# 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:

Contact: anders.damsgaard@geo.au.dk
Webpage: http://users-cs.au.dk/adc
Manuscript, last revision: May 9, 2012.

---

[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 ($t_{\text{minus}}$) 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:

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

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

The intersection distance in vector steps ($t_{\text{minus}}$) 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
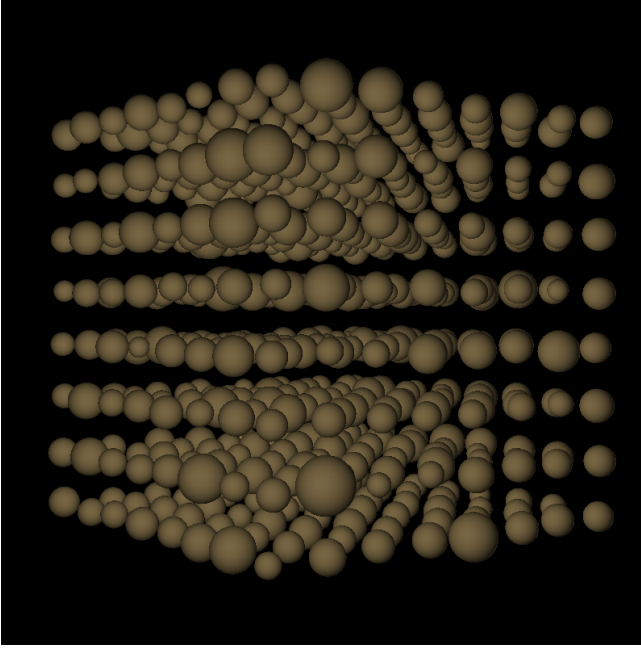
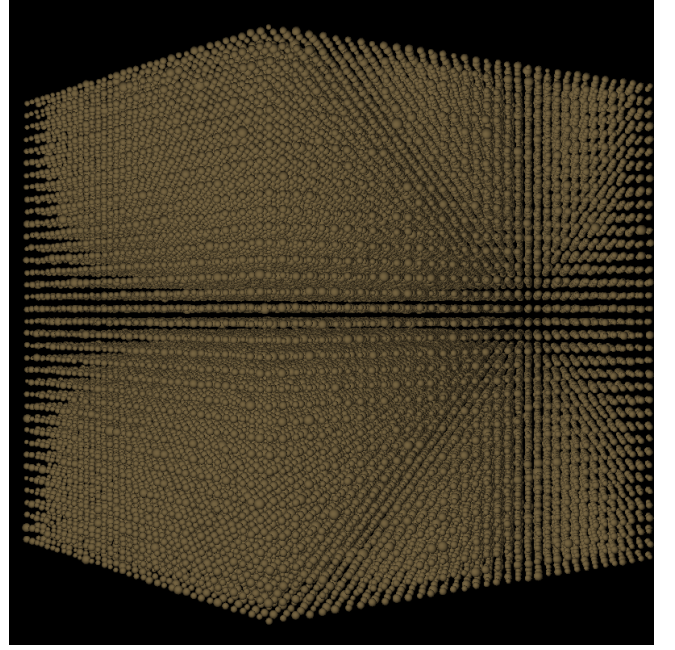Fig. 1. Sample output of GPU raytracer rendering of 512 particles.



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.

---

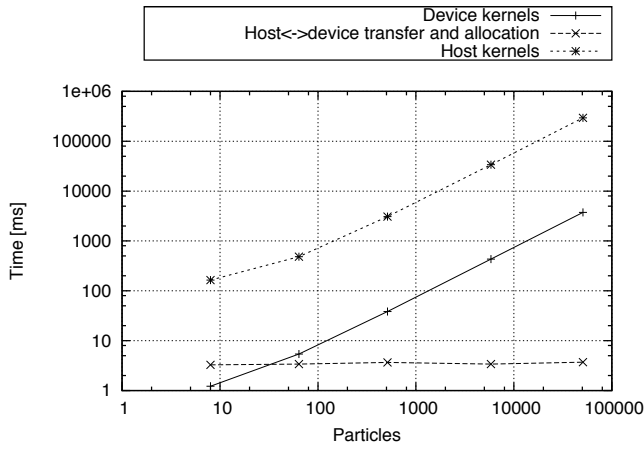[3]http://paulbourke.net/dataformats/ppm/

