# CUDA raytracing algorithm for visualizing discrete element model output

Anders Damsgaard Christensen, *20062213*

*Abstract*—A raytracing algorithm is constructed using the CUDA API for visualizing output from a CUDA discrete element model, which outputs spatial information in dynamic particle systems. The raytracing algorithm is optimized with constant memory and compilation flags, and performance is measured as a function of the number of particles and the number of pixels. The execution time is compared to equivalent CPU code, and the speedup under a variety of conditions is found to have a mean value of 55.6 times.

*Index Terms*—CUDA, discrete element method, raytracing

## I. INTRODUCTION

VISUALIZING systems containing many spheres using traditional object-order graphics rendering can often result in very high computational requirements, as the usual automated approach is to construct a meshed surface with a specified resolution for each sphere. The memory requirements are thus quite high, as each surface will consist of many vertices. Raytracing [1] is a viable alternative, where spheric entities are saved as data structures with a centre coordinate and a radius. The rendering is performed on the base of these values, which results in a perfectly smooth surfaced sphere. To accelerate the rendering, the algorithm is constructed utilizing the CUDA API [2], where the problem is divided into $n \times m$ threads, corresponding to the desired output image resolution. Each thread iterates through all particles and applies a simple shading model to determine the final RGB values of the pixel.

Previous studies of GPU or CUDA implementations of ray tracing algorithms reported major speedups, compared to corresponding CPU applications (e.g. [3], [4], [5], [6]). None of the software was however found to be open-source and GPL licensed, so a simple raytracer was constructed, customized to render particles, where the data was stored in a specific data format.

### A. Discrete Element Method

The input particle data to the raytracer is the output of a custom CUDA-based Discrete Element Method (DEM) application currently in development. The DEM model is used to numerically simulate the response of a drained, soft, granular sediment bed upon normal stresses and shearing velocities similar to subglacial environments under ice streams [7]. In contrast to laboratory experiments on granular material, the discrete element method [8] approach allows close monitoring of the progressive deformation, where all involved physical

parameters of the particles and spatial boundaries are readily available for continuous inspection.

The discrete element method (DEM) is a subtype of molecular dynamics (MD), and discretizes time into sufficiently small timesteps, and treats the granular material as discrete grains, interacting through contact forces. Between time steps, the particles are allowed to overlap slightly, and the magnitude of the overlap and the kinematic states of the particles is used to compute normal- and shear components of the contact force. The particles are treated as spherical entities, which simplifies the contact search. The spatial simulation domain is divided using a homogeneous, uniform, cubic grid, which greatly reduces the amount of possible contacts that are checked during each timestep. The grid-particle list is sorted using Thrust[1], and updated each timestep. The new particle positions and kinematic values are updated by inserting the resulting force and torque into Newton's second law, and using a Taylor-based second order integration scheme to calculate new linear and rotational accelerations, velocities and positions.

### B. Application usage

The CUDA DEM application is a command line executable, and writes updated particle information to custom binary files with a specific interval. This raytracing algorithm is constructed to also run from the command line, be non-interactive, and write output images in the PPM image format. This format is chosen to allow rendering to take place on cluster nodes with CUDA compatible devices.

Both the CUDA DEM and raytracing applications are open-source[2], although still under heavy development.

This document consists of a short introduction to the basic mathematics behind the ray tracing algorithm, an explaination of the implementation using the CUDA API [2] and a presentation of the results. The CUDA device source code and C++ host source code for the ray tracing algorithm can be found in the appendix, along with instructions for compilation and execution of the application.

## II. RAY TRACING ALGORITHM

The goal of the ray tracing algorithm is to compute the shading of each pixel in the image [9]. This is performed by creating a viewing ray from the eye into the scene, finding the closest intersection with a scene object, and computing the resulting color. The general structure of the program is demonstrated in the following pseudo-code:

[1]http://code.google.com/p/thrust/
[2]http://users-cs.au.dk/adc/files/sphere.tar.gz

```
for each pixel do
   compute viewing ray origin and direction
   iterate through objects and find the closest hit
   set pixel color to value computed from hit ↩
      point, light, n
```

The implemented code does not utilize recursive rays, since the modeled material grains are matte in appearance.

### A. Ray generation

The rays are in vector form defined as:

$$\mathbf{p}(t) = \mathbf{e} + t(\mathbf{s} - \mathbf{e}) \tag{1}$$

The perspective can be either *orthograpic*, where all viewing rays have the same direction, but different starting points, or use *perspective projection*, where the starting point is the same, but the direction is slightly different [9]. For the purposes of this application, a perspective projection was chosen, as it results in the most natural looking image. The ray data structures were held flexible enough to allow an easy implementation of orthographic perspective, if this is desired at a later point.

The ray origin $\mathbf{e}$ is the position of the eye, and is constant. The direction is unique for each ray, and is computed using:

$$\mathbf{s} - \mathbf{e} = -d\mathbf{w} + u\mathbf{u} + v\mathbf{v} \tag{2}$$

where $\{\mathbf{u}, \mathbf{v}, \mathbf{w}\}$ are the orthonormal bases of the camera coordinate system, and $d$ is the focal length [9]. The camera coordinates of pixel $(i, j)$ in the image plane, $u$ and $v$, are calculated by:

$$u = l + (r - l)(i + 0.5)/n$$

$$v = b + (t - b)(j + 0.5)/m$$

where $l$, $r$, $t$ and $b$ are the positions of the image borders (left, right, top and bottom) in camera space. The values $n$ and $m$ are the number of pixels in each dimension.

