

Module 3 Lab

CUDA Image Color to Grayscale

GPU Teaching Kit – Accelerated Computing

OBJECTIVE

The purpose of this lab is to convert an RGB image into a gray scale image. The input is an RGB triple of float values and the student will convert that triple to a single float grayscale intensity value. A pseudo-code version of the algorithm is shown below:

```
for ii from 0 to height do
  for jj from 0 to width do
    idx = ii * width + jj
    # here channels is 3
    r = input[3*idx]
    g = input[3*idx + 1]
    b = input[3*idx + 2]
    grayImage[idx] = (0.21*r + 0.71*g + 0.07*b)
  end
end
```

PREREQUISITES

Before starting this lab, make sure that:

- You have completed the required module videos

IMAGE FORMAT

For people who are developing on their own system, the input image is stored in PPM P6 format while the output grayscale image is stored in PPM P5 format. Students can create their own input images by exporting their image into PPM images. The easiest way to create image is via external tools. On Unix, `bmptoppm` converts BMP images to PPM images.

INSTRUCTIONS

Edit the code in the code tab to perform the following:

- allocate device memory
- copy host memory to device
- initialize thread block and kernel grid dimensions
- invoke CUDA kernel
- copy results from device to host
- deallocate device memory

Instructions about where to place each part of the code is demarcated by the `//@@` comment lines.

LOCAL SETUP INSTRUCTIONS

The most recent version of source code for this lab along with the build-scripts can be found on the [Bitbucket repository](#). A description on how to use the [CMake](#) tool in along with how to build the labs for local development found in the [README](#) document in the root of the repository.

The executable generated as a result of compiling the lab can be run using the following command:

```
./ImageColorToGrayscale.Template -e <expected.pbm> \
-i <input.ppm> -o <output.pbm> -t image`.
```

where `<expected.pbm>` is the expected output, `<input.ppm>` is the input dataset, and `<output.pbm>` is an optional path to store the results. The datasets can be generated using the dataset generator built as part of the compilation process.

QUESTIONS

- (1) How many floating operations are being performed in your color conversion kernel? EXPLAIN.

ANSWER: **There are 3 floating point multiplications and 2 additions**

- (2) Which format would be more efficient for color conversion: a 2D matrix where each entry is an RGB value or a 3D matrix where each slice in the Z axis represents a color. I.e. is it better to have color interleaved in this application? can you name an application where the opposite is true?

ANSWER: **In this case an interleaved representation is better since the data accesses are coalesced.**

- (3) How many global memory reads are being performed by your kernel? EXPLAIN.

ANSWER: **There is one global memory read per pixel (i.e. three values are being read).**

- (4) How many global memory writes are being performed by your kernel?
EXPLAIN.

ANSWER: **The is one global memory write per pixel. We write one floating point value, since the output is a summary of the RGB values.**

- (5) Describe what possible optimizations can be implemented to your kernel to achieve a performance speedup.

ANSWER: **One can operate in texture memory.**

- (6) Name three applications for color conversion.

ANSWER: **Color conversion is a requirement to many algorithms. Some examples are image histogram, edge detection, and image effects.**

CODE TEMPLATE

The following code is suggested as a starting point for students. The code handles the import and export as well as the checking of the solution. Students are expected to insert their code in the sections demarcated with `//@@`. Students expected the other code unchanged. The tutorial page describes the functionality of the `wb*` methods.

```

1  #include <wb.h>
2
3  #define wbCheck(stmt)                                     \
4      do {                                                 \
5          cudaError_t err = stmt;                          \
6          if (err != cudaSuccess) {                        \
7              wbLog(ERROR, "Failed to run stmt ", #stmt);  \
8              wbLog(ERROR, "Got CUDA error ... ", cudaGetErrorString(err)); \
9              return -1;                                    \
10         }                                                 \
11     } while (0)
12
13  //@@ INSERT CODE HERE
14
15  int main(int argc, char *argv[]) {
16      wbArg_t args;
17      int imageChannels;
18      int imageWidth;
19      int imageHeight;
20      char *inputImageFile;
21      wbImage_t inputImage;
22      wbImage_t outputImage;
23      float *hostInputImageData;
24      float *hostOutputImageData;
25      float *deviceInputImageData;
26      float *deviceOutputImageData;
27
28      args = wbArg_read(argc, argv); /* parse the input arguments */
29
30      inputImageFile = wbArg_getInputFile(args, 0);

