



Real-time Rendering of 3D Fractal Geometry

Final Year Dissertation

21/04/2022

by

*Solomon Baarda
Meng Software Engineering
Heriot-Watt University*

Supervisor: Benjamin Kenwright

Second Reader: Ali Muzaffar

Declaration

I, Solomon Baarda confirm that this work submitted for assessment is my own and is expressed in my own words. Any uses made within it of the works of other authors in any form (e.g., ideas, equations, figures, text, tables, programs) are properly acknowledged at any point of their use. A list of the references employed is included.

Signed:

A handwritten signature in black ink that reads "Solomon Baarda". The signature is written in a cursive style with a clear distinction between the first name and the last name.

Date: 21/04/2022

Abstract

A *fractal* is a pattern which remains detailed at any arbitrary scale. As such, fractals appear everywhere in the natural world from the structure of clouds to the roughness of mountains to the edges of our coastlines. Many disciplines require software capable of displaying fractals, since these patterns can be used when modelling the growth of cancerous cells, when describing losses and gains in the stock market, and when viewing complex data sets. Most of the currently available 3D fractal viewing applications are designed for outputting images or videos of high visual quality, using realistic lighting and shading techniques, while there are few applications which can accurately display changes to 3D fractals in real time. This project developed an application capable of updating and rendering 3D fractal geometry in real-time, computing each frame using brute force through the utilisation of the parallel nature of the current generation of graphical processors.

Acknowledgements

I'd like to take a moment to express my gratitude to those people without whom this project would not have been successful. Firstly, I'd like to thank my supervisor *Ben Kenwright* for his continued support over the duration of this project. His advice and expertise in this field helped guide me throughout the duration of my project. I'd also like to thank *Ali Muzaffar* for his valuable feedback on my initial research report, which encouraged me to complete more thorough research in several areas. Finally, I'd like to thank *Chris Hulme*, *Callum Scott*, and *Fraser Orr*, for allowing me to use their computers to gather results for my project, without which this work would be less significant.

Tables

Table of Contents

Declaration.....	i
Abstract	ii
Acknowledgements	iii
Tables	iv
Table of Contents.....	iv
Table of Figures.....	vii
Table of Tables.....	viii
Table of Equations.....	viii
Table of Graphs.....	viii
Definitions and Abbreviations	ix
Common Definitions.....	ix
Common Abbreviations	ix
1 Introduction	1
1.1 Project Description.....	1
1.2 Aims & Objectives	2
1.3 Scope.....	3
1.4 Additional Resources.....	4
1.5 Document Structure.....	4
2 Literature Review	5
2.1 Fractals	5
2.1.1 Sierpiński Tetrahedron.....	5
2.1.2 Sierpiński Cube.....	6
2.1.3 Mandelbulb	7
2.1.4 3D Fractals Summary.....	9
2.2 Fractal Rendering Methods.....	10
2.2.1 Rasterization	10
2.2.2 Ray Tracing	11
2.2.3 Ray Marching.....	11
2.3 Surface Shading.....	16
2.3.1 Phong Reflection Model.....	16
2.4 GPU Computing.....	17
2.4.1 CUDA.....	18
2.4.2 OpenCL.....	18

2.4.3	GPU Computing Summary.....	19
2.5	Review of Existing Applications.....	19
2.5.1	Fragmentarium.....	20
2.5.2	Synthclipse	20
2.5.3	Mandelbulb3D.....	20
2.5.4	Mandelbulber.....	21
2.5.5	FractalLab	21
2.5.6	Existing Applications Summary.....	22
2.6	Literature Review Summary.....	22
3	Development.....	23
3.1	Development Environment.....	23
3.2	Application Design.....	23
3.2.1	High Level Overview	23
3.2.2	Application Design Principles	24
3.2.3	Kernel Design Principles	26
3.3	Technologies	28
3.4	Development Strategy.....	29
3.5	Documentation.....	32
3.6	Interface Design.....	32
4	Evaluation	33
4.1	Unit Testing	33
4.2	Benchmarking Framework	34
4.3	Fractals Implemented	36
4.4	Evaluation of Application.....	37
4.4.1	Optimisations.....	37
4.4.2	Efficiency.....	42
4.4.3	Unmeasurable Features.....	43
4.4.4	Measurable Features.....	44
4.4.5	Application Scalability.....	48
4.5	Evaluation of Requirements Specification.....	51
4.6	Evaluation of Aims and Objectives.....	51
4.7	Comparison with Literature Results.....	53
4.8	Future Work	56
4.8.1	Features.....	56
4.8.2	Algorithmic Optimisations.....	56
4.8.3	Design Optimisations.....	57
4.8.4	Fractals.....	58

