

# GPU Accelerated Computing

## Assignment 03 Report

Syed Muhammad Ali — Saad Abdul Hakim

April 6, 2024

## 1 Introduction

Path and RayTracers have a lot of room for parallelization since they compute color values for each pixel independently. In this assignment we have implemented small pathtracer from Kevin Pearson using CUDA.

## 2 Approach

### 2.1 Identifying points where parallelization can be done

The first task involved identifying the parts of code where parallelization can be done. Since we need to assign a value to each pixel of the image, we can parallelize the loop which iterates over all the pixels serially on CPU (or minutely in parallel using CPU multithreading) to gain significant speedup.

### 2.2 The Approach

We started with the 99 liner cpp implementation of smallpt from the provided reference and then started converting function and classes into device functions. Then we made a cuda kernel for converting the serial assignment of pixel values to parallel. Finally, this kernel is called inside the main function.

### 2.3 Code parts

#### 2.3.1 Converting Vec to float3

Since CUDA supports the float3 data type, we can convert all the instance of Vec class in smallpt to float3 instances. Additionally we use vecutils.h, header which provides device functions for different operations such as dot, cross product using float3 object.

#### 2.3.2 Conversion of struct methods to device functions

Now we can convert the structs such as Ray and Sphere so that they can be initialized on the GPU

```
1 struct Ray
2 {
3     float3 o, d;
4     __device__ Ray(float3 o_, float3 d_) : o(o_), d(d_) {}
5 };
```

Listing 1: Converting structs and its methods into device class/functions

Also, we convert other miscellaneous functions such as clamp and toInt into host,device functions

```
1 inline __host__ __device__ double
2 clamp(double x)
3 {
4     return x < 0 ? 0 : x > 1 ? 1
5         : x;
```

```

6 }
7
8 inline __host__ __device__ int toInt(double x) { return int(pow(clamp(x), 1 / GAMMA) *
    255 + .5); }

```

Listing 2: Converting structs and its methods into device class/functions

## 2.4 Random Initialization

We use curand for random uniform initialization instead of erand48() or std uniform generator since we need to initialize random values on device.

### 2.4.1 The Radiance Function

The main part of the ray tracer is the radiance function which actually is responsible for tracing the ray. We first tried to make a recursive device function just like in smallpt. However, we found that due to nondeterministic stack initializations, it performs very poorly, and in some cases randomly crashes. Therefore, we converted the recursive radiance function to a loop-based radiance function. This significantly improved the performance of our code. However, still, we had to recursively call the radiance function if the object material type is transparent since two recursive calls were being made in the smallpt radiance function. This again reduces performance by a bit, but is necessary for getting the exact same output. If we remove this recursive call, then the performance improves, but in scenes with transparent object, they slightly differ from smallpt rendered scene.

```

1  __device__ float3 radiance(Ray &r, curandState state, int d = 0)
2  { // returns ray color
3
4      float3 pixelcolor = make_float3(0.0f, 0.0f, 0.0f);
5      float3 mask = make_float3(1.0f, 1.0f, 1.0f);
6
7      for (int depth = d; depth < 10; depth++)
8      {
9
10         double t; // distance to closest intersection
11         int id = 0; // index of closest intersected sphere
12
13         // test ray for intersection with scene
14         if (!intersect(r, t, id))
15         {
16             return pixelcolor + (BACKGROUND * mask); // for scene 3
17         }
18
19         const Sphere &obj = spheres[id];
20         float3 x = r.o + r.d * t;
21         float3 n = normalize(x - obj.p);
22         float3 nl = dot(n, r.d) < 0 ? n : n * -1;
23         float3 f = obj.c;
24         double p = f.x > f.y && f.x > f.z ? f.x : f.y > f.z ? f.y
25                 : f.z;
26
27         // Russian Roulette. Can be turned on but it reduces code performance (due to
28             branch divergence maybe)
29         // if (depth > 5)
30         // {
31         //     if (curand_uniform(&state) < p)
32         //     {
33         //         f = f * (1 / p);
34         //     }
35         //     else
36         //     {
37         //         pixelcolor += obj.e * mask;
38         //         break;
39         //     }
40         // }
41
42         // create 2 random numbers
43         float r1 = 2 * M_PI * curand_uniform(&state);