```

```

31
32 inputImage = wbImport(inputImageFile);
33
34 imageWidth = wbImage_getWidth(inputImage);
35 imageHeight = wbImage_getHeight(inputImage);
36 // For this lab the value is always 3
37 imageChannels = wbImage_getChannels(inputImage);
38
39 // Since the image is monochromatic, it only contains one channel
40 outputImage = wbImage_new(imageWidth, imageHeight, 1);
41
42 hostInputImageData = wbImage_getData(inputImage);
43 hostOutputImageData = wbImage_getData(outputImage);
44
45 wbTime_start(GPU, "Doing GPU Computation (memory + compute)");
46
47 wbTime_start(GPU, "Doing GPU memory allocation");
48 cudaMalloc((void **)&deviceInputImageData,
49             imageWidth * imageHeight * imageChannels * sizeof(float));
50 cudaMalloc((void **)&deviceOutputImageData,
51             imageWidth * imageHeight * sizeof(float));
52 wbTime_stop(GPU, "Doing GPU memory allocation");
53
54 wbTime_start(Copy, "Copying data to the GPU");
55 cudaMemcpy(deviceInputImageData, hostInputImageData,
56             imageWidth * imageHeight * imageChannels * sizeof(float),
57             cudaMemcpyHostToDevice);
58 wbTime_stop(Copy, "Copying data to the GPU");
59
60 //////////////////////////////////////
61 wbTime_start(Compute, "Doing the computation on the GPU");
62 //@@ INSERT CODE HERE
63
64 wbTime_stop(Compute, "Doing the computation on the GPU");
65
66 //////////////////////////////////////
67 wbTime_start(Copy, "Copying data from the GPU");
68 cudaMemcpy(hostOutputImageData, deviceOutputImageData,
69             imageWidth * imageHeight * sizeof(float),
70             cudaMemcpyDeviceToHost);
71 wbTime_stop(Copy, "Copying data from the GPU");
72
73 wbTime_stop(GPU, "Doing GPU Computation (memory + compute)");
74
75 wbSolution(args, outputImage);
76
77 cudaFree(deviceInputImageData);
78 cudaFree(deviceOutputImageData);
79
80 wbImage_delete(outputImage);
81 wbImage_delete(inputImage);
82

```

```

83     return 0;
84 }

```

CODE SOLUTION

The following is a possible implementation of the lab. This solution is intended for use only by the teaching staff and should not be distributed to students.

```

1  #include <wb.h>
2
3  #define wbCheck(stmt) \
4      do { \
5          cudaError_t err = stmt; \
6          if (err != cudaSuccess) { \
7              wbLog(ERROR, "Failed to run stmt ", #stmt); \
8              wbLog(ERROR, "Got CUDA error ... ", cudaGetErrorString(err)); \
9              return -1; \
10         } \
11     } while (0)
12
13  /// INSERT CODE HERE
14
15  #define TILE_WIDTH 16
16  __global__ void rgb2gray(float *grayImage, float *rgbImage, int channels,
17                          int width, int height) {
18      int x = threadIdx.x + blockIdx.x * blockDim.x;
19      int y = threadIdx.y + blockIdx.y * blockDim.y;
20
21      if (x < width && y < height) {
22          // get 1D coordinate for the grayscale image
23          int grayOffset = y * width + x;
24          // one can think of the RGB image having
25          // CHANNEL times columns than the gray scale image
26          int rgbOffset = grayOffset * channels;
27          float r      = rgbImage[rgbOffset]; // red value for pixel
28          float g      = rgbImage[rgbOffset + 1]; // green value for pixel
29          float b      = rgbImage[rgbOffset + 2]; // blue value for pixel
30          // perform the rescaling and store it
31          // We multiply by floating point constants
32          grayImage[grayOffset] = 0.21f * r + 0.71f * g + 0.07f * b;
33      }
34  }
35
36  int main(int argc, char *argv[]) {
37      wbArg_t args;
38      int imageChannels;
39      int imageWidth;
40      int imageHeight;
41      char *inputImageFile;
42      wbImage_t inputImage;
43      wbImage_t outputImage;