### B. Ray-sphere intersection

Given a sphere with a center $\mathbf{c}$, and radius $R$, a equation can be constrained, where $\mathbf{p}$ are all points placed on the sphere surface:

$$(\mathbf{p} - \mathbf{c}) \cdot (\mathbf{p} - \mathbf{c}) - R^2 = 0 \tag{3}$$

By substituting the points $\mathbf{p}$ with ray equation 1, and rearranging the terms, a quadratic equation emerges:

$$(\mathbf{d} \cdot \mathbf{d})t^2 + 2\mathbf{d} \cdot (\mathbf{e} - \mathbf{c})t + (\mathbf{e} - \mathbf{c}) \cdot (\mathbf{e} - \mathbf{c}) - R^2 = 0 \tag{4}$$

The number of ray steps $t$ is the only unknown, so the number of intersections is found by calculating the determinant:

$$\Delta = (2(\mathbf{d} \cdot (\mathbf{e} - \mathbf{c})))^2 - 4(\mathbf{d} \cdot \mathbf{d})((\mathbf{e} - \mathbf{c}) \cdot (\mathbf{e} - \mathbf{c}) - R^2 \tag{5}$$

A negative value denotes no intersection between the sphere and the ray, a value of zero means that the ray touches the sphere at a single point (ignored in this implementation), and a positive value denotes that there are two intersections, one when the ray enters the sphere, and one when it exits. In the code, a conditional branch checks wether the determinant is

positive. If this is the case, the distance to the intersection in ray "steps" is calculated using:

$$t = \frac{-\mathbf{d} \cdot (\mathbf{e} - \mathbf{c}) \pm \sqrt{\eta}}{(\mathbf{d} \cdot \mathbf{d})} \tag{6}$$

where

$$\eta = (\mathbf{d} \cdot (\mathbf{e} - \mathbf{c}))^2 - (\mathbf{d} \cdot \mathbf{d})((\mathbf{e} - \mathbf{c}) \cdot (\mathbf{e} - \mathbf{c}) - R^2)$$

Only the smallest intersection ($t_{\mathrm{minus}}$) is calculated, since this marks the point where the sphere enters the particle. If this value is smaller than previous intersection distances, the intersection point $\mathbf{p}$ and surface normal $\mathbf{n}$ at the intersection point is calculated:

$$\mathbf{p} = \mathbf{e} + t_{\mathrm{minus}}\mathbf{d} \tag{7}$$

$$\mathbf{n} = 2(\mathbf{p} - \mathbf{c}) \tag{8}$$

The intersection distance in vector steps ($t_{\mathrm{minus}}$) is saved in order to allow comparison of the distance with later intersections.

### C. Pixel shading

The pixel is shaded using *Lambertian* shading [9], where the pixel color is proportional to the angle between the light vector ($\mathbf{l}$) and the surface normal. An ambient shading component is added to simulate global illumination, and prevent that the spheres are completely black:

$$L = k_a I_a + k_d I_d \max(0, (\mathbf{n} \cdot \mathbf{l})) \tag{9}$$

where the $a$ and $d$ subscripts denote the ambient and diffusive (Lambertian) components of the ambient/diffusive coefficients ($k$) and light intensities ($I$). The pixel color $L$ is calculated once per color channel.

### D. Computational implementation

The above routines were first implemented in CUDA for device execution, and afterwards ported to a CPU C++ equivalent, used for comparing performance. The CPU raytracing algorithm was optimized to shared-memory parallelism using OpenMP [10]. The execution method can be chosen when launching the raytracer from the command line, see the appendix for details. In the CPU implementation, all data was stored in linear arrays of the right size, ensuring 100% memory efficiency.

## III. CUDA IMPLEMENTATION

When constructing the algorithm for execution on the GPGPU device, the data-parallel nature of the problem (SIMD: single instruction, multiple data) is used to deconstruct the rendering task into a single thread per pixel. Each thread iterates through all particles, and ends up writing the resulting color to the image memory.

The application starts by reading the discrete element method data from a custom binary file. The particle data, consisting of position vectors in three-dimensional Euclidean space ($\mathbf{R}^3$) and particle radii, is stored together in a `float4` array, with the particle radius in the `w` position. This has
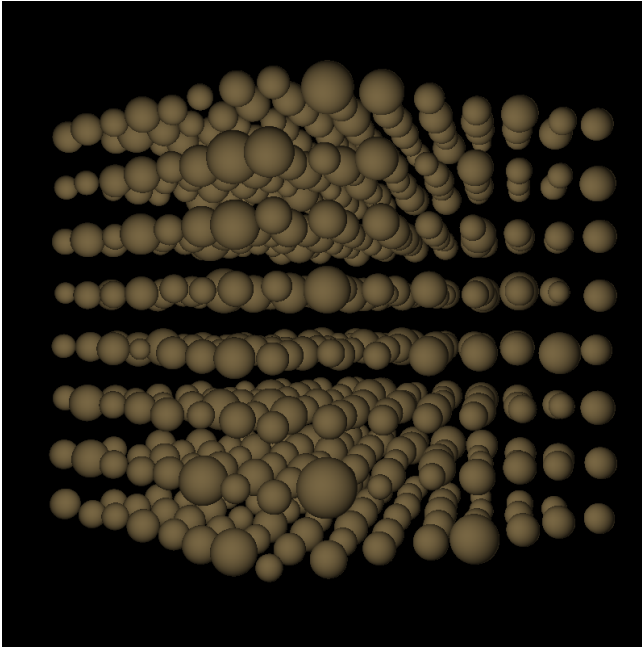
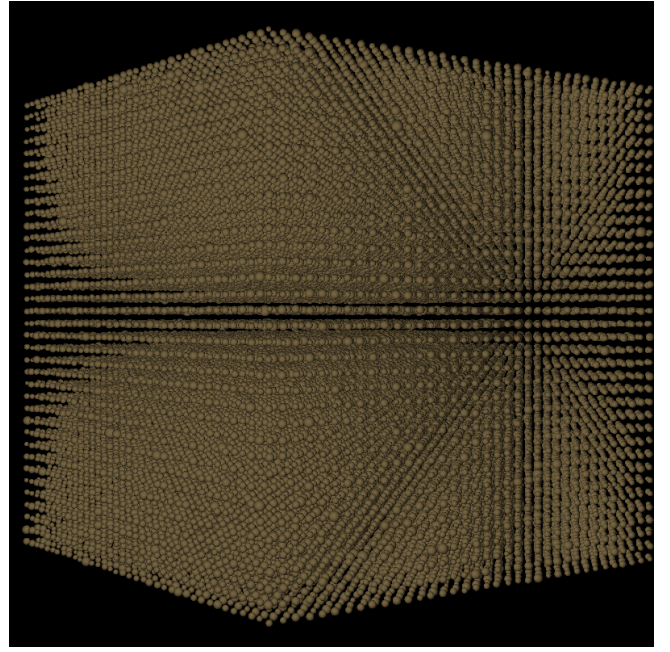Fig. 1.    Sample output of GPU raytracer rendering of 512 particles.



Fig. 2.    Sample output of GPU raytracer rendering of 50 653 particles.

large advantages to storing the data in separate `float3` and `float` arrays; Using `float4` (instead of `float3`) data allows coalesced memory access [11] to the arrays of data in device memory, resulting in efficient memory requests and transfers [12], and the data access pattern is coherent and convenient. Other three-component vectors were also stored as `float4` for the same reasons, even though this sometimes caused a slight memory redundancy. The image data is saved in a three-channel linear `unsigned char` array. Global memory access are coalesced whenever possible. Divergent branches in the kernel code were avoided as much as possible [11].