4.8.5	Evaluation.....	58
5	Conclusion.....	59
5.1	Achievements.....	59
5.2	Limitations	59
5.3	Future Work.....	60
5.4	Final Thoughts.....	60
6	Appendix.....	61
6.1	Experimentation Renders.....	61
6.2	Sphere SDF Example	62
6.3	Smooth SDF Combinations.....	62
6.4	SDF Surface Normal Calculations	63
6.5	Existing Solution Renders.....	64
6.5.1	Fragmentarium Renders	64
6.5.2	Synthclipse Renders	65
6.5.3	Mandelbulb3D Renders	66
6.5.4	Mandelbulber Renders.....	67
6.5.5	FractalLab Renders.....	67
6.6	Application Class Diagrams.....	68
6.7	Comparison of -cl-fast-relaxed-math Optimisation	69
6.8	Comparison of Glow with Bounding Volume.....	70
6.9	Requirements Specification.....	71
6.10	Use Cases	73
6.11	Project Plan.....	74
6.11.1	Professional, Legal, Ethical & Social Issues.....	74
6.11.2	Risk Analysis.....	74
6.11.3	Project Timeline.....	76
7	References.....	79

Table of Figures

Figure 2.1.1 Sierpiński triangle (left) [8] and tetrahedron (right) [8] both of recursive depth 5...	6
Figure 2.1.2 Sierpiński carpet (left) [9] and cube (right) [10] both of recursive depth 4	6
Figure 2.1.3 Mandelbrot set overview (left), antenna (middle), and seahorse (right).....	8
Figure 2.1.4 Mandelbulb power of two (left), three (middle), and eight (right) [17]	9
Figure 2.2.1 Ray marching diagram.....	12
Figure 2.2.2 Ray marched sphere and box scene experiment union (left), subtraction (middle) and smooth union (right)	15
Figure 2.3.1 Phong reflection model [26].....	16
Figure 3.2.1 Application high level overview.....	24
Figure 4.4.1 Nvidia Nsight Systems output for Mandelbulb benchmark.....	42
Figure 4.4.2 Windows GPU usage for Mandelbulb benchmark.....	43
Figure 4.4.3 Basic shading (left) and Blinn-Phong shading (right)	46
Figure 4.4.4 Geometry glow applied to Mandelbulb.....	46
Figure 4.4.5 Hard shadows (left) and soft shadows (right)	47
Figure 6.1.1 Render of the Mandel bulb fractal (left) and cross section (right) using equation from [49]	61
Figure 6.1.2 Surface normal visualised (left) and phong shading experiment (right)	61
Figure 6.2.1 Sphere SDF diagram.....	62
Figure 6.5.1 Fragmentarium interface when viewing the Mandelbulb fractal	64
Figure 6.5.2 A recursive scene (left) and Mandel bulb (right).....	64
Figure 6.5.3 Tree fractal.....	65
Figure 6.5.4 Synthclipse interface when viewing an ice cube shader.....	65
Figure 6.5.5 Mandelbulb3D interface when viewing the Mandelbulb fractal	66
Figure 6.5.6 Mandel bulb fractal with natural looking colouring.....	66
Figure 6.5.7 Mandelbulber interface when viewing the Mandelbulb fractal	67
Figure 6.5.8 FractalLab interface when viewing the Mandelbulb fractal.....	67
Figure 6.6.1 Application public methods	68
Figure 6.6.2 Dependency graph of the application	68
Figure 6.6.3 Dependency graph for example scene hello_world.cl.....	68
Figure 6.7.1 Mandelbulb.cl scene with -cl-fast-relaxed-math enabled	69
Figure 6.7.2 Mandelbulb.cl with -cl-fast-relaxed-math disabled	69
Figure 6.8.1 Mandelbulb scene with glow without bounding volume (left) and with bounding volume (right)	70
Figure 6.10.1 Use case diagram for the application	73
Figure 6.11.1 Initial Gantt chart for deliverable one (top) and deliverable two (bottom)	77
Figure 6.11.2 Revised Gantt chart for deliverable two	78