```

```

43     float r2 = curand_uniform(&state);
44     float r2s = sqrtf(r2);
45
46     if (obj.refl == DIFF)
47     {
48         float3 w = nl;
49         float3 u = normalize(cross((fabs(w.x) > .1 ? make_float3(0, 1, 0) :
50             make_float3(1, 0, 0)), w));
51         float3 v = cross(w, u);
52         float3 d = normalize(u * cos(r1) * r2s + v * sin(r1) * r2s + w * sqrtf(1 -
53             r2));
54         pixelcolor += mask * obj.e;
55         r.o = x + nl * 0.05f;
56         r.d = d;
57     }
58
59     else if (obj.refl == SPEC)
60     {
61         r.o = x + nl * 0.05f;
62         r.d = r.d - n * 2 * dot(n, r.d);
63         pixelcolor += mask * obj.e;
64     }
65
66     else if (obj.refl == REFR)
67     {
68         Ray reflRay(x + nl * 0.05f, r.d - n * 2 * dot(n, r.d)); // Ideal
69         dielectric REFRACTION
70         bool into = dot(n, nl) > 0; // Ray from
71         outside going in?
72         double nc = 1, nt = 1.5, nnt = into ? nc / nt : nt / nc, ddn = dot(r.d, nl
73             ), cos2t;
74         if ((cos2t = 1 - nnt * nnt * (1 - ddn * ddn)) < 0) // Total internal
75             reflection
76         {
77             pixelcolor += mask * obj.e;
78             r.o = reflRay.o;
79             r.d = reflRay.d;
80         }
81         else
82         {
83             float3 tdir = normalize((r.d * nnt - n * ((into ? 1 : -1) * (ddn * nnt
84                 + sqrt(cos2t)))));
85             double a = nt - nc, b = nt + nc, R0 = a * a / (b * b), c = 1 - (into ?
86                 -ddn : dot(tdir, n));
87             double Re = R0 + (1 - R0) * c * c * c * c * c, Tr = 1 - Re, P = .25 +
88                 .5 * Re, RP = Re / P, TP = Tr / (1 - P);
89             if (depth > 2)
90             {
91                 if (curand_uniform(&state) < P)
92                 {
93                     pixelcolor += mask * obj.e * RP;
94                     r.o = reflRay.o;
95                     r.d = reflRay.d;
96                     // mask *= RP;
97                 }
98                 else
99                 {
100                     pixelcolor += mask * obj.e * TP;
101                     r.o = x + nl * 0.05f;
102                     r.d = tdir;
103                     // mask *= TP;
104                 }
105             }
106         }
107     }
108
109     else
110     {
111         pixelcolor += mask * obj.e * Re;
112         // pixelcolor += radiance(reflRay, state, depth + 1) * Re; // a
113         recursive call. Reduces the performance of code but is
114         important.
115         // r.o = reflRay.o;
116         // r.d = reflRay.d;
117         r.o = x;

```

```

103         r.d = tdir;
104         pixelcolor += mask * obj.e * Tr;
105     }
106 }
107
108 // multiply with colour of object
109 mask *= f;
110 }
111
112 return pixelcolor;
113 }
114 }

```

Listing 3: Converting structs and its methods into device class/functions

## 2.5 pathTracer Kernel

Finally, to take advantage of thread level parallelism, we make a pathtracer kernel to simultaneously calculate pixel values. After calculation, we copy the image on the host.