The algorithm starts by allocating memory on the device for the particle data, the ray parameters, and the image RGB values. Afterwards, all particle data is transfered from the host-to the device memory.

All pixel values are initialized to $[R, G, B] = [0, 0, 0]$, which serves as the image background color. Afterwards, a kernel is executed with a thread for each pixel, testing for intersections between the pixel's viewing ray and all particles, and returning the closest particle. This information is used when computing the shading of the pixel.

After all pixel values have been computed, the image data is transfered back to the host memory, and written to the disk. The application ends by liberating dynamically allocated memory on both the device and the host.

### A. Thread and block layout

The thread/block layout passed during kernel launches is arranged in the following manner:

```
dim3 threads(16, 16);
dim3 blocks((width+15)/16, (height+15)/16);
```

The image pixel position of the thread can be determined from the thread- and block index and dimensions. The layout corresponds to a thread tile size of 256, and a dynamic number of blocks, ensured to fit the image dimensions with only small eventual redundancy [13]. Since this method will initialize extra threads in most situations, all kernels (with return type `void`) start by checking wether the thread-/block index actually falls inside of the image dimensions:

```
int i = threadIdx.x + blockIdx.x * ↵
    blockDim.x;
int j = threadIdx.y + blockIdx.y * ↵
    blockDim.y;
unsigned int mempos = x + y * blockDim.x ↵
    * gridDim.x;
if (mempos > pixels)
  return;
```

The linear memory position (`mempos`) is used as the index when reading or writing to the linear arrays residing in global device memory.

### B. Image output

After completing all pixel shading computations on the device, the image data is transfered back to the host memory, and together with a header written to a PPM[3] image file. This file is converted to the PNG format using ImageMagick.

### C. Performance

Since this simple raytracing algorithm generates a single non-recursive ray for each pixel, which in turn checks all spheres for intersection, the application is expected to scale in the form of $O(n \times m \times N)$, where $n$ and $m$ are the output image dimensions in pixels, and $N$ is the number of particles.

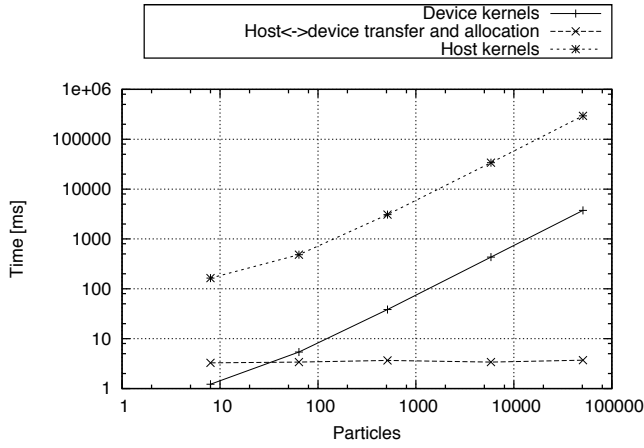[3]http://paulbourke.net/dataformats/ppm/

Fig. 3. Performance scaling with varying particle numbers at image dimensions 800 by 800 pixels.
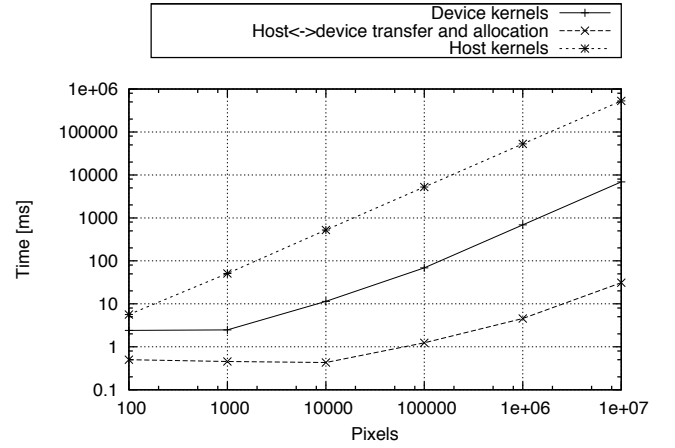


Fig. 4. Performance scaling with varying image dimensions ($n \times m$) with 5832 particles.

The data transfer between the host and device is kept at a bare minimum, as the intercommunication is considered a bottleneck in relation to the potential device performance[11]. Thread synchronization points are only inserted were necessary, and the code is optimized by the compilers to the target architecture (see appendix).

The host execution time was profiled using a `clock()` based CPU timer from `time.h`, which was normalized using the constant `CLOCKS_PER_SEC`.

The device execution time was profiled using two `cudaEvent_t` timers, one measuring the time spent in the entire device code section, including device memory allocation, data transfer to- and from the host, execution of the kernels, and memory deallocation. The other timer only measured time spent in the kernels. The threads were synchronized before stopping the timers. A simple CPU timer using `clock()` will *not* work, since control is returned to the host thread before the device code has completed all tasks.

Figures 3 and 4 show the profiling results, where the number of particles and the image dimensions were varied. With exception of executions with small image dimensions, the kernel execution time results agree with the $O(n \times m \times N)$ scaling predicion.

The device memory allocation and data transfer was also profiled, and turns out to be only weakly dependant on the particle numbers (fig. 3), but more strongly correlated to image dimensions (fig. 4). As with kernel execution times, the execution time converges against an overhead value at small image dimensions.

The CPU time spent in the host kernels proves to be linear with the particle numbers, and linear with the image dimensions. This is due to the non-existant overhead caused by initialization of the device code, and reduced memory transfer.

The ratio between CPU computational times and the sum of the device kernel execution time and the host—device memory transfer and additional memory allocation was calculated, and had a mean value of $55.6$ and a variance of $739$ out of the 11 comparative measurements presented in the figures. It should be noted, that the smallest speedups were recorded when using very small image dimensions, probably unrealistic in real use.

As the number of particles are not known by compilation, it is not possible to store particle positions and -radii in constant memory. Shared memory was also on purpose avoided, since the memory per thread block (64 kb) would not be sufficient in rendering of simulations containing containing more than $16\,000$ particles ($16\,000$ `float4` values). The constant memory was however utilized for storing the camera related parameters; the orthonormal base vectors, the observer position, the image dimensions, the focal length, and the light vector.

