

Procedural Rendering with Ray Marching

Christian A. Robles
University of Southern California
Los Angeles, USA
roblesch@usc.edu

Samuel Yin
University of Southern California
Los Angeles, USA
slyin@usc.edu

Hoseung Lee
University of Southern California
Los Angeles, USA
hoseungl@usc.edu

Vansh Dhar
University of Southern California
Los Angeles, USA
vdhar@usc.edu

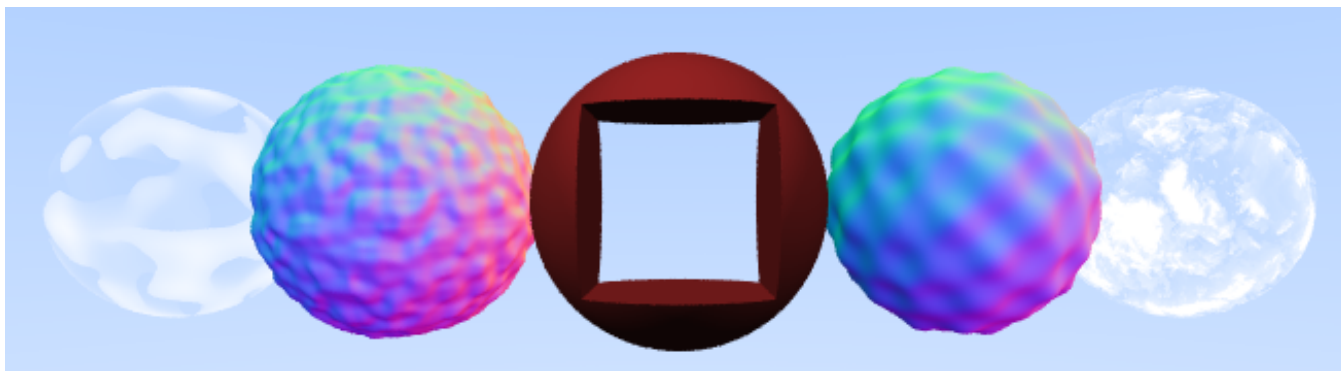


Figure 1: Procedural Perturbations, Constructive Solid Geometry, Marched Depth Clouds

ABSTRACT

Ray tracing requires evaluation of explicit ray-surface intersections. Complex shapes can be difficult to model as they must be broken down into an aggregate of many primitives, and ray tracing scales poorly with the number of objects in the scene. Ray marching allows us to model complex shapes more efficiently by implicitly identifying intersections with surfaces through the use of a Signed Distance Function. With ray marching, we can efficiently model and render complex shapes, create new shapes with constructive solid geometry, and render procedural materials and surfaces.

1 INTRODUCTION

Let a scene consist of a set of surfaces in a 3d world coordinate space with a camera at some position. We will generate an image by sending rays through each pixel coordinate in the camera's viewport, determining if an intersection occurs, and evaluating a material at that intersection.

Consider a ray tracing renderer. We determine the color of each pixel by comparing this ray to each object in the scene and determining if an intersection occurs, then evaluating material properties. Consider a simple sphere primitive. Ray-sphere intersection requires solving $|P(t) - C|^2 = r^2$ for some ray $P(t) = A + tb$ and some sphere with center C and radius r . We must additionally decide which intersection point is nearest and in front of the camera, if we are viewing from the inside or outside of the sphere, and calculate a surface normal to calculate the appropriate material color at that intersection. Even for a relatively simple primitive

such as a sphere, determining a ray-surface intersection for every surface in a scene is a computationally expensive task.

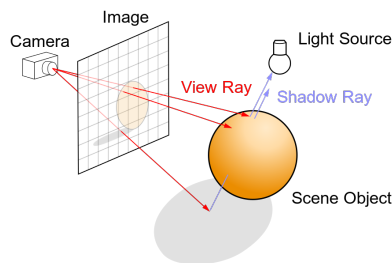


Figure 2: Evaluating a scene with ray tracing, via Wikimedia Commons. (<https://w.wiki/55Pu>).