## 2.6 Constant Scene

Since the scenes are constant and are not modified, we load them on the constant memory. This improves the performance in comparison to using global memory.

```

1  --constant-- Sphere spheres[] = {
2      {1e5f, {1e5f + 1.0f, 40.8f, 81.6f}, {0.0f, 0.0f, 0.0f}, {0.75f, 0.25f, 0.25f},
3          DIFF}, // Left
4      {1e5f, {-1e5f + 99.0f, 40.8f, 81.6f}, {0.0f, 0.0f, 0.0f}, {.25f, .25f, .75f}, DIFF},
5          // Right
6      {1e5f, {50.0f, 40.8f, 1e5f}, {0.0f, 0.0f, 0.0f}, {.75f, .75f, .75f}, DIFF},
7          // Back
8      {1e5f, {50.0f, 40.8f, -1e5f + 600.0f}, {0.0f, 0.0f, 0.0f}, {1.00f, 1.00f, 1.00f},
9          DIFF}, // Frnt
10     {1e5f, {50.0f, 1e5f, 81.6f}, {0.0f, 0.0f, 0.0f}, {.75f, .75f, .75f}, DIFF},
11         // Botm
12     {1e5f, {50.0f, -1e5f + 81.6f, 81.6f}, {0.0f, 0.0f, 0.0f}, {.75f, .75f, .75f}, DIFF},
13         // Top
14     {16.5f, {27.0f, 16.5f, 47.0f}, {0.0f, 0.0f, 0.0f}, {1.0f, 1.0f, 1.0f}, SPEC},
15         // small sphere 1
16     {16.5f, {73.0f, 16.5f, 78.0f}, {0.0f, 0.0f, 0.0f}, {1.0f, 1.0f, 1.0f}, REFR},
17         // small sphere 2
18     {600.0f, {50.0f, 681.6f - .27f, 81.6f}, {12.0f, 12.0f, 12.0f}, {0.0f, 0.0f, 0.0f},
19         DIFF} // Light
20 };

```

Listing 4: Using Constant memory for scenes

Moreover, we had to reconstruct scene 3 completely and scene2 partially since we could not find their exact source codes. Also, it is important to note that we needed to recalibrate the gamma correction and background color values for different scenes to produce the same output as provided in the assignment.

## 3 Results on CPU

## 4 Results on GPU

Following are the results obtained on parallel pathtracer at 5000 samples.

## 5 Performance Comparison

At smaller number of samples, the CPU and GPU have comparable performance, however, when the number of samples are increased, the GPU code outperforms CPU by a significant margin. Following