Table of Tables

Table 1.1.1 Common definitions.....	ix
Table 1.1.2 Common abbreviations	ix
Table 3.2.1 Class responsibilities	25
Table 3.2.2 Key scene methods to implement.....	26
Table 3.3.1 Packages used by the application.....	29
Table 4.2.1 Benchmark framework results	35
Table 4.2.2 Results calculated from benchmark.....	35
Table 4.4.1 Graphics cards used by the benchmarking systems.....	48
Table 4.7.1 Comparison matrix for literature results.....	55
Table 6.9.1 Functional requirement specification	71
Table 6.9.2 Non-functional requirement specification	73
Table 6.11.1 Risk rating matrix.....	75
Table 6.11.2 Risk analysis matrix	75

Table of Equations

Equation 2.1.i White and Nylander's formula for the nth power of a point in 3D space	8
Equation 2.1.ii Addition of two points in 3D space	9
Equation 2.2.i Equation for transforming a point in space	14
Equation 2.2.ii Equations for scaling a point down (left) and up (right)	14
Equation 2.2.iii Equations for union (left), subtraction (middle) and intersection (right) of two SDFs.....	15

Table of Graphs

Graph 4.4.1 Performance of optimisations in the Mandelbulb scene	39
Graph 4.4.2 Comparison of FPS for Planet scene with C and OpenCL C noise function	41
Graph 4.4.3 Comparison of percentage FPS increase for Planet scene with C and OpenCL C noise function.....	41
Graph 4.4.4 Computational Cost of Features in Mandelbulb Scene.....	45
Graph 4.4.5 Scene Performance on Devices of Varying Spec	49
Graph 4.4.6 Mandelbulb Scene Frame Time for Varying Resolutions	50

Definitions and Abbreviations

Common Definitions

Table 1.1.1 Common definitions

Word	Definition
Complex number	Number system containing an imaginary unit, defined as $i = \sqrt{-1}$
Convex polyhedron	3D equivalent of a regular 2D polygon
Euclidian geometry	A geometry containing basic objects like points and lines
Fractal	Recursively created never-ending pattern that is usually self-similar
Frame	One of many still images that make up a moving image
Geometry	Branch of mathematics concerned with the properties of space, distance, shape, size, and positions
Method overriding	Object-oriented programming technique which allows a subclass to change the implementation of a method that a parent class is providing
Object-oriented programming	Programming style that organises its software design around reusable objects
Polygon	2D shape with three or more sides
Polyhedrons	3D shape with six or more faces
Quaternion	Extension of the complex numbers, often used for storing position, translation, and scale values in 3D space
Ray	Line in 3D space with a beginning and an end
Regular polygon	Shape where all edges are the same length, and all angles between vertices are also equal
Render	Process of creating a still image using a computer
Vector	Mathematical object representing a position in space relative to another

Common Abbreviations

Table 1.1.2 Common abbreviations

Word	Abbreviation
CPU	Central Processing Unit
DE	Distance estimation function
FPS	Frames per second
GPU	Graphics Processing Unit
PC	Personal computer
SDF	Signed distance function

1 Introduction

1.1 Project Description

A *fractal* is a pattern which remains detailed at any arbitrary scale. As such, fractals appear everywhere in the natural world from the structure of clouds to the roughness of mountains to the edges of our coastlines. Benoit Mandelbrot, inventor of the concept of fractal geometry, famously wrote "Clouds are not spheres, mountains are not cones, coastlines are not circles, bark is not smooth, nor does lightning travel in a straight line" [1]. Fractal geometry describes the more non-uniform shapes found in nature, from the patterns created by the structure of our own DNA molecules to the shapes created by the formation of galaxies, and everything between.