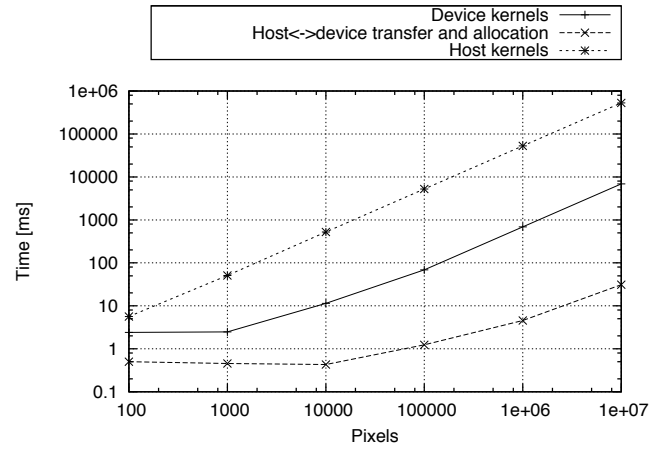Fig. 3. Performance scaling with varying particle numbers at image dimensions 800 by 800 pixels.



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  #include "header.h"
6
7  // Constants
8  __constant__ float4 const_u;
9  __constant__ float4 const_v;
10 __constant__ float4 const_w;
11 __constant__ float4 const_eye;
12 __constant__ float4 const_imgplane;
13 __constant__ float  const_d;
14 __constant__ float4 const_light;
15 __constant__ unsigned int const_pixels;
16 __constant__ Inttype const_np;
17
18
19 // Host prototype functions
20
21 extern "C"
22 void cameraInit(float4 eye, float4 lookat, float ←
       imgw, float hw_ratio,
23            unsigned int pixels, Inttype np);
24
25 extern "C"
26 void checkForCudaErrors(const char* ←
       checkpoint_description);
27
28 extern "C"
29 int rt(float4* p, const Inttype np,
30         rgb* img, const unsigned int width, const ←
           unsigned int height,
31         f3 origo, f3 L, f3 eye, f3 lookat, float imgw,
32         const int visualize, float* color, const float ←
           max_val);
33
34 #endif
```

Listing 2. rt-kernel.cu

```
1  #include <iostream>
2  #include <cutil_math.h>
3  #include "header.h"
4  #include "rt−kernel.h"
5  #include "colorbar.h"
6
7  unsigned int iDivUp(unsigned int a, unsigned int b) {
8    return (a % b != 0) ? (a / b + 1) : (a / b);
9  }
10
11 __inline__ __host__ __device__ float3 f4_to_f3(float4 ←
       in)
12 {
13   return make_float3(in.x, in.y, in.z);
14 }
15
16 __inline__ __host__ __device__ float4 f3_to_f4(float3 ←
       in)
17 {
18   return make_float4(in.x, in.y, in.z, 0.0f);
19 }
```