In ray tracing, we are limited to only those surface primitives where we can explicitly calculate a ray-surface intersection. Planar surfaces can be solved with a ray-plane intersection and inside-outside test. Various shapes and orientations of surfaces can be derived by model transformations of translations rotations and scales. Higher complexity surfaces that would be otherwise too difficult to evaluate directly can be modeled by the composition of many planar primitives, but performance scales poorly with the number of surfaces in the scene. Complex materials additionally require even more rays to be cast for lighting, scattering, and transmission. Effects like soft shadows and anti-aliasing further degrade performance, requiring a polynomial total number of rays to be cast to determine the color of an individual pixel.

Consider a ray marching renderer. We again determine the color of each pixel by evaluating the scene in the direction of a camera ray. Rather than evaluating the intersection between this ray and any surface in the scene, we instead let the scene provide a *Signed Distance Function* which evaluates some point $p(x, y, z)$ and returns the approximate distance d between p and the nearest surface in the scene. If the value is positive, we are d units away from the nearest surface. If it is negative, we are d units inside the surface. If d is close to zero, then we can say we are intersecting the nearest surface. Starting from $t = 0$, we will for each ray evaluate the distance to the closest surface in the scene d and march d units in the direction of the ray. We will evaluate the scene again and repeat this process until d is near-zero, or until the ray terminates due to a march limit or a maximum depth along the ray t_{max} .

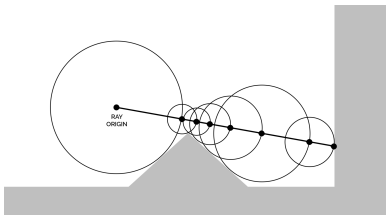


Figure 3: Marching along a ray using the SDF of the scene, via Wikimedia Commons. (<https://w.wiki/55Ps>).

In the ray marching model, ray-surface intersections no longer require explicit evaluation of intersection formulae. Instead, we implicitly model ray-surface intersections by marching down the ray until the distance to the scene is near-zero. With these implicit intersections, we are no longer constrained to simple surface primitives. We can render any surface that has a defined distance function. Additionally, we can model scenes of complex functions like fractals or noise functions by mapping their values to distance values. We can combine distance functions for multiple surface primitives by union or intersection, and interpolate between them with smoothing functions. We can model 3D materials by stepping through them and evaluating functions along each step. Surface normals no longer require explicit evaluation and can instead be calculated by the gradient of the distance function.

2 RENDERER IMPLEMENTATION

We implement our renderer from reference code provided by Peter Shirley's *Ray Tracing in One Weekend*. We use the template code Peter provides up to Chapter 4: *Rays, a Simple Camera, and Background*. We also use the Anti Aliasing techniques from Chapter 7, and the Positionable Camera from Chapter 11. We follow Peter's interface patterns for surfaces and materials, with some modifications. We rename *hittable* to *surface*, and *hittable_list* to *scene*. We discard all ray tracing functionality and replace it with our implementations of ray marching.

2.1 Interfaces

2.1.1 Surface. The *surface* class is an interface for intersectable objects in the scene. Any class that implements *surface* must accept

in its constructor a valid pointer to a *material*. It must provide an implementation for

```
double distance(const vec3& p)
```

The *distance* function returns the signed distance from some point $p(x, y, z)$ and the surface. If the result is negative, p is considered inside the surface.

2.1.2 Material. A class that implements the *material* interface must implement

```
vec3 color(ray &r, vec3 p, vec3 N,
           vector<light> lights)
```

The *color* function accepts an incoming ray, an intersection point, a surface normal, and a list of directional lights. It returns the color of the material evaluated by those parameters.