Previous GPU implementations often rely on k-D trees, constructed as an sorting method for static scene objects[3], [5]. A k-D tree implementation would drastically reduce the global memory access induced by each thread, so it is therefore the next logical step with regards to optimizing the ray tracing algorithm presented here.

## IV. CONCLUSION

This document presented the implementation of a basic ray tracing algorithm, utilizing the highly data-parallel nature of the problem when porting the work load to CUDA. Performance tests showed the expected, linear correlation between image dimensions, particle numbers and execution time. Comparisons with an equivalent CPU algorithm showed large speedups, typically up to two orders of magnitude. This speedup did not come at a cost of less correct results.

The final product will come into good use during further development and completion of the CUDA DEM particle model, and is ideal since it can be used for offline rendering on dedicated, heterogeneous GPU-CPU computing nodes. The included device code will be the prefered method of execution, whenever the host system allows it.

### REFERENCES

[1] T. Whitted, "An improved illumination model for shaded display," *Communications of the ACM*, vol. 23, no. 6, pp. 343–349, 1980.
[2] NVIDIA, *CUDA C Programming Guide*, 3rd ed., NVIDIA Corporation: Santa Clara, CA, USA, Oct 2010.
[3] D. Horn, J. Sugerman, M. Houston, and P. Hanrahan, "Interactive k-D tree GPU raytracing," *Association for Computing Machinery, Inc.*, pp. 167–174 pp., 2007.

[4] M. Shih, Y. Chiu, Y. Chen, and C. Chang, "Real-time ray tracing with cuda," *Algorithms and Architectures for Parallel Processing*, pp. 327–337, 2009.

[5] S. Popov, J. Günther, H. Seidel, and P. Slusallek, "Stackless kd-tree traversal for high performance gpu ray tracing," in *Computer Graphics Forum*, vol. 26, no. 3. Wiley Online Library, 2007, pp. 415–424.

[6] D. Luebke and S. Parker, "Interactive ray tracing with cuda," *NVIDIA Technical Presentation, SIGGRAPH*, 2008.

[7] D. Evans, E. Phillips, J. Hiemstra, and C. Auton, "Subglacial till: formation, sedimentary characteristics and classification," *Earth-Science Reviews*, vol. 78, no. 1-2, pp. 115–176, 2006.

[8] P. Cundall and O. Strack, "A discrete numerical model for granular assemblies," *Géotechnique*, vol. 29, pp. 47–65, 1979.

[9] P. Shirley, M. Ashikhmin, M. Gleicher, S. Marschner, E. Reinhard, K. Sung, W. Thompson, and P. Willemsen, *Fundamentals of computer graphics*, 3rd ed. AK Peters, 2009.

[10] B. Chapman, G. Jost, and R. Van Der Pas, *Using OpenMP: portable shared memory parallel programming*. The MIT Press, 2007, vol. 10.

[11] NVIDIA, *CUDA C Best Practices Guide Version*, 3rd ed., NVIDIA Corporation: Santa Clara, CA, USA, Aug 2010, cA Patent 95,050.

[12] L. Nyland, M. Harris, and J. Prins, "Fast n-body simulation with cuda," *GPU gems*, vol. 3, pp. 677–695, 2007.

[13] J. Sanders and E. Kandrot, *CUDA by example*. Addison-Wesley, 2010.

# APPENDIX A
## TEST ENVIRONMENT

The raytracing algorithm was developed, tested and profiled on a mid 2010 Mac Pro with a 2.8 Ghz Quad-Core Intel Xeon CPU and a NVIDIA Quadro 4000 for Mac, dedicated to CUDA applications. The CUDA driver was version 4.0.50, the CUDA compilation tools release 4.0, V0.2.1221. The GCC tools were version 4.2.1. Each CPU core is multithreaded by two threads for a total of 8 threads.

The CUDA source code was compiled with `nvcc`, and linked to `g++` compiled C++ source code with `g++`. For all benchmark tests, the code was compiled with the following commands:

```
g++ −c −Wall −O3 −arch x86_64 −fopenmp ...
nvcc −c −use_fast_math −gencode ←
    arch=compute_20 , code=sm_20 −m64 ...
g++ −arch x86_64 −lcuda −lcudart −fopenmp ←
    *.o −o rt
```

When profiling device code performance, the application was executed two times, and the time of the second run was noted. This was performed to avoid latency caused by device driver initialization.

The host system was measured to have a memory bandwidth of 4642.1 MB/s when transfering data from the host to the device, and 3805.6 MB/s when transfering data from the device to the host.

# APPENDIX B
## SOURCE CODE

The entire source code, as well as input data files, can be found in the following archive http://users-cs.au.dk/adc/files/sphere-rt.tar.gz The source code is built and run with the commands:

```
make
make run
```

With the `make run` command, the Makefile uses ImageMagick to convert the PPM file to PNG format, and the OS X command `open` to display the image. Other input data files are included with other particle number magnitudes. The syntax for the raytracer application is the following:

```
./rt <CPU | GPU> <sphere−binary.bin> ←
    <width> <height> <output−image.ppm>
```

This appendix contains the following source code files:

### LISTINGS

### A. CUDA raytracing source code

Listing 1. rt_kernel.h

```
1   #ifndef RT_KERNEL_H_
2   #define RT_KERNEL_H_
3
4   #include <vector_functions.h>
5
6   // Constants
7   __constant__ float4 const_u;
8   __constant__ float4 const_v;
9   __constant__ float4 const_w;
10  __constant__ float4 const_eye;
11  __constant__ float4 const_imgplane;
12  __constant__ float  const_d;
13  __constant__ float4 const_light;
14
15
16  // Host prototype functions
17
18  extern "C"
19  void cameraInit(float4 eye, float4 lookat, float ←
        imgw, float hw_ratio);
20
21  extern "C"
22  void checkForCudaErrors(const char* ←
        checkpoint_description);
23
24  extern "C"
25  int rt(float4* p, const unsigned int np,
26      rgb* img, const unsigned int width, const ←
            unsigned int height,
27      f3 origo, f3 L, f3 eye, f3 lookat, float imgw);
28
29  #endif
```