```

```

44     float *hostInputImageData;
45     float *hostOutputImageData;
46     float *deviceInputImageData;
47     float *deviceOutputImageData;
48
49     args = wbArg_read(argc, argv); /* parse the input arguments */
50
51     inputImageFile = wbArg_getInputFile(args, 0);
52
53     inputImage = wbImport(inputImageFile);
54
55     imageWidth  = wbImage_getWidth(inputImage);
56     imageHeight = wbImage_getHeight(inputImage);
57     // For this lab the value is always 3
58     imageChannels = wbImage_getChannels(inputImage);
59
60     // Since the image is monochromatic, it only contains one channel
61     outputImage = wbImage_new(imageWidth, imageHeight, 1);
62
63     hostInputImageData = wbImage_getData(inputImage);
64     hostOutputImageData = wbImage_getData(outputImage);
65
66     wbTime_start(GPU, "Doing GPU Computation (memory + compute)");
67
68     wbTime_start(GPU, "Doing GPU memory allocation");
69     cudaMalloc((void **)&deviceInputImageData,
70               imageWidth * imageHeight * imageChannels * sizeof(float));
71     cudaMalloc((void **)&deviceOutputImageData,
72               imageWidth * imageHeight * sizeof(float));
73     wbTime_stop(GPU, "Doing GPU memory allocation");
74
75     wbTime_start(Copy, "Copying data to the GPU");
76     cudaMemcpy(deviceInputImageData, hostInputImageData,
77               imageWidth * imageHeight * imageChannels * sizeof(float),
78               cudaMemcpyHostToDevice);
79     wbTime_stop(Copy, "Copying data to the GPU");
80
81     ////////////////////////////////////
82     wbTime_start(Compute, "Doing the computation on the GPU");
83     ///@@ INSERT CODE HERE
84     dim3 dimGrid(ceil((float)imageWidth / TILE_WIDTH),
85                  ceil((float)imageHeight / TILE_WIDTH));
86     dim3 dimBlock(TILE_WIDTH, TILE_WIDTH, 1);
87     rgb2gray<<<dimGrid, dimBlock>>>(deviceOutputImageData,
88                                     deviceInputImageData, imageChannels,
89                                     imageWidth, imageHeight);
90     wbTime_stop(Compute, "Doing the computation on the GPU");
91
92     ////////////////////////////////////
93     wbTime_start(Copy, "Copying data from the GPU");
94     cudaMemcpy(hostOutputImageData, deviceOutputImageData,
95               imageWidth * imageHeight * sizeof(float),
96               cudaMemcpyDeviceToHost);

```

```
97     wbTime_stop(Copy, "Copying data from the GPU");
98
99     wbTime_stop(GPU, "Doing GPU Computation (memory + compute)");
100
101     wbSolution(args, outputImage);
102
103     cudaFree(deviceInputImageData);
104     cudaFree(deviceOutputImageData);
105
106     wbImage_delete(outputImage);
107     wbImage_delete(inputImage);
108
109     return 0;
110 }
```

© ⓘ ⓘ This work is licensed by UIUC and NVIDIA (2016) under a [Creative Commons Attribution-NonCommercial 4.0 License](#).