While fractal patterns look aesthetically pleasing and have featured in popular pieces of artwork, they have also led to technical breakthroughs and have a range of uses in industry [1]–[4]. In medicine, fractals have been used to help distinguish between cancerous cells which grow abnormally, and healthy human blood vessels which typically grow in fractal patterns. Fractal patterns are a key design choice in the design of efficient Wi-Fi and cell phone antennas, as some of these patterns maximise the area of antenna that can transmit and receive signals. In computer science, fractal compression is an efficient method for compressing images and other files, which uses the fractal characteristic that parts of a file will resemble other parts of the same file. Fractals have been used to describe losses and gains in the stock market, and they have also been used visualise complex data sets such as Twitter posts, which can be visualised using branching trees. Some use cases require software capable of displaying a visual depiction of the fractal being modelled, such as in the discipline of fluid mechanics, when viewing complex turbulence flows and the structure of porous materials.

Most of the currently available fractal viewing applications are designed for outputting images or videos of high visual quality, using realistic lighting and shading techniques, while there are few applications which can accurately display changes to 3D fractals in real time. This is likely

due to performance limitations of commercially available Central Processing Units (CPUs) and Graphics Processing Units (GPUs), which for many years were not powerful enough for computing images of 3D fractals in real-time. However, over the last couple years, the parallel degree of commercially available GPUs has increased significantly and Nvidia's current generation may be powerful enough to allow us to render fractals in real time using brute force, through the utilisation of the parallel nature of the current generation of graphical processors. This project set out to investigate if the current generation of Nvidia's GPUs are powerful enough to render 3D fractal geometry in real time, using brute force.

1.2 Aims & Objectives

The aim of this project was to develop an application which can update and render 3D fractal geometry in real-time. The render used common optical effects including the Blinn-Phong lighting and reflection model, hard and soft shadows, and geometry glow. In addition, scenes can be created and viewed using the application, and this process is as straight forward as possible.

Listed below are the key objectives that this project set out to achieve. These objectives helped guide the project in the correct direction as to achieve its aim.

Objective 1: Research topic

Background research of the project area to gain a better understanding of the chosen topic and the scope of the project. Once this first stage was completed, the research was kept up to date and any significant developments in the project research area were explored.

Objective 2: Investigate existing solutions

Several relevant existing solutions exist, and an analysis of their strengths and flaws helped guide the project in the correct direction. This information was kept up to date and any relevant newly released applications were added and reviewed.

Objective 3: Core functionality

Implement the core functionality of a non-real-time 3D fractal renderer. This objective formed the safe core of the project, and the following objectives built upon this.

Objective 4: Additional functionality

This objective involved adding additional functionality to the renderer, such as making the application capable of real-time rendering, adding a game-loop, adding a controllable camera, making scenes dynamic, and adding optical effects and lighting. The scope of this objective was increased and decreased as necessary, and several additional stretch goals were included in the requirements specification in appendix 6.9.

Objective 5: Evaluation

Benchmark the performance of the application across various systems and evaluate how successfully the project aim was achieved.

Objective 6: Create user documentation

Create documentation to assist users and developers of the application.

These objectives formed the main tasks to be completed during the duration of this project, and the requirements specification in appendix 6.9 and the Gantt chart in appendix 6.11.3 were structured around these. It was necessary to complete some of the objectives in the order they are specified, as they build upon previous objectives. Objectives one and two ran for the entire duration of the project, to ensure that the project stayed up to date with current developments. Objectives three to six, however, relate to specific functionality to be implemented, and were completed in order. These objectives have been created bearing the SMART properties in mind. SMART stands for Specific, Measurable, Achievable, Realistic and Time constrained.