Listing 2. rt_kernel.cu

```
1   #include <iostream>
2   #include <cutil_math.h>
3   #include "header.h"
4   #include "rt_kernel.h"
5
6   __inline__ __host__ __device__ float3 f4_to_f3(float4 ←
        in)
7   {
8     return make_float3(in.x, in.y, in.z);
9   }
10
11  __inline__ __host__ __device__ float4 f3_to_f4(float3 ←
        in)
12  {
13    return make_float4(in.x, in.y, in.z, 0.0f);
14  }
15
16  // Kernel for initializing image data
17  __global__ void imageInit(unsigned char* _img, ←
        unsigned int pixels)
18  {
19    // Compute pixel position from threadIdx/blockIdx
20    int x = threadIdx.x + blockIdx.x * blockDim.x;
21    int y = threadIdx.y + blockIdx.y * blockDim.y;
22    unsigned int mempos = x + y * blockDim.x * gridDim.x;
23    if (mempos > pixels)
24      return;
```

```
25
26     _img[mempos*4]       = 0;    // Red channel
27     _img[mempos*4 + 1] = 0;      // Green channel
28     _img[mempos*4 + 2] = 0;      // Blue channel
29  }
30
31  // Calculate ray origins and directions
32  __global__ void rayInitPerspective(float4* _ray_origo,
33                           float4* _ray_direction,
34                           float4 eye,
35                                         unsigned int width,
36                           unsigned int height)
37  {
38      // Compute pixel position from threadIdx/blockIdx
39      int i = threadIdx.x + blockIdx.x * blockDim.x;
40      int j = threadIdx.y + blockIdx.y * blockDim.y;
41      unsigned int mempos = i + j * blockDim.x * gridDim.x;
42      if (mempos > width*height)
43        return;
44
45      // Calculate pixel coordinates in image plane
46      float p_u = const_imgplane.x + (const_imgplane.y -
            const_imgplane.x)
47                  * (i + 0.5f) / width;
48      float p_v = const_imgplane.z + (const_imgplane.w -
            const_imgplane.z)
49                  * (j + 0.5f) / height;
50
51      // Write ray origo and direction to global memory
52      _ray_origo[mempos]      = const_eye;
53      _ray_direction[mempos] = -const_d*const_w +
            p_u*const_u + p_v*const_v;
54  }
55
56  // Check wether the pixel's viewing ray intersects
          with the spheres,
57  // and shade the pixel correspondingly
58  __global__ void rayIntersectSpheres(float4* _ray_origo,
59                                        float4*
                                          _ray_direction,
60                           float4* _p,
61                   unsigned char* _img,
62                   unsigned int pixels,
63                   unsigned int np)
64  {
65      // Compute pixel position from threadIdx/blockIdx
66      int x = threadIdx.x + blockIdx.x * blockDim.x;
67      int y = threadIdx.y + blockIdx.y * blockDim.y;
68      unsigned int mempos = x + y * blockDim.x * gridDim.x;
69      if (mempos > pixels)
70        return;
71
72      // Read ray data from global memory
73      float3 e = f4_to_f3(_ray_origo[mempos]);
74      float3 d = f4_to_f3(_ray_direction[mempos]);
75      //float  step = length(d);
76
77      // Distance, in ray steps, between object and eye
            initialized with a large value
78      float tdist = 1e10f;
79
80      // Surface normal at closest sphere intersection
81      float3 n;
82
83      // Intersection point coordinates
84      float3 p;
85
86      // Iterate through all particles
87      for (unsigned int i=0; i<np; i++) {
88
89          // Read sphere coordinate and radius
90          float3 c = f4_to_f3(_p[i]);
91          float  R = _p[i].w;
92
93          // Calculate the discriminant: d = B^2 - 4AC
94          float Delta =
                (2.0f*dot(d,(e-c)))*(2.0f*dot(d,(e-c)))   // B^2
95                  - 4.0f*dot(d,d)     // -4*A
96                  * (dot((e-c),(e-c)) - R*R);   // C
97
98          // If the determinant is positive, there are two
                solutions
99          // One where the line enters the sphere, and one
                where it exits
100         if (Delta > 0.0f) {
101
102             // Calculate roots, Shirley 2009 p. 77
103             float t_minus = ((dot(-d,(e-c)) - sqrt(
                    dot(d,(e-c))*dot(d,(e-c)) - dot(d,d)
104                     * (dot((e-c),(e-c)) - R*R) ) ) /
                        dot(d,d));
105
106             // Check wether intersection is closer than
                    previous values
107             if (fabs(t_minus) < tdist) {
108         p = e + t_minus*d;
109         tdist = fabs(t_minus);
110         n = normalize(2.0f * (p - c));   // Surface normal
111             }
112
113         } // End of solution branch
114
115     } // End of particle loop
116
117     // Write pixel color
118     if (tdist < 1e10) {
119
120         // Lambertian shading parameters
121         float dotprod = fabs(dot(n, f4_to_f3(const_light)));
122         float I_d = 40.0f;  // Light intensity
123         float k_d = 5.0f;   // Diffuse coefficient
124
125         // Ambient shading
126         float k_a = 10.0f;
127         float I_a = 5.0f;
128
129         // Write shading model values to pixel color
                channels
130         _img[mempos*4]       = (unsigned char) ((k_d * I_d
                * dotprod
131                         + k_a * I_a)*0.48f);
132         _img[mempos*4 + 1] = (unsigned char) ((k_d * I_d
                * dotprod
133                         + k_a * I_a)*0.41f);
134         _img[mempos*4 + 2] = (unsigned char) ((k_d * I_d
                * dotprod
135                         + k_a * I_a)*0.27f);
136     }
137 }
138
139
140 extern "C"
141 __host__ void cameraInit(float4 eye, float4 lookat,
        float imgw, float hw_ratio)
142 {
143     // Image dimensions in world space (l, r, b, t)
144     float4 imgplane = make_float4(-0.5f*imgw,
            0.5f*imgw, -0.5f*imgw*hw_ratio,
            0.5f*imgw*hw_ratio);
145