2.1.3 Scene. The *scene* class manages objects and lights in a scene, and maintains internal lists of pointers to instances of these classes. It provides the following key functions:

```
bool near_zero(double d)
```

Returns *true* if d is smaller than machine-epsilon: the difference between 1.0 and the next value representable by the floating-point type.

```
double distance_estimator(vec3 p)
```

Iterates over the list of scene surfaces and returns the distance from some point $p(x, y, z)$ and the surface nearest that point.

```
vec3 normal(vec3 &p)
```

Returns the "surface normal" at some point $p(x, y, z)$. We take a small step along each of the x, y, z axes and return the unit length vector in the direction of the gradient of the distance function.

```
bool march(ray& r, hit_record& rec)
```

Beginning at $t = 0$, *march* evaluates the distance from the point along the ray r at t to the nearest surface in the scene. If the distance is near zero, *march* will terminate and return true, and store the information about the intersection and the material of the surface intersected on the hit record. If the distance is not near zero, it will increase t by d with a minimum step size of 0.001 and evaluate the distance estimator again. Will terminate if a surface is hit or if $t \geq t_{max}$.

```
vec3 ray_color(ray& r)
```

Evaluates a ray by performing marching. If a surface is hit, queries the material properties for the color evaluated at the point of intersection. Otherwise, returns the background color of the scene.

2.2 Surface Primitives

The following classes provide implementations for *surface*. In each formulation, p refers to some point $p(x, y, z)$ along a ray. In our notation $\|v\|$ denotes the Euclidean norm for some vector v , and $|v|$ denotes a component-wise absolute value for some vector v . *Min* and *Max* operations are also component-wise.

2.2.1 *Sphere*. An instance of *sphere* is specified by a center C and a radius r .

$$distance = ||(p - C)|| - r$$

2.2.2 *Box*. An instance of *box* is specified by a vector B representing the distance from the center to the edges of the box -

$$q = |p| - B$$

$$distance = ||\max(q, 0)|| + \min(\max(q_x, \max(q_y, q_z)), 0)$$

2.2.3 *Equilateral Triangular Prism*. An instance of a triangular prism is specified by a center C , prism length L , and triangle height H .

$$q = |p - C|$$

$$distance = \max(q_z - H_y, \max(q_x * 0.866 + p_y/2, -p_y) - H_x/2)$$

2.2.4 *Cylinder*. An instance of a cylinder is specified by a center C , height H , and radius r .

$$q = p - C$$

$$d_x = |length(q_{xz}) - r|$$

$$d_y = |length(q_y) - H|$$

$$distance = \min(\max(d_x, d_y), 0) + ||\max(d, 0)||$$

2.2.5 *Pyramid*. An instance of a pyramid with a constant square 1×1 base is specified by a center C and height H .

$$a = H^2 + 1/4$$

$$p = p - C$$

$$p_{xz} = |p_{xz}|$$

$$p_z > p_x ? p_{xz} = p_{zx}$$

$$p_{xz} = p_{xz} - 1/2$$

$$b = (p_z, H * p_y - p_x/2, H * p_x + p_y/2)$$

$$c = \max(-b_x, 0)$$

$$d = clamp((b_y - p_z/2)/(a + 1/4), 0, 1)$$

$$e = a * (b_x + s)^2 + b_y^2$$

$$f = a * (b_x + d/2)^2 + (b_y - a * d)^2$$

$$\min(q_y, -q_x * a - q_y/2) > 0 ? g = 0 : g = \min(a, b)$$

$$distance = \sqrt{(g + b_z^2)/a} * sign(\max(b_z, -p_y))$$

2.2.6 *Infinite Cylinder*. An instance of an infinite cylinder is specified by a center C and radius R . Its distance function is identical to a cylinder without comparing the y coordinate.



Figure 4: Triangular Prism, Cylinder, Pyramid, Infinite Cylinder

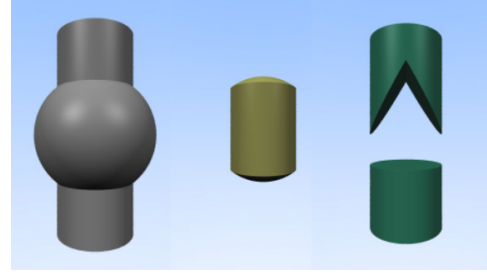


Figure 5: Union, Intersection, Difference

2.3 Materials

2.3.1 *Diffuse*. The *diffuse* class implements the *color* function to return the color C as

$$C = (Ks \sum_L (le(R \cdot E)^s)) + (Kd \sum_L (le(N \cdot L))) + (Ka * la)$$

2.3.2 *Clouds*. We implement the materials *perlin_cloud_2d*, *perlin_cloud_3d*, *gardner_cloud_2d* and *gardner_cloud_3d* by mapping surface coordinates to a noise function and attenuating output color by march depth. Described in section 4.