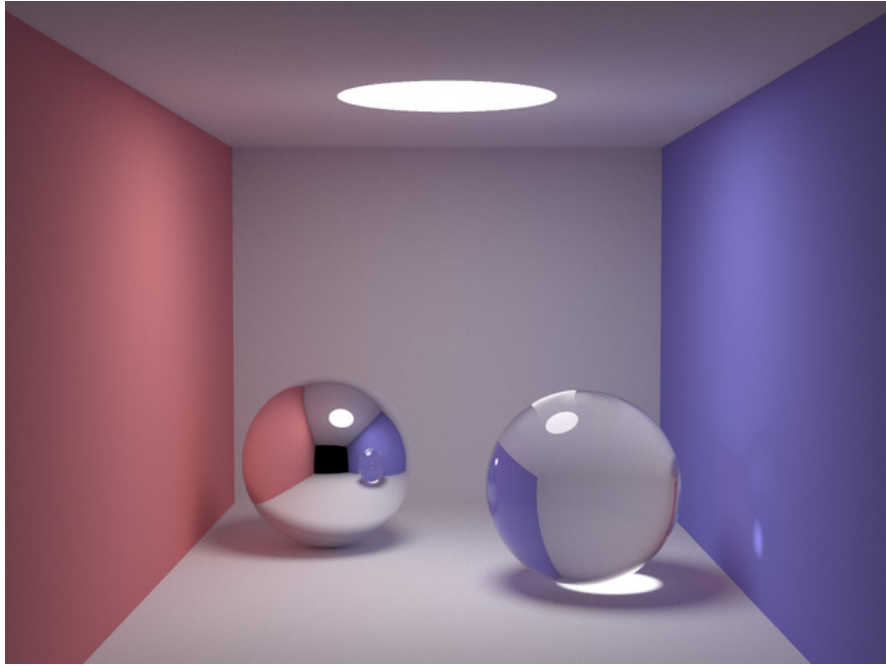


Figure 1: Scene 1 on CPU

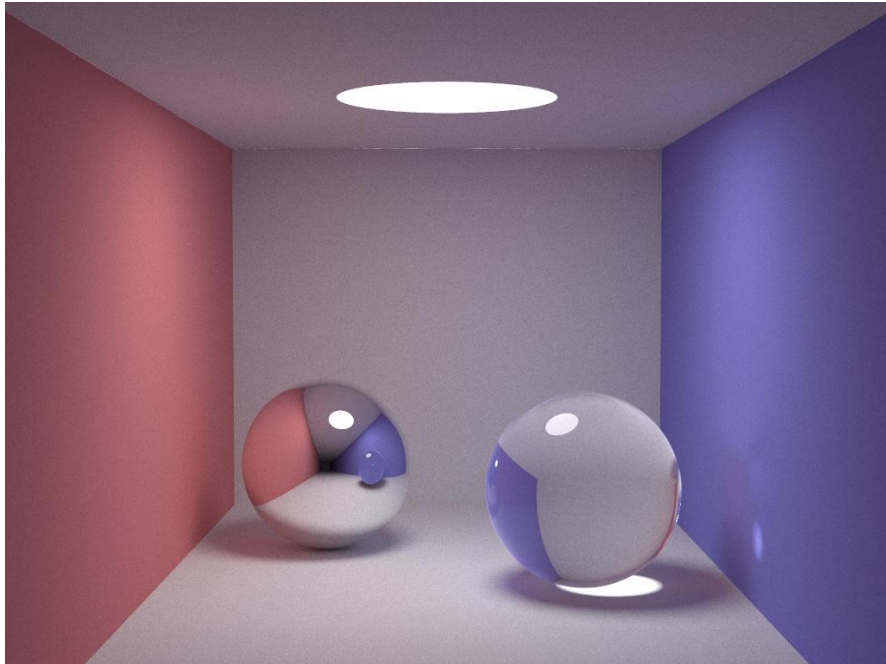


Figure 2: Scene 1 on GPU

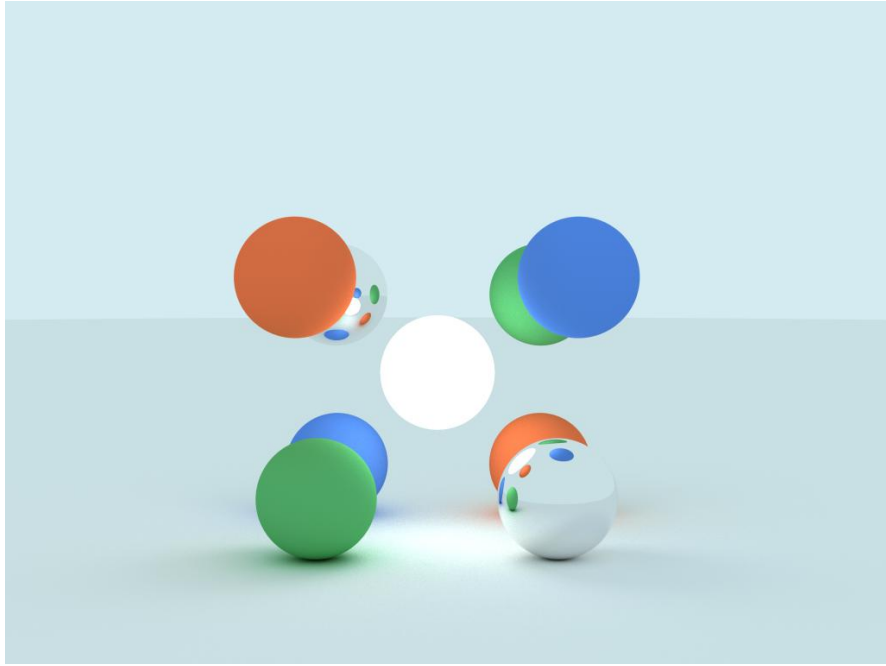


Figure 3: Scene 2: 9 sphere scene on GPU

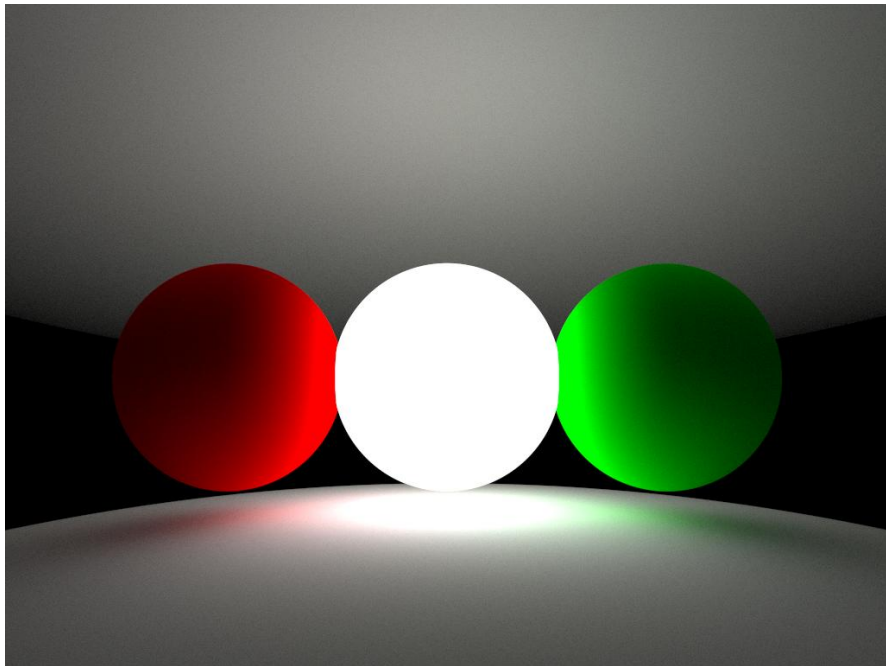