146     // The view vector
147     float4 view = eye - lookat;
148
149     // Construct the camera view orthonormal base
150     float4 v = make_float4(0.0f, 1.0f, 0.0f, 0.0f);  //
            v: Pointing upward
151     float4 w = -view/length(view);              // w:
            Pointing backwards
152     float4 u = make_float4(cross(make_float3(v.x, v.y,
            v.z),
153                 make_float3(w.x, w.y, w.z)),
154             0.0f); // u: Pointing right
155
156     // Focal length 20% of eye vector length
157     float d = length(view)*0.8f;
158
159     // Light direction (points towards light source)
160     float4 light =
            normalize(-1.0f*eye*make_float4(1.0f, 0.2f,
            0.6f, 0.0f));
161
162     std::cout << "  Transfering camera values to
            constant memory\n";
163
164     cudaMemcpyToSymbol("const_u", &u, sizeof(u));
165     cudaMemcpyToSymbol("const_v", &v, sizeof(v));
166     cudaMemcpyToSymbol("const_w", &w, sizeof(w));
167     cudaMemcpyToSymbol("const_eye", &eye, sizeof(eye));
168     cudaMemcpyToSymbol("const_imgplane", &imgplane,
            sizeof(imgplane));
169     cudaMemcpyToSymbol("const_d", &d, sizeof(d));
```

```
170    cudaMemcpyToSymbol("const_light", &light, ←
           sizeof(light));
171  }
172
173  // Check for CUDA errors
174  extern "C"
175  __host__ void checkForCudaErrors(const char* ←
         checkpoint_description)
176  {
177    cudaError_t err = cudaGetLastError();
178    if (err != cudaSuccess) {
179      std::cout << "\nCuda_error_detected,_checkpoint:_←
             " << checkpoint_description
180                << "\nError_string:_" << ←
                       cudaGetErrorString(err) << "\n";
181      exit(EXIT_FAILURE);
182    }
183  }
184
185
186  // Wrapper for the rt kernel
187  extern "C"
188  __host__ int rt(float4* p, unsigned int np,
189                   rgb* img, unsigned int width, ←
                        unsigned int height,
190           f3 origo, f3 L, f3 eye, f3 lookat, float ←
                  imgw) {
191
192    using std::cout;
193
194    cout << "Initializing_CUDA:\n";
195
196    // Initialize GPU timestamp recorders
197    float t1, t2;
198    cudaEvent_t t1_go, t2_go, t1_stop, t2_stop;
199    cudaEventCreate(&t1_go);
200    cudaEventCreate(&t2_go);
201    cudaEventCreate(&t2_stop);
202    cudaEventCreate(&t1_stop);
203
204    // Start timer 1
205    cudaEventRecord(t1_go, 0);
206
207    // Allocate memory
208    cout << "__Allocating_device_memory\n";
209    static float4 *_p;              // Particle positions ←
             (x,y,z) and radius (w)
210    static unsigned char *_img;    // RGBw values in ←
             image
211    static float4 *_ray_origo;     // Ray origo (x,y,z)
212    static float4 *_ray_direction; // Ray direction ←
             (x,y,z)
213    cudaMalloc((void**)&_p, np*sizeof(float4));
214    cudaMalloc((void**)&_img, ←
           width*height*4*sizeof(unsigned char));
215    cudaMalloc((void**)&_ray_origo, ←
           width*height*sizeof(float4));
216    cudaMalloc((void**)&_ray_direction, ←
           width*height*sizeof(float4));
217
218    // Transfer particle data
219    cout << "__Transfering_particle_data:_host_→_←
             device\n";
220    cudaMemcpy(_p, p, np*sizeof(float4), ←
           cudaMemcpyHostToDevice);
221
222    // Check for errors after memory allocation
223    checkForCudaErrors("CUDA_error_after_memory_←
             allocation");
224
225    // Arrange thread/block structure
226    unsigned int pixels = width*height;
227    float hw_ratio = (float)height/(float)width;
228    dim3 threads(16,16);
229    dim3 blocks((width+15)/16, (height+15)/16);
230
231    // Start timer 2
232    cudaEventRecord(t2_go, 0);
233
234    // Initialize image to background color
235    imageInit<<< blocks, threads >>>(_img, pixels);
236
237    // Initialize camera
238    cameraInit(make_float4(eye.x, eye.y, eye.z, 0.0f),
239                make_float4(lookat.x, lookat.y, ←
                      lookat.z, 0.0f),
240                imgw, hw_ratio);
241    checkForCudaErrors("CUDA_error_after_cameraInit");
242
243    // Construct rays for perspective projection
244    rayInitPerspective <<< blocks, threads >>>(
245        _ray_origo, _ray_direction,
246        make_float4(eye.x, eye.y, eye.z, 0.0f),
247        width, height);
248
249    // Find closest intersection between rays and spheres
250    rayIntersectSpheres <<< blocks, threads >>>(
251        _ray_origo, _ray_direction,
252        _p, _img, pixels, np);
253
254
255    // Make sure all threads are done before continuing ←
             CPU control sequence
256    cudaThreadSynchronize();
257
258    // Check for errors
259    checkForCudaErrors("CUDA_error_after_kernel_←
             execution");
260
261    // Stop timer 2
262    cudaEventRecord(t2_stop, 0);
263    cudaEventSynchronize(t2_stop);
264
265    // Transfer image data from device to host
266    cout << "__Transfering_image_data:_device_→_host\n";
267    cudaMemcpy(img, _img, ←
           width*height*4*sizeof(unsigned char), ←
           cudaMemcpyDeviceToHost);
268
269    // Free dynamically allocated device memory
270    cudaFree(_p);
271    cudaFree(_img);
272    cudaFree(_ray_origo);
273    cudaFree(_ray_direction);
274
275    // Stop timer 1
276    cudaEventRecord(t1_stop, 0);
277    cudaEventSynchronize(t1_stop);
278
279    // Calculate time spent in t1 and t2
280    cudaEventElapsedTime(&t1, t1_go, t1_stop);
281    cudaEventElapsedTime(&t2, t2_go, t2_stop);
282
283    // Report time spent
284    cout << "__Time_spent_on_entire_GPU_routine:_"
285         << t1 << "_ms\n";
286    cout << "__-_Kernels:_" << t2 << "_ms\n"
287         << "__-_Memory_alloc._and_transfer:_" << t1-t2 ←
                << "_ms\n";
288
289    // Return successfully
290    return 0;
291  }
```