2.3.3 *Normals*. Maps a normal N to RGB as

$$C = (N + 1)/2$$

2.3.4 *Flat*. Returns a constant color.

3 CONSTRUCTIVE SOLID GEOMETRY

Constructive Solid Geometry is the technique of combining primitives to form a more complex object using Boolean operators such as *union*, *difference*, and *intersection*. Ray marching makes CSG easy. When a signed distance function returns a positive value, it means it is outside of the object. A negative SDF means we are inside of the object. When can use this property to implement the Boolean operations with combination and negation of the SDFs for two or more surfaces.

3.0.1 *Union*. The union of two objects is computed simply by the minimum of the distances of two objects. Since we are already returning the surface nearest to the point in our distance estimator function, this operation is already handled by the ray marching algorithm.

3.0.2 *Intersection*. The intersection of two objects is the area in which both objects converge. To get the intersection of two objects, distance to both objects should be less than or equal to zero. This creates an object that only exists at the intersection of the two objects.

3.0.3 *Difference*. Finally, the difference of two objects is the maximum of the distance functions with one function negated. Flipping the sign inverts the object so everything that was considered inside an object is now outside and vice versa. The object now only exists if it is inside the non-inverted object and outside of the inverted object, resulting in one area "cutting" into the other.

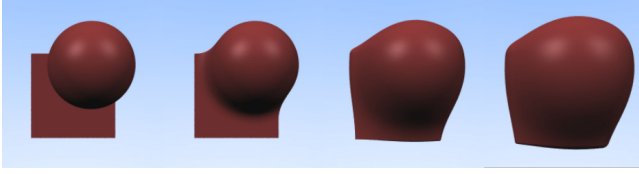


Figure 6: Smooth union with $k = 0, 0.2, 0.6, 0.9$

3.0.4 *Smoothing*. We can blend the intersection of two objects by interpolating where the objects approach each other. We can apply smoothing to union as

```
h = max(k - abs(d1 - d2), 0.0);
return min(d1, d2) - h * h * 0.25/k;
```

Where k is some value between 0 and 1 and controls the amount of smoothing. For intersection and difference, we apply the same logic as before: using max instead of min and flipping the sign for the difference.

4 PROCEDURAL TEXTURES

Procedural generation of 3D clouds is a complex challenge in rendering. Clouds in real-time applications are often drawn by mapping pre-rendered 2D textures at distances or on volumes not interactable by the camera. However, these approaches fall apart when viewed at shallow angles or positioned near the camera. We will examine how ray marching can be leveraged to generate 3D textures on marchable depth surfaces to approximate the appearance of clouds.

We will examine two noise functions for generating cloud textures. The first is proposed by Geoffrey Y. Gardner in his paper *Visual Simulation of Clouds*. The second is the famous Perlin Noise, originally proposed by Ken Perlin in his paper *An image synthesizer*.

4.1 Gardner Noise

In *Visual Simulation of Clouds*, Gardner proposes a texturing function that approximates the appearance of clouds. We refer to this function as *Gardner Noise*. It models textures as a product of sums of $3 < n < 8$ sine waves

$$T(X, Y, Z) = k \sum_{i=1}^n [C_i \sin(FX_i x + PX_i)] \times \sum_{i=1}^n [C_i \sin(FY_i y + PY_i) + T_0]$$

Where frequencies and coefficients are chosen by

$$\begin{aligned} FX_{i+1} &= 2FX_i \\ FY_{i+1} &= 2FY_i \\ C_{i+1} &= .707C_i \end{aligned}$$

And PX_i and PY_i are determined by

$$\begin{aligned} PX_i &= \pi/2 \sin(.5FY_{i-1}Y) + \pi \sin(FX_i Z/2) & \text{for } i > 1 \\ PY_i &= \pi/2 \sin(.5FX_{i-1}iX) + \pi \sin(FX_i Z/2) & \text{for } i > 1 \end{aligned}$$

T_0 is a parameter controlling contrast, and k is computed such that the maximum of $T(X, Y, Z) \approx 1$.