```
20
21   // Kernel for initializing image data
22   __global__ void imageInit(unsigned char* _img, ←
         unsigned int pixels)
23   {
24     // Compute pixel position from threadIdx/blockIdx
25     unsigned int mempos = threadIdx.x + blockIdx.x * ←
           blockDim.x;
26     if (mempos > pixels)
27       return;
28
29     _img[mempos*4]     = 255; // Red channel
30     _img[mempos*4 + 1] = 255; // Green channel
31     _img[mempos*4 + 2] = 255; // Blue channel
32   }
33
34   // Calculate ray origins and directions
35   __global__ void rayInitPerspective(float4* _ray_origo,
36                          float4* _ray_direction,
37                          float4 eye,
38                                        unsigned int width,
39                          unsigned int height)
40   {
41     // Compute pixel position from threadIdx/blockIdx
42     unsigned int mempos = threadIdx.x + blockIdx.x * ←
           blockDim.x;
43     if (mempos > width*height)
44       return;
45
46     // Calculate 2D position from linear index
47     unsigned int i = mempos % width;
48     unsigned int j = (int)floor((float)mempos/width) % ←
           width;
49
50     // Calculate pixel coordinates in image plane
51     float p_u = const_imgplane.x + (const_imgplane.y − ←
           const_imgplane.x)
52                 * (i + 0.5f) / width;
53     float p_v = const_imgplane.z + (const_imgplane.w − ←
           const_imgplane.z)
54                 * (j + 0.5f) / height;
55
56     // Write ray origo and direction to global memory
57     _ray_origo[mempos]     = const_eye;
58     _ray_direction[mempos] = −const_d*const_w + ←
           p_u*const_u + p_v*const_v;
59   }
60
61   // Check wether the pixel's viewing ray intersects ←
         with the spheres,
62   // and shade the pixel correspondingly
63   __global__ void rayIntersectSpheres(float4* _ray_origo,
64                              float4* ←
                                 _ray_direction,
65                              float4* _p,
66                  unsigned char* _img)
67   {
68     // Compute pixel position from threadIdx/blockIdx
69     unsigned int mempos = threadIdx.x + blockIdx.x * ←
           blockDim.x;
70     if (mempos > const_pixels)
71       return;
72
73     // Read ray data from global memory
74     float3 e = f4_to_f3(_ray_origo[mempos]);
75     float3 d = f4_to_f3(_ray_direction[mempos]);
76     //float  step = length(d);
77
78     // Distance, in ray steps, between object and eye ←
           initialized with a large value
79     float tdist = 1e10f;
80
81     // Surface normal at closest sphere intersection
82     float3 n;
83
84     // Intersection point coordinates
85     float3 p;
86
87     // Iterate through all particles
88     for (Inttype i=0; i<const_np; ++i) {
89
90       // Read sphere coordinate and radius
91       float3 c   = f4_to_f3(_p[i]);
92       float  R   = _p[i].w;
93
94       // Calculate the discriminant: d = B^2 − 4AC
```

```
95       float Delta = ←
             (2.0f*dot(d,(e−c)))*(2.0f*dot(d,(e−c)))   // B^2
96                 − 4.0f*dot(d,d)    // −4*A
97                 * (dot((e−c),(e−c)) − R*R);   // C
98
99       // If the determinant is positive, there are two ←
             solutions
100      // One where the line enters the sphere, and one ←
             where it exits
101      if (Delta > 0.0f) {
102
103        // Calculate roots, Shirley 2009 p. 77
104        float t_minus = ((dot(−d,(e−c)) − sqrt( ←
             dot(d,(e−c))*dot(d,(e−c)) − dot(d,d)
105                 * (dot((e−c),(e−c)) − R*R) ) ) / ←
                     dot(d,d));
106
107        // Check wether intersection is closer than ←
             previous values
108        if (fabs(t_minus) < tdist) {
109      p = e + t_minus*d;
110      tdist = fabs(t_minus);
111      n = normalize(2.0f * (p − c));    // Surface normal
112        }
113
114      } // End of solution branch
115
116    } // End of particle loop
117
118    // Write pixel color
119    if (tdist < 1e10f) {
120
121      // Lambertian shading parameters
122      float dotprod = fmax(0.0f,dot(n, ←
             f4_to_f3(const_light)));
123      float I_d = 40.0f;  // Light intensity
124      float k_d = 5.0f;   // Diffuse coefficient
125
126      // Ambient shading
127      float k_a = 10.0f;
128      float I_a = 5.0f; // 5.0 for black background
129
130      // Write shading model values to pixel color ←
             channels
131      _img[mempos*4]     = (unsigned char) ((k_d * I_d ←
             * dotprod
132                 + k_a * I_a)*0.48f);
133      _img[mempos*4 + 1] = (unsigned char) ((k_d * I_d ←
             * dotprod
134                 + k_a * I_a)*0.41f);
135      _img[mempos*4 + 2] = (unsigned char) ((k_d * I_d ←
             * dotprod
136                 + k_a * I_a)*0.27f);
137
138    }
139  }
140
141  // Check wether the pixel's viewing ray intersects ←
         with the spheres,
142  // and shade the pixel correspondingly using a colormap
143  __global__ void rayIntersectSpheresColormap(float4* ←
         _ray_origo,
144                                      float4* ←
                                         _ray_direction,
145                  float4* _p,
146                  float* _color,
147                  unsigned char* _img,
148                  float max_val)
149  {
150    // Compute pixel position from threadIdx/blockIdx
151    unsigned int mempos = threadIdx.x + blockIdx.x * ←
           blockDim.x;
152    if (mempos > const_pixels)
153      return;
154
155    // Read ray data from global memory
156    float3 e = f4_to_f3(_ray_origo[mempos]);
157    float3 d = f4_to_f3(_ray_direction[mempos]);
158    //float  step = length(d);
159
160    // Distance, in ray steps, between object and eye ←
           initialized with a large value
161    float tdist = 1e10f;
162
163    // Surface normal at closest sphere intersection
164    float3 n;
```