1.3 Scope

The scope of the project was carefully considered, and several stretch goals were included in the requirements specification to be implemented if good progress was made. Objective 4 left large

amounts of flexibility in the scope of the project as the objective could be extended or cut back depending on time constraints.

Common rendering optimisations such as the various types of culling, level of detail, and billboarding were out of scope of the project and have not been discussed. This project instead focussed on implementing the core functionality of a brute force 3D fractal renderer as a proof of concept, and building upon it from there.

1.4 Additional Resources

All project files including source code, builds of the application, screenshots, and some of the papers cited in this report can be found in the project [GitHub repository](#) [5]. Additionally, user and developer documentation is located in the application [documentation website](#) [6].

Some images in this report may be blurry due to image compression by Microsoft Word. A gallery containing screenshots and videos of the application can be accessed through the project GitHub repository readme file.

1.5 Document Structure

Continuing from section 1, section 2 discusses relevant background literature for the project. Section 3 describes the work completed during the duration of the project and discusses relevant technical design decisions. Section 4 evaluates the application created and the work completed during the duration of the project. Section 5 concludes the document by evaluating the significance of the work completed and discusses future developments for the application.

Appendix 6 contains additional equations, graphs, and images referenced throughout the report, along with the original project requirements specification, risk analysis, and project plan. Finally, section 7 contains references to all sources cited in this document.

2 Literature Review

This literature review contains a breakdown of information relevant for understanding the complexity of this project. While this review contains explanations of all relevant key concepts, basic knowledge of recursion, vector maths, and complex numbers is assumed. **Throughout this section, the notation $a \times b$ has been used to represent a multiplied by b , and not the cross product of a and b .**

This review is split into several sections, first the theory of 2D fractals and their 3D counterparts will be discussed. Then the key concepts of ray marching, the chosen method of rendering fractals, will be outlined. This is followed by a brief introduction to surface shading techniques, specifically the Phong reflection model. Then the core concepts of GPU parallel programming will be broken down, and an analysis of the available parallel programming frameworks will be completed. Finally, the literature review finishes with an analysis of relevant existing solutions.

2.1 Fractals

As discussed in the introduction, a fractal is a pattern that remains detailed at any scale. This concept defines fractal geometry, which describes up the more non-uniform shapes found in nature, like clouds, mountains, and coastlines. There exist several ways of artificially creating a fractal, from manually defining a simple repeating pattern to studying the convergence of equations. This section will discuss several common 3D fractals (and their 2D counterparts) and the methods used for creating them.

2.1.1 Sierpiński Tetrahedron

The Sierpiński tetrahedron, is a 3D representation of the famous 2D Sierpiński triangle fractal, named after the Polish mathematician Wacław Sierpiński [7]. The Sierpiński triangle is one of the most simple and elegant fractals and has been a popular decorative pattern for centuries. This pattern is created by recursively each splitting each solid equilateral triangle into four smaller equilateral triangles and removing the middle one. Theoretically, these steps are

repeated forever, but in practice when creating this fractal using a computer, some maximum depth must be specified as computers only have finite memory.

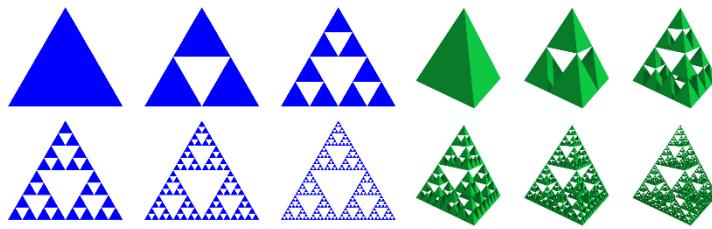


Figure 2.1.1 Sierpiński triangle (left) [8] and tetrahedron (right) [8] both of recursive depth 5

As the recursive depth of the fractal increases, so does the number of objects (either triangles or tetrahedrons) in the scene. The total number of objects in the Sierpiński triangle increases by a factor of three each iteration and the tetrahedron by a factor of four. This is the limiting factor when rendering the Sierpiński tetrahedron, as computer memory is finite and can only store limited number of objects.