4.2 Perlin Noise

In his paper *An image synthesizer*, Ken Perlin describes a noise function that accepts a n -vector and returns a value with the qualities of statistical invariance under rotation and translation and a narrow bandpass limit in frequency. We use 3-noise, which is constructed by

1. Construct a 3 dimensional grid of size x by y by z .
2. Assign to each grid intersection a pseudo-random unit length gradient vector $v_{x_i, y_i}(x, y, z)$
3. For some input $p(x, y, z)$, determine which grid cell contains p . Identify the corners of that cell and calculate an offset vector from each corner to p .
4. For each corner, calculate the dot product of its offset vector and its gradient vector.
5. Interpolate the grid corner dot products at p and return this value as the result.

4.3 Procedural Clouds

Using noise to generate textures is simple. We evaluate a ray-surface intersection and pass the point $p(x, y, z)$ to one of these functions and use the $[0, 1.0]$ result as a texture value T . Gardner phase and coefficient parameters can be tuned to modify the size and variance of the cloud-like textures, and Perlin noise can be tuned to map input coordinates to larger or smaller gradient grids for fuzzier or granular textures.

We can additionally march through surfaces to generate more convincing 3D textures. We accept as a material parameter the approximate depth of a surface and use this value to determine a marching step bound and step size. We use a cutoff threshold to march past low texture values and attenuate the final result by a combination of step size and step count.

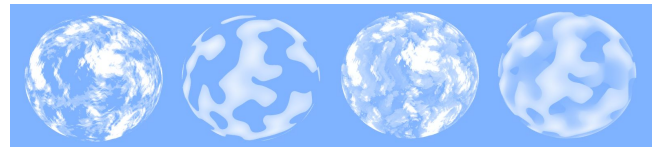


Figure 7: Surface Mapped and Depth Marched Gardner & Perlin Clouds

5 PROCEDURAL SURFACES

In section 3 we saw how we can alter surfaces by altering the SDF. We can create a new class of surfaces we call procedural surfaces by combining this technique with our noise functions.

Procedural surfaces are analogous to the bump mapped textures. Bump maps use lighting effects to provide the appearance of surface variation, but do not actually affect surface geometry. With procedural surfaces, we can create more convincing effects by altering the surface directly. Additionally, we can render fractals and other complex functions by mapping values to an SDF over some defined set of inputs.

5.1 Displacement Surfaces

Sphere perturbations are simple - for some displacement function D , we add its value to the SDF of a sphere

$$distance = (p - C) - r + D$$

We implement two displacement surfaces. *Perlin Sphere* uses Perlin Noise (4.1.2) to evaluate $D(x, y, z)$. *Perturbed Sphere* implements D as a sum of sine waves

$$D(x, y, z) = c * \sin(\varphi + x) * \sin(\varphi + y) * \sin(\varphi + z)$$

Where c is some coefficient that modulates the intensity of displacement and φ is some constant phase shift.

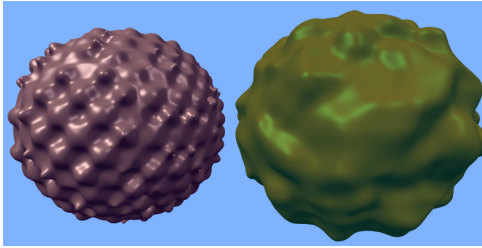


Figure 8: Diffuse Perturbed and Perlin Spheres

5.2 Fractals

Fractals are infinitely complex mathematical shapes that recurse as you zoom into an area of a fractal. These infinitely complex fractals could not be rendered in ray tracing due to having no analytical intersection. With ray marching, we can use distance functions to render visualizations of complex relations such as the Mandelbulb fractal. Our implementation of a Mandelbulb is formulated as

$$r = \sqrt{x^2 + y^2 + z^2}$$

$$\phi = \arctan \frac{y}{x} = \arg(x + yi)$$

$$\theta = \arctan \frac{\sqrt{x^2 + y^2}}{z} = \arccos \frac{z}{r}$$

$$\mathbf{v}^n := r^n \langle \sin(f(\theta, \phi)) \cos(g(\theta, \phi)), \sin(f(\theta, \phi)) \sin(g(\theta, \phi)), \cos(f(\theta, \phi)) \rangle$$

