Caption

References

- [1] NVIDIA Corporation. *K-Means Clustering Algorithm* , <https://www.nvidia.com/en-us/glossary/k-means/>. Oct 11, 2017.
- [2] Renan. *Difference between cuda.h, cuda_runtime.h, cuda_runtime_api.h* , <https://stackoverflow.com/questions/6302695/difference-between-cuda-h-cuda-runtime-h-cuda-runtime-api-h>. Jun 10, 2011.