2.1.2 Sierpiński Cube

The Sierpiński cube, also known as the Menger cube or Menger Sponge, is another 3D representation of one of Sierpiński's 2D fractals. This 2D fractal is known as the Sierpiński carpet, which follows very similar recursive rules to the Sierpiński triangle but uses squares instead of triangles.

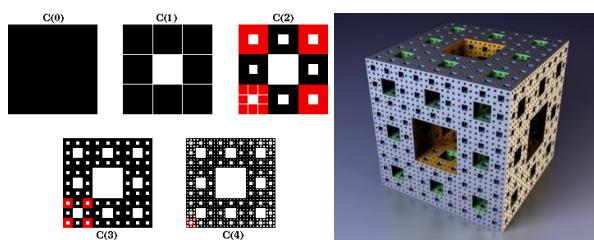


Figure 2.1.2 Sierpiński carpet (left) [9] and cube (right) [10] both of recursive depth 4

The number of objects required to create these fractals at various recursive depths increases similarly to the Sierpiński triangle and tetrahedron, but at a larger rate. The Sierpiński carpet increases by a factor of eight each iteration while the Sierpiński cube by a factor of 20. A Sierpiński cube at recursive depth n is made up of 20^n smaller cubes. As with the Sierpiński

tetrahedron, the limiting factor when trying to create this fractal is the exponential growth of the number of objects in the scene as the recursive depth increases.

In addition to the Sierpiński tetrahedron and cube, Sierpiński variations of all other convex polyhedrons exist. A convex polyhedron is the 3D equivalent of a 2D regular polygon. While there are an infinite number of regular polygons, there are only five possible convex polyhedrons. These are the tetrahedron, cube, octahedron, dodecahedron, and icosahedron. These polyhedrons are known as the platonic solids, and their fractal counterparts are called the platonic solid fractals. [This Wikipedia page](#) contains a concise list of these shapes [11].

2.1.3 Mandelbulb

While the platonic solid fractals are created by making many copies of a primitive shape, another method for creating fractals is to plot the convergence of values for certain mathematical equations. This can instead generate much more “natural” looking fractal patterns. One of the first fractals of this type to be discovered was the Mandelbrot set, discovered by Adrien Douady and named after Benoit Mandelbrot, the inventor of the concept of fractal geometry. The Mandelbrot fractal is defined as the set of complex numbers c for which the iteration from $z = 0$ in the equation $z_{n+1} = z_n^2 + c$ remains bounded (does not diverge to infinity) [12]. This definition is commonly written in the form $f_c(z) = z^2 + c$, where the value of c is varied. While this equation is staggeringly simple, when plotted on a graph using the value of c (a complex number in the form $x + iy$) with x and y corresponding to the position on the x and y axis, a colour can be assigned based on how quickly that value of c tends towards infinity. This results in a beautiful shape that when zoomed in on, shows more fractal shapes within. [This Wikipedia page](#) contains a fantastic image gallery showing common shapes found within the Mandelbrot set [13], a few of which are shown in Figure 2.1.3.

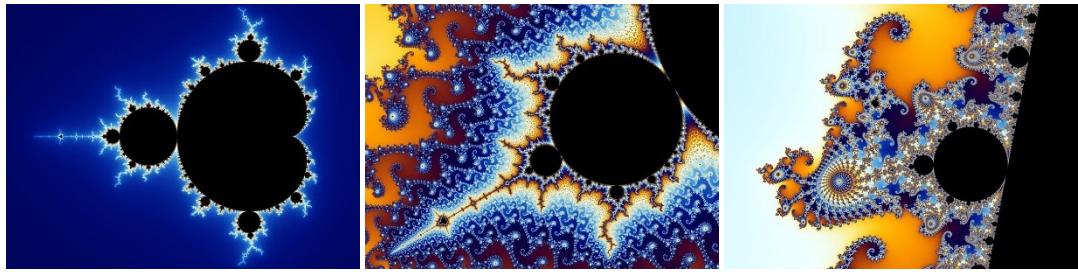


Figure 2.1.3 Mandelbrot set overview (left), antenna (middle), and seahorse (right)