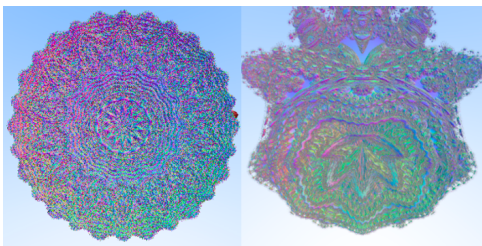


Figure 9: Mandelbulbs at different view angles

For complex shapes, it is important to supersample rays to reduce aliasing artifacts introduced by high variance within pixel coordinates.

6 PARALLELIZATION

A ray marching rendering system requires more computational resources than ray tracing due to the simple fact that in ray tracing, computation is done only at the point of object intersection. In ray marching, distance functions are evaluated for every iteration of the march. Thus, to improve the execution time of our renderer, we utilize the power of parallel computing with the help of CUDA programming and NVIDIA GPU.

6.1 Platform Used

Hardware specification: The CUDA program was tested on a 64-bit operating system with 16 GB of Ram and an AMD Ryzen 9 5900HX processor, along with a Nvidia GeForce RTX 3070 GPU (Laptop version).

Software Specification: The hardware ran on Windows 11 OS and we used Microsoft Visual Studio as our programming IDE to develop the CUDA program. We also used Nvidia Nsight Compute to profile our CUDA program to better understand the resource utilization of the GPU while the program was running.

6.2 CUDA Programming Overview

In GPU programming, there are several parameters on both the software side and the hardware side that affect performance. Whenever a part of the program is outsourced to a GPU to be computed, On the software side: the total number of threads/ processors used are represented by a "Grid", which in our case will represent the image being rendered, therefore, the total number of threads will be equal to the number of pixels in our image. These threads are further divided into blocks which are independent units of execution with no communication possible between different blocks.

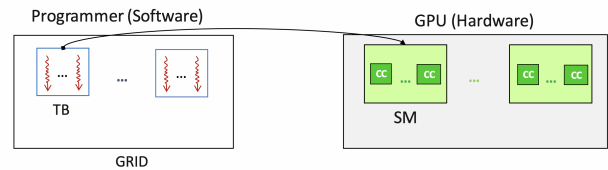


Figure 10: Diagram depicting software and hardware abstraction for GPU programming

Now on the hardware side: The total number of cores/processors in the GPU are divided between Streaming Multiprocessors (SMs) such that each block can be executed only on a single SM with the possibility of several blocks being executed concurrently on the same SM if the hardware resource allows it. Thus, fine-tuning parameters like block size for an image of specific resolution being rendered based on the resource utilization of the hardware can lead to a significant increase in SpeedUp.

$$SpeedUp = S/P \tag{1}$$

here, S: is the Program Serial Execution Time and P: is the Program Parallel Execution time

Resolution	Block Size	Number of Samples	Parallel Time	Serial Time	Speed Up	Branch Efficiency	A/T Occupancy
256x256	8x8	8	0.027s	29.483s	1091.963	99.27	83.42/87.5
256x256	16x16	8	0.031s	29.483s	951.065	99.27	80.14/87.5
256x256	32x8	8	0.042s	29.483s	701.976	97.54	76.59/83.33
256x256	128x1	8	0.038s	29.483s	775.868	99.27	77.37/83.33
256x256	256x1	8	0.033s	29.483s	893.424	98.65	79.28/87.5
512x512	8x8	8	0.073s	117.817s	1613.932	99.38	82.74/87.5
512x512	16x16	8	0.081s	117.817s	1454.531	99.27	78.94/87.5
512x512	32x8	8	0.096s	117.817s	1227.260	99.27	72.41/83.33
512x512	128x1	8	0.092s	117.817s	1280.620	98.65	74.53/87.5
512x512	256x1	8	0.093s	117.817s	1266.849	98.65	74.68/87.5
1024x1024	8x8	8	0.191s	476.936s	2497.047	99.38	83.65/87.5
1024x1024	16x16	8	0.209s	476.936s	2281.990	98.65	81.12/87.5
1024x1024	32x8	8	0.243s	476.936s	1962.700	99.38	74.64/83.33
1024x1024	128x1	8	0.241s	476.936s	1978.988	98.65	75.23/83.33
1024x1024	256x1	8	0.244s	476.936s	1954.656	98.65	75.87/83.33