```
165
166       // Intersection point coordinates
167       float3 p;
168
169       //float fieldval;
170       unsigned int hitidx;
171
172       // Iterate through all particles
173       for (Inttype i=0; i<const_np; ++i) {
174
175          // Read sphere coordinate and radius
176          float3 c       = f4_to_f3(_p[i]);
177          float  R       = _p[i].w;
178          //float  fieldval_tmp = _linarr[i];
179
180          // Calculate the discriminant: d = B^2 - 4AC
181          float Delta = ←
              (2.0f*dot(d,(e-c)))*(2.0f*dot(d,(e-c)))  // B^2
182                  - 4.0f*dot(d,d)    // -4*A
183                  * (dot((e-c),(e-c)) - R*R);  // C
184
185          // If the determinant is positive, there are two ←
                 solutions
186          // One where the line enters the sphere, and one ←
                 where it exits
187          if (Delta > 0.0f) {
188          //if (Delta > 0.0f && fieldval_tmp/max_value > ←
                 0.75f) { // Only render particles with an ←
                 upper 75% value
189
190             // Calculate roots, Shirley 2009 p. 77
191             float t_minus = ((dot(-d,(e-c)) - sqrt( ←
                 dot(d,(e-c))*dot(d,(e-c)) - dot(d,d)
192                     * (dot((e-c),(e-c)) - R*R) ) ) / ←
                         dot(d,d));
193
194             // Check wether intersection is closer than ←
                    previous values
195             if (fabs(t_minus) < tdist) {
196          p = e + t_minus*d;
197          tdist = fabs(t_minus);
198          n = normalize(2.0f * (p - c));   // Surface normal
199          //fieldval = fieldval_tmp;
200          hitidx = i;
201             }
202
203          } // End of solution branch
204
205       } // End of particle loop
206
207       // Write pixel color
208       if (tdist < 1e10) {
209
210          // Fetch particle data used for color
211          float ratio = _color[hitidx] / max_val;
212
213          // Make sure the ratio doesn't exceed the 0.0-1.0 ←
                 interval
214          if (ratio < 0.01f)
215             ratio = 0.01f;
216          if (ratio > 0.99f)
217             ratio = 0.99f;
218
219          // Lambertian shading parameters
220          float dotprod = fmax(0.0f,dot(n, ←
                 f4_to_f3(const_light)));
221          float I_d = 40.0f;  // Light intensity
222          float k_d = 5.0f;   // Diffuse coefficient
223
224          // Ambient shading
225          float k_a = 10.0f;
226          float I_a = 5.0f;
227
228          // Write shading model values to pixel color ←
                 channels
229          _img[mempos*4]     = (unsigned char) ((k_d * I_d ←
                 * dotprod
230                     + k_a * I_a)*red(ratio));
231          _img[mempos*4 + 1] = (unsigned char) ((k_d * I_d ←
                 * dotprod
232                     + k_a * I_a)*green(ratio));
233          _img[mempos*4 + 2] = (unsigned char) ((k_d * I_d ←
                 * dotprod
234                     + k_a * I_a)*blue(ratio));
235       }
236    }
237
238    extern "C"
239    __host__ void cameraInit(float4 eye, float4 lookat, ←