## B. CPU raytracing source code

Listing 3.  rt_kernel_cpu.h
```
1  #ifndef RT_KERNEL_CPU_H_
2  #define RT_KERNEL_CPU_H_
3
4  #include <vector_functions.h>
5
6  // Host prototype functions
7
8  void cameraInit(float3 eye, float3 lookat, float ←
        imgw, float hw_ratio);
9
10 int rt_cpu(float4* p, const unsigned int np,
11       rgb* img, const unsigned int width, const ←
             unsigned int height,
12       f3 origo, f3 L, f3 eye, f3 lookat, float imgw);
13
14 #endif
```

Listing 4.  rt_kernel_cpu.cpp
```
1  #include <iostream>
2  #include <cstdio>
3  #include <cmath>
4  #include <time.h>
5  #include <cuda.h>
```

```
6   #include <cutil_math.h>
7   #include <string.h>
8   #include "header.h"
9   #include "rt_kernel_cpu.h"
10
11  // Constants
12  float3 constc_u;
13  float3 constc_v;
14  float3 constc_w;
15  float3 constc_eye;
16  float4 constc_imgplane;
17  float   constc_d;
18  float3 constc_light;
19
20  __inline__ float3 f4_to_f3(float4 in)
21  {
22    return make_float3(in.x, in.y, in.z);
23  }
24
25  __inline__ float4 f3_to_f4(float3 in)
26  {
27    return make_float4(in.x, in.y, in.z, 0.0f);
28  }
29
30  __inline__ float lengthf3(float3 in)
31  {
32    return sqrt(in.x*in.x + in.y*in.y + in.z*in.z);
33  }
34
35  // Kernel for initializing image data
36  void imageInit_cpu(unsigned char* _img, unsigned int ←
        pixels)
37  {
38    for (unsigned int mempos=0; mempos<pixels; ←
          mempos++) {
39      _img[mempos*4]     = 0; // Red channel
40      _img[mempos*4 + 1] = 0; // Green channel
41      _img[mempos*4 + 2] = 0; // Blue channel
42    }
43  }
44
45  // Calculate ray origins and directions
46  void rayInitPerspective_cpu(float3* _ray_origo,
47                  float3* _ray_direction,
48              float3 eye,
49                      unsigned int width,
50              unsigned int height)
51  {
52    int i;
53    #pragma omp parallel for
54    for (i=0; i<width; i++) {
55      for (unsigned int j=0; j<height; j++) {
56
57        unsigned int mempos = i + j*height;
58
59        // Calculate pixel coordinates in image plane
60        float p_u = constc_imgplane.x + ←
              (constc_imgplane.y - constc_imgplane.x)
61      * (i + 0.5f) / width;
62        float p_v = constc_imgplane.z + ←
              (constc_imgplane.w - constc_imgplane.z)
63      * (j + 0.5f) / height;
64
65        // Write ray origo and direction to global memory
66        _ray_origo[mempos]     = constc_eye;
67        _ray_direction[mempos] = -constc_d*constc_w + ←
              p_u*constc_u + p_v*constc_v;
68      }
69    }
70  }
71
72  // Check wether the pixel's viewing ray intersects ←
        with the spheres,
73  // and shade the pixel correspondingly
74  void rayIntersectSpheres_cpu(float3* _ray_origo,
75                        float3* _ray_direction,
76                        float4* _p,
77              unsigned char* _img,
78              unsigned int pixels,
79              unsigned int np)
80  {
81    int mempos;
82    #pragma omp parallel for
83    for (mempos=0; mempos<pixels; mempos++) {
84
85      // Read ray data from global memory
86      float3 e = _ray_origo[mempos];
87      float3 d = _ray_direction[mempos];
88      //float   step = lengthf3(d);
89
90      // Distance, in ray steps, between object and eye ←
            initialized with a large value
91      float tdist = 1e10f;
92
93      // Surface normal at closest sphere intersection
94      float3 n;
95
96      // Intersection point coordinates
97      float3 p;
98
99      // Iterate through all particles
100     for (unsigned int i=0; i<np; i++) {
101
102       // Read sphere coordinate and radius
103       float3 c = f4_to_f3(_p[i]);
104       float   R = _p[i].w;
105
106       // Calculate the discriminant: d = B^2 - 4AC
107       float Delta = ←
            (2.0f*dot(d,(e-c)))*(2.0f*dot(d,(e-c))) ←
                // B^2
108     - 4.0f*dot(d,d) // -4*A
109     * (dot((e-c),(e-c)) - R*R);   // C
110
111       // If the determinant is positive, there are ←
            two solutions
112       // One where the line enters the sphere, and ←
            one where it exits
113       if (Delta > 0.0f) {
114
115     // Calculate roots, Shirley 2009 p. 77
116     float t_minus = ((dot(-d,(e-c)) - sqrt( ←
            dot(d,(e-c))*dot(d,(e-c)) - dot(d,d)
117     * (dot((e-c),(e-c)) - R*R) ) ) / dot(d,d));
118
119     // Check wether intersection is closer than ←
            previous values
120     if (fabs(t_minus) < tdist) {
121       p = e + t_minus*d;
122       tdist = fabs(t_minus);
123       n = normalize(2.0f * (p - c));   // Surface normal
124     }
125
126       } // End of solution branch
127
128     } // End of particle loop
129
130     // Write pixel color
131     if (tdist < 1e10) {
132
133       // Lambertian shading parameters
134       float dotprod = fabs(dot(n, constc_light));
135       float I_d = 40.0f;  // Light intensity
136       float k_d = 5.0f;  // Diffuse coefficient
137
138       // Ambient shading
139       float k_a = 10.0f;
140       float I_a = 5.0f;
141
142       // Write shading model values to pixel color ←
            channels
143       _img[mempos*4]     = (unsigned char) ((k_d * ←
            I_d * dotprod
144     + k_a * I_a)*0.48f);
145       _img[mempos*4 + 1] = (unsigned char) ((k_d * ←
            I_d * dotprod
146     + k_a * I_a)*0.41f);
147       _img[mempos*4 + 2] = (unsigned char) ((k_d * ←
            I_d * dotprod
148     + k_a * I_a)*0.27f);
149     }
150   }
151 }
152
153
154 void cameraInit_cpu(float3 eye, float3 lookat, float ←
        imgw, float hw_ratio)
155 {
156   // Image dimensions in world space (l, r, b, t)
157   float4 imgplane = make_float4(-0.5f*imgw, ←
            0.5f*imgw, -0.5f*imgw*hw_ratio, ←
            0.5f*imgw*hw_ratio);
```

```
158
159        // The view vector
160        float3 view = eye - lookat;
161
162        // Construct the camera view orthonormal base
163        float3 v = make_float3(0.0f, 1.0f, 0.0f);   // v: ←
                Pointing upward
164        float3 w = -view/lengthf3(view);            // w: ←
                Pointing backwards
165        float3 u = cross(make_float3(v.x, v.y, v.z), ←
                make_float3(w.x, w.y, w.z)); // u: Pointing right
166
167        // Focal length 20% of eye vector length
168        float d = lengthf3(view)*0.8f;
169
170        // Light direction (points towards light source)
171        float3 light = ←
                normalize(-1.0f*eye*make_float3(1.0f, 0.2f, ←
                0.6f));
172
173        std::cout << "__Transfering_camera_values_to_←
                constant_memory\n";
174
175        constc_u = u;
176        constc_v = v;
177        constc_w = w;
178        constc_eye = eye;
179        constc_imgplane = imgplane;
180        constc_d = d;
181        constc_light = light;
182    }
183
184
185    // Wrapper for the rt algorithm
186    int rt_cpu(float4* p, unsigned int np,
187            rgb* img, unsigned int width, unsigned int ←
                height,
188            f3 origo, f3 L, f3 eye, f3 lookat, float imgw) {
189
190        using std::cout;
191
192        cout << "Initializing_CPU_raytracer:\n";
193
194        // Initialize GPU timestamp recorders
195        float t1_go, t2_go, t1_stop, t2_stop;
196
197        // Start timer 1
198        t1_go = clock();
199
200        // Allocate memory
201        cout << "__Allocating_device_memory\n";
202        static unsigned char *_img;        // RGBw values in ←
                image
203        static float3* _ray_origo;         // Ray origo (x,y,z)
204        static float3* _ray_direction;     // Ray direction ←
                (x,y,z)
205        _img         = new unsigned char[width*height*4];
206        _ray_origo   = new float3[width*height];
207        _ray_direction = new float3[width*height];
208
209        // Arrange thread/block structure
210        unsigned int pixels = width*height;
211        float hw_ratio = (float)height/(float)width;
212
213        // Start timer 2
214        t2_go = clock();
215
216        // Initialize image to background color
217        imageInit_cpu(_img, pixels);
218
219        // Initialize camera
220        cameraInit_cpu(make_float3(eye.x, eye.y, eye.z),
221                        make_float3(lookat.x, lookat.y, ←
                            lookat.z),
222                    imgw, hw_ratio);
223
224        // Construct rays for perspective projection
225        rayInitPerspective_cpu(
226            _ray_origo, _ray_direction,
227            make_float3(eye.x, eye.y, eye.z),
228            width, height);
229
230        // Find closest intersection between rays and spheres
231        rayIntersectSpheres_cpu(
232            _ray_origo, _ray_direction,
233            p, _img, pixels, np);
234
235        // Stop timer 2
236        t2_stop = clock();
237
238        memcpy(img, _img, sizeof(unsigned char)*pixels*4);
239
240        // Free dynamically allocated device memory
241        delete [] _img;
242        delete [] _ray_origo;
243        delete [] _ray_direction;
244
245        // Stop timer 1
246        t1_stop = clock();
247
248        // Report time spent
249        cout << "__Time_spent_on_entire_CPU_raytracing_←
                routine:_"
250            << (t1_stop-t1_go)/CLOCKS_PER_SEC*1000.0 << "_←
                ms\n";
251        cout << "___Functions:_" << ←
                (t2_stop-t2_go)/CLOCKS_PER_SEC*1000.0 << "_ms\n";
252
253        // Return successfully
254        return 0;
255    }
```