The Mandelbulb is a commonly used 3D representation of the 2D Mandelbrot fractal, created by Daniel White and Paul Nylander. For many years, it was thought that a true 3D representation of the Mandelbrot fractal did not exist, since there is no 3D representation of the 2D space of complex numbers, on which the Mandelbrot fractal is built upon [14]. While this is still the case, White and Nylander made a significant breakthrough which resulted in the creation of a 3D fractal bearing similar characteristics to the 2D Mandelbrot set.

White and Nylander considered some of the geometrical properties of the complex numbers. The multiplication of two complex numbers is a kind of rotation, and the addition is a kind of transformation. White and Nylander experimented with ways of preserving these characteristics when converting from 2D to 3D, and their solution was to change the squaring part of the formula to instead use a higher power, a practice sometimes used with the 2D Mandelbrot fractal to produce snowflake type results [15]. This change leads us to the equation $f(z) = z^n + c$, where z and c are triplex numbers, representing a point with an x , y , and z coordinate. The value of n is varied to give different results, where $n = 8$ is commonly used as this results in a good amount of fractal detail when zooming in.

White and Nylander's formula for the n th power of a point in 3D space [15], [16] is given as:

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix}^n = r^n \begin{bmatrix} \sin(n\theta) \cos(n\varphi) \\ \sin(n\theta) \sin(n\varphi) \\ \cos(n\theta) \end{bmatrix}$$

where $r = \sqrt{x^2 + y^2 + z^2}$, $\theta = \text{atan2}(\sqrt{x^2 + y^2}, z)$, $\varphi = \text{atan2}(y, x)$

Equation 2.1.i White and Nylander's formula for the n th power of a point in 3D space

And the addition of two points is given as:

$$\begin{bmatrix} x_1 \\ y_1 \\ z_1 \end{bmatrix} + \begin{bmatrix} x_2 \\ y_2 \\ z_2 \end{bmatrix} = \begin{bmatrix} x_1 + x_2 \\ y_1 + y_2 \\ z_1 + z_2 \end{bmatrix}$$

Equation 2.1.ii Addition of two points in 3D space

Using both formulas, the equation $f(z) = z^n + c$ can now easily be solved and when rendered in 3D results in some beautiful images. Some renders from [Daniel White's website](#) are shown in Figure 2.1.4.

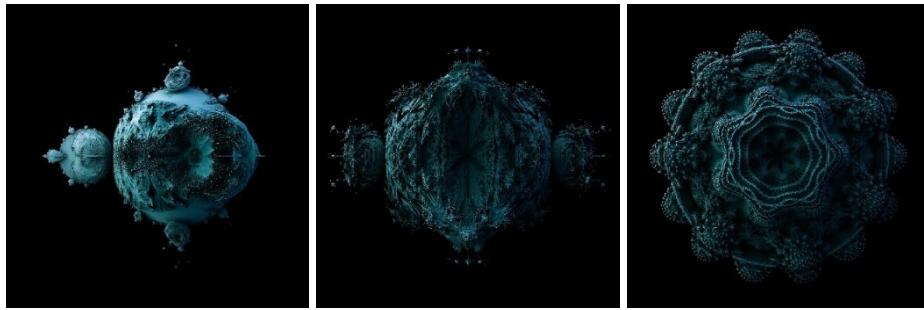


Figure 2.1.4 Mandelbulb power of two (left), three (middle), and eight (right) [17]

The power of two version of the function (left) results in a filled 3D Mandelbrot set, with little fractal detail. Power of three (middle) contains more detail, and power of eight (right) is the sweet spot for this formula, in the trade-off between fractal detail and performance