              float imgw, float hw_ratio,
240                    unsigned int pixels, Inttype np)
241    {
242       // Image dimensions in world space (l, r, b, t)
243       float4 imgplane = make_float4(-0.5f*imgw, ←
              0.5f*imgw, -0.5f*imgw*hw_ratio, ←
              0.5f*imgw*hw_ratio);
244
245       // The view vector
246       float4 view = eye - lookat;
247
248       // Construct the camera view orthonormal base
249       float4 up = make_float4(0.0f, 1.0f, 0.0f, 0.0f);  ←
              // Pointing upward along +y
250       float4 w = -view/length(view);          // w: ←
              Pointing backwards
251       float4 u = make_float4(cross(make_float3(up.x, ←
              up.y, up.z),
252                    make_float3(w.x, w.y, w.z)), 0.0f)
253              / length(cross(make_float3(up.x, up.y, ←
                  up.z), make_float3(w.x, w.y, w.z)));
254       float4 v = make_float4(cross(make_float3(w.x, w.y, ←
              w.z), make_float3(u.x, u.y, u.z)), 0.0f);
255
256       // Focal length 20% of eye vector length
257       float d = length(view)*0.8f;
258
259       // Light direction (points towards light source)
260       float4 light = ←
              normalize(-1.0f*eye*make_float4(1.0f, 0.2f, ←
              0.6f, 0.0f));
261
262       std::cout << "  Transfering camera values to ←
              constant memory\n";
263
264       cudaMemcpyToSymbol("const_u", &u, sizeof(u));
265       cudaMemcpyToSymbol("const_v", &v, sizeof(v));
266       cudaMemcpyToSymbol("const_w", &w, sizeof(w));
267       cudaMemcpyToSymbol("const_eye", &eye, sizeof(eye));
268       cudaMemcpyToSymbol("const_imgplane", &imgplane, ←
              sizeof(imgplane));
269       cudaMemcpyToSymbol("const_d", &d, sizeof(d));
270       cudaMemcpyToSymbol("const_light", &light, ←
              sizeof(light));
271       cudaMemcpyToSymbol("const_pixels", &pixels, ←
              sizeof(pixels));
272       cudaMemcpyToSymbol("const_np", &np, sizeof(np));
273    }
274
275    // Check for CUDA errors
276    extern "C"
277    __host__ void checkForCudaErrors(const char* ←
           checkpoint_description)
278    {
279       cudaError_t err = cudaGetLastError();
280       if (err != cudaSuccess) {
281          std::cout << "\nCuda error detected, checkpoint: ←
                 " << checkpoint_description
282                     << "\nError string: " << ←
                         cudaGetErrorString(err) << "\n";
283          exit(EXIT_FAILURE);
284       }
285    }
286
287
288    // Wrapper for the rt kernel
289    extern "C"
290    __host__ int rt(float4* p, Inttype np,
291                    rgb* img, unsigned int width, ←
                        unsigned int height,
292          f3 origo, f3 L, f3 eye, f3 lookat, float imgw,
293          int visualize, float* color, float max_val)
294    {
295       using std::cout;
296
297       cout << "Initializing CUDA:\n";
298
299       // Initialize GPU timestamp recorders
300       float t1, t2;
301       cudaEvent_t t1_go, t2_go, t1_stop, t2_stop;
302       cudaEventCreate(&t1_go);
303       cudaEventCreate(&t2_go);
304       cudaEventCreate(&t2_stop);
```