## Appendix C
## Bloopers

This section contains pictures of the unfinished raytracer in a malfunctioning state, which can provide interesting, however unusable, results.
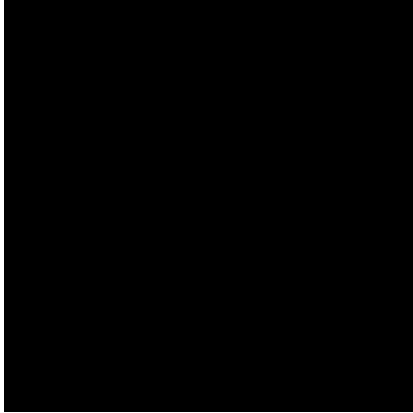


Fig. 5. An end result encountered way too many times, e.g. after inadvertently dividing the focal length with zero.
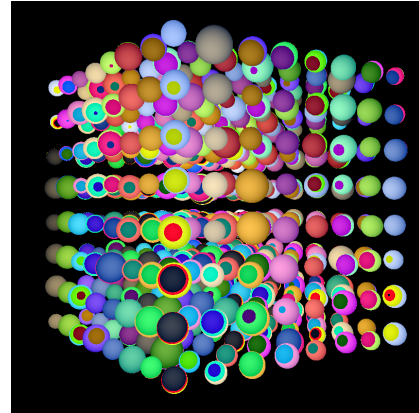


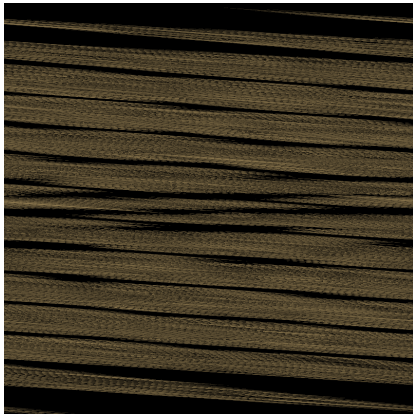Fig. 8. Colorful result caused by forgetting to normalize the normal vector calculation.



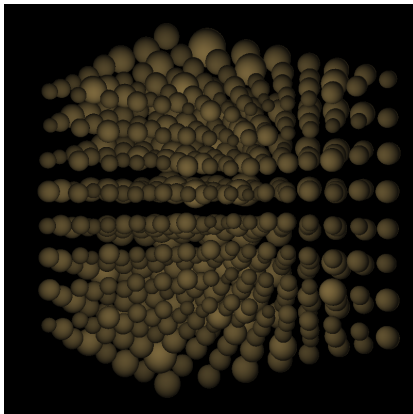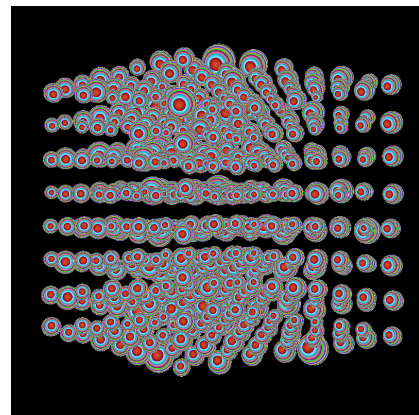Fig. 6. Twisted appearance caused by wrong thread/block layout.



Fig. 9. Another colorful result with wrong order of particles, cause by using `abs()` (meant for integers) instead of `fabs()` (meant for floats).



Fig. 7. A problem with the distance determination caused the particles to be displayed in a wrong order.