Table 1: In the above table, we list the different parameter configurations of our CUDA program and the SpeedUp achieved for those configurations. We also list the Branch Efficiency and Achieved/Theoretical Occupancy for each configuration.

6.3 Parallel Algorithm

This section describes the algorithm/pseudo code run by each thread of the program.

Algorithm 1: Algorithm executed by each Thread/Pixel

```

i ← ThreadId.x + BlockDim.x * BlockIdx.x
j ← ThreadId.y + BlockDim.y * BlockIdx.y
Ensure: i < Image_Width and j < Image_Height
N ← No_of_Samples
Colorij ← 0
while N ≠ 0 do
  i+ = Random(−1, 1)
  j+ = Random(−1, 1)
  ray = Compute_Ray(Camera_Info, i, j)
  Colorij+ = Compute_Color(3D_Scene, ray)
  N ← N − 1
end while
FrameBuffer[j, i] = Colorij/No_of_Samples

```

6.4 Experiments and Results

As discussed in section 6.2, several parameters can affect the execution time of our rendering system. The focus of our experiments will be to understand the range of the parameters for which we can achieve maximum speedup. The main parameters of focus will be: Block size (Number of Threads in each block), Grid size (Image Resolution). To better understand Resource utilization of our GPU based on these parameters we will also look at Branch Efficiency and Achieved Occupancy.

Branch Efficiency is the ratio of similar work done by threads in a block to the total work done by threads in a block. (Note: Threads in a block run concurrently only if they do similar work, thus, this ratio should be high to achieve maximum performance). Achieved Occupancy is the ratio of number of active processors to the total

number of processors in a Streaming Multiprocessor. (Note: This ratio should be as close to theoretical occupancy as possible to achieve maximum performance).

Thus, for a range of parameter values the Program Parallel Execution Time, Program Serial Execution Time and SpeedUp can be seen in Table 1. By analysing the results in Table 1, we can see that SpeedUp scales with Image Resolution (with the highest SpeedUp achieved being $\approx 2500x$), for any Image Resolution Block Size of '8x8' yield highest SpeedUp. Branch Efficiency remains high for any parameter configuration, therefore, divergent behaviour in our program is bare minimum. Although the Theoretical Occupancy is not close to 100% due to the register utilization of our CUDA program, the achieved occupancy for various parameter configuration remains close to theoretical occupancy, thus extracting high performance benefit from the GPU.

7 CONCLUSION

In this project we developed a Ray Marching rendering system and demonstrated the unique properties of the Signed Distance Function with Constructive Solid Geometry, Procedural Materials, Displacement Surfaces and Fractals. Finally, we utilized the power of parallel computing via Nvidia GPU and CUDA programming to achieve a Speed Up factor of 2500.

REFERENCES

- [1] Peter Shirley. Ray Tracing In One Weekend. <https://raytracing.github.io>
- [2] Geoffrey Y. Gardner. Visual Simulation of Clouds. SIGGRAPH '85
- [3] Ken Perlin. 1985. An Image Synthesizer. SIGGRAPH '85
- [4] <https://developer.nvidia.com/blog/accelerated-ray-tracing-cuda>
- [5] <https://www.skytopia.com/project/fractal/mandelbulb.html>
- [6] <https://iquilezles.org/articles/distfunctions>
- [7] <http://jamie-wong.com/2016/07/15/ray-marching-signed-distance-functions>
- [8] <https://michaelwalczyk.com/blog-ray-marching.html>