```
305     cudaEventCreate(&t1_stop);
306
307     // Start timer 1
308     cudaEventRecord(t1_go, 0);
309
310     // Allocate memory
311     cout << "__Allocating_device_memory\n";
312     static float4 *_p;            // Particle positions ←
                                        (x,y,z) and radius (w)
313     static float  *_color;       // Array for ←
                                        linear values to color the particles after
314     static unsigned char *_img;  // RGBw values in ←
                                        image
315     static float4 *_ray_origo;   // Ray origo (x,y,z)
316     static float4 *_ray_direction; // Ray direction ←
                                        (x,y,z)
317     cudaMalloc((void**)&_p, np*sizeof(float4));
318     cudaMalloc((void**)&_color, np*sizeof(float)); // 0 ←
                                        size if visualize = 0;
319     cudaMalloc((void**)&_img, ←
                width*height*4*sizeof(unsigned char));
320     cudaMalloc((void**)&_ray_origo, ←
                width*height*sizeof(float4));
321     cudaMalloc((void**)&_ray_direction, ←
                width*height*sizeof(float4));
322
323     // Transfer particle data
324     cout << "__Transfering_particle_data:_host_->_←
                device\n";
325     cudaMemcpy(_p, p, np*sizeof(float4), ←
                cudaMemcpyHostToDevice);
326     if (visualize == 1)
327       cudaMemcpy(_color, color, np*sizeof(float), ←
                cudaMemcpyHostToDevice);
328
329     // Check for errors after memory allocation
330     checkForCudaErrors("CUDA_error_after_memory_←
                allocation");
331
332     // Arrange thread/block structure
333     unsigned int pixels = width*height;
334     float hw_ratio = (float)height/(float)width;
335     //dim3 threads(16,16);
336     const unsigned int threadsPerBlock = 256;
337     //dim3 blocks((width+15)/16, (height+15)/16);
338     const unsigned int blocksPerGrid = iDivUp(pixels, ←
                threadsPerBlock);
339
340     // Start timer 2
341     cudaEventRecord(t2_go, 0);
342
343     // Initialize image to background color
344     imageInit<<< blocksPerGrid, threadsPerBlock ←
                >>>(_img, pixels);
345
346     // Initialize camera
347     cameraInit(make_float4(eye.x, eye.y, eye.z, 0.0f),
348                make_float4(lookat.x, lookat.y, ←
                        lookat.z, 0.0f),
349           imgw, hw_ratio, pixels, np);
350     checkForCudaErrors("CUDA_error_after_cameraInit");
351
352     // Construct rays for perspective projection
353     rayInitPerspective<<< blocksPerGrid, ←
                threadsPerBlock >>>(
354       _ray_origo, _ray_direction,
355       make_float4(eye.x, eye.y, eye.z, 0.0f),
356       width, height);
357
358     cudaThreadSynchronize();
359
360     // Find closest intersection between rays and spheres
361     if (visualize == 1) { // Visualize pressure
362       rayIntersectSpheresColormap<<< blocksPerGrid, ←
                threadsPerBlock >>>(
363         _ray_origo, _ray_direction,
364         _p, _color, _img, max_val);
365     } else { // Normal visualization
366       rayIntersectSpheres<<< blocksPerGrid, ←
                threadsPerBlock >>>(
367       _ray_origo, _ray_direction,
368       _p, _img);
369     }
370
371     // Make sure all threads are done before continuing ←
                CPU control sequence
```