Figure 4: Scene 3: Wikipedia 3 spheres on GPU

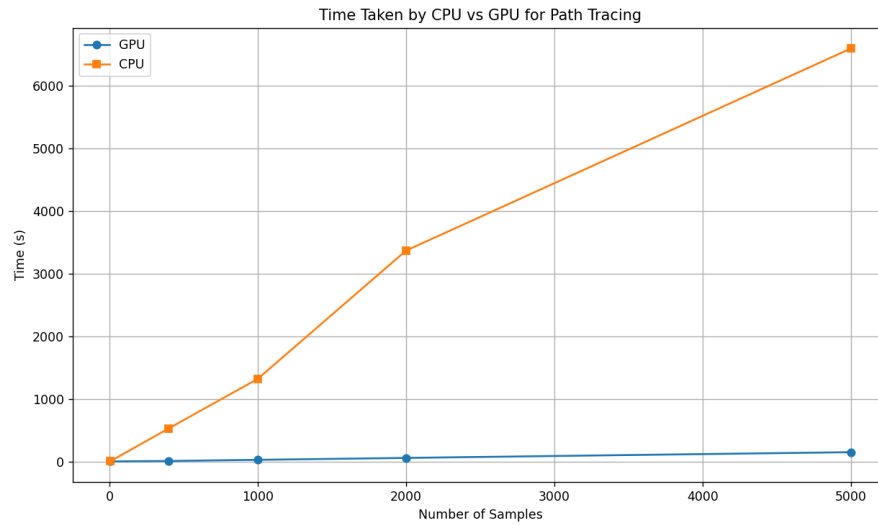


Figure 5: Performnace of GPU vs CPU

are the time comparisons of rendering scene 1 on CPU (with OpenMP) and GPU (Cuda) for different no of samples.

Hence, the GPU gives a significant speedup compared to GPU.