```
372     cudaThreadSynchronize();
373
374     // Check for errors
375     checkForCudaErrors("CUDA_error_after_kernel_←
                execution");
376
377     // Stop timer 2
378     cudaEventRecord(t2_stop, 0);
379     cudaEventSynchronize(t2_stop);
380
381     // Transfer image data from device to host
382     cout << "__Transfering_image_data:_device_->_host\n";
383     cudaMemcpy(img, _img, ←
                width*height*4*sizeof(unsigned char), ←
                cudaMemcpyDeviceToHost);
384
385     // Free dynamically allocated device memory
386     cudaFree(_p);
387     cudaFree(_color);
388     cudaFree(_img);
389     cudaFree(_ray_origo);
390     cudaFree(_ray_direction);
391
392     // Stop timer 1
393     cudaEventRecord(t1_stop, 0);
394     cudaEventSynchronize(t1_stop);
395
396     // Calculate time spent in t1 and t2
397     cudaEventElapsedTime(&t1, t1_go, t1_stop);
398     cudaEventElapsedTime(&t2, t2_go, t2_stop);
399
400     // Report time spent
401     cout << "__Time_spent_on_entire_GPU_routine:_"
402          << t1 << "_ms\n";
403     cout << "__-_Kernels:_" << t2 << "_ms\n"
404          << "__-_Memory_alloc._and_transfer:_" << t1-t2 ←
                << "_ms\n";
405
406     // Return successfully
407     return 0;
408   }
```

## 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]     = 255;   // Red channel
 40        _img[mempos*4 + 1] = 255;   // Green channel
 41        _img[mempos*4 + 2] = 255;   // 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,j;
 53    unsigned int mempos;
 54    float p_u, p_v;
 55    #pragma omp parallel for private(mempos,j,p_u,p_v)
 56    for (i=0; i<(int)width; ++i) {
 57      for (j=0; j<(int)height; ++j) {
 58
 59        mempos = i + j*width;
 60
 61        // Calculate pixel coordinates in image plane
 62        p_u = constc_imgplane.x + (constc_imgplane.y − ↩
                 constc_imgplane.x)
 63      * (i + 0.5f) / width;
 64        p_v = constc_imgplane.z + (constc_imgplane.w − ↩
                 constc_imgplane.z)
 65      * (j + 0.5f) / height;
 66
 67        // Write ray origo and direction to global memory
 68        _ray_origo[mempos]     = constc_eye;
 69        _ray_direction[mempos] = −constc_d*constc_w + ↩
                 p_u*constc_u + p_v*constc_v;
 70      }
 71    }
 72  }
 73
 74  // Check wether the pixel's viewing ray intersects ↩
          with the spheres,
 75  // and shade the pixel correspondingly
 76  void rayIntersectSpheres_cpu(float3* _ray_origo,
 77                          float3* _ray_direction,
 78                          float4* _p,
 79               unsigned char* _img,
 80               unsigned int pixels,
 81               unsigned int np)
 82  {
 83    long int mempos;
 84    float3 e, d, n, p, c;
 85    float tdist, R, Delta, t_minus, dotprod, I_d, k_d, ↩
           k_a, I_a;
 86    Inttype i;
 87    #pragma omp parallel for ↩
           private(e,d,n,p,c,tdist,R,Delta,t_minus,dotprod,I_d,k_d,k_a,I_a,i)
 88    for (mempos=0; mempos<pixels; ++mempos) {
 89
 90      // Read ray data from global memory
 91      e = _ray_origo[mempos];
 92      d = _ray_direction[mempos];
 93
 94      // Distance, in ray steps, between object and eye ↩
             initialized with a large value
 95      tdist = 1e10f;
 96
 97      // Iterate through all particles
 98      for (i=0; i<np; ++i) {
 99
100        // Read sphere coordinate and radius
101        c = f4_to_f3(_p[i]);
102        R = _p[i].w;
103
104        // Calculate the discriminant: d = B^2 − 4AC
105        Delta = (2.0f*dot(d,(e−c)))*(2.0f*dot(d,(e−c))) ↩
                  // B^2
106      − 4.0f*dot(d,d) // −4*A
107      * (dot((e−c),(e−c)) − R*R);   // C
108
109        // If the determinant is positive, there are ↩
                  two solutions
110        // One where the line enters the sphere, and ↩
                  one where it exits
111        if (Delta > 0.0f) {
112
113      // Calculate roots, Shirley 2009 p. 77
114      t_minus = ((dot(−d,(e−c)) − sqrt( ↩
                  dot(d,(e−c))*dot(d,(e−c)) − dot(d,d)
115        * (dot((e−c),(e−c)) − R*R) ) ) / dot(d,d));
116
117      // Check wether intersection is closer than ↩
                  previous values
118      if (fabs(t_minus) < tdist) {
119        p = e + t_minus*d;
120        tdist = fabs(t_minus);
121        n = normalize(2.0f * (p − c));   // Surface normal
122      }
123
124        } // End of solution branch
125
126      } // End of particle loop
127
128      // Write pixel color
129      if (tdist < 1e10) {
130
131        // Lambertian shading parameters
132        //float dotprod = fabs(dot(n, constc_light));
133        dotprod = fmax(0.0f,dot(n, constc_light));
134        I_d = 40.0f;  // Light intensity
135        k_d = 5.0f;  // Diffuse coefficient
136
137        // Ambient shading
138        k_a = 10.0f;
139        I_a = 5.0f;
140
141        // Write shading model values to pixel color ↩
                  channels
142        _img[mempos*4]     = (unsigned char) ((k_d * ↩
                  I_d * dotprod
143        + k_a * I_a)*0.48f);
144        _img[mempos*4 + 1] = (unsigned char) ((k_d * ↩
                  I_d * dotprod
145        + k_a * I_a)*0.41f);
146        _img[mempos*4 + 2] = (unsigned char) ((k_d * ↩
                  I_d * dotprod
147        + k_a * I_a)*0.27f);
148      }
149    }
150  }
151
152
153  void cameraInit_cpu(float3 eye, float3 lookat, float ↩
          imgw, float hw_ratio)
154  {
155    // Image dimensions in world space (l, r, b, t)
156    float4 imgplane = make_float4(−0.5f*imgw, ↩
             0.5f*imgw, −0.5f*imgw*hw_ratio, ↩
             0.5f*imgw*hw_ratio);
157
158    // The view vector
159    float3 view = eye − lookat;
160
161    // Construct the camera view orthonormal base
162    float3 up = make_float3(0.0f, 1.0f, 0.0f);  // ↩
             Pointing upward along +y
163    float3 w = −view/length(view);             // w: ↩
             Pointing backwards
164    float3 u = cross(up, w) / length(cross(up, w));
165    float3 v = cross(w, u);
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        std::cout << "Rendering␣image...";
184    }
185
186
187    // Wrapper for the rt algorithm
188    int rt_cpu(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␣CPU␣raytracer:\n";
195
196        // Initialize GPU timestamp recorders
197        float t1_go, t2_go, t1_stop, t2_stop;
198
199        // Start timer 1
200        t1_go = clock();
201
202        // Allocate memory
203        cout << "␣␣Allocating␣device␣memory\n";
204        static unsigned char ∗_img;          // RGBw values in ←
                image
205        static float3∗ _ray_origo;        // Ray origo (x,y,z)
206        static float3∗ _ray_direction;    // Ray direction ←
                (x,y,z)
207        _img       = new unsigned char[width∗height∗4];
208        _ray_origo     = new float3[width∗height];
209        _ray_direction = new float3[width∗height];
210
211        // Arrange thread/block structure
212        unsigned int pixels = width∗height;
213        float hw_ratio = (float)height/(float)width;
214
215        // Start timer 2
216        t2_go = clock();
217
218        // Initialize image to background color
219        imageInit_cpu(_img, pixels);
220
221        // Initialize camera
222        cameraInit_cpu(make_float3(eye.x, eye.y, eye.z),
223                    make_float3(lookat.x, lookat.y, ←
                        lookat.z),
224                imgw, hw_ratio);
225
226        // Construct rays for perspective projection
227        rayInitPerspective_cpu(
228            _ray_origo, _ray_direction,
229            make_float3(eye.x, eye.y, eye.z),
230            width, height);
231
232        // Find closest intersection between rays and spheres
233        rayIntersectSpheres_cpu(
234            _ray_origo, _ray_direction,
235            p, _img, pixels, np);
236
237        // Stop timer 2
238        t2_stop = clock();
239
240        memcpy(img, _img, sizeof(unsigned char)∗pixels∗4);
241
242        // Free dynamically allocated device memory
243        delete [] _img;
244        delete [] _ray_origo;
245        delete [] _ray_direction;
246
247        // Stop timer 1
248        t1_stop = clock();
249
250        // Report time spent
251        cout << "␣done.\n"
252            << "␣␣Time␣spent␣on␣entire␣CPU␣raytracing␣←
                routine:␣"
253            << (t1_stop−t1_go)/CLOCKS_PER_SEC∗1000.0 << "␣←
                ms\n";
254        cout << "␣␣−␣Functions:␣" << ←
                (t2_stop−t2_go)/CLOCKS_PER_SEC∗1000.0 << "␣ms\n";
255
256        // Return successfully
257        return 0;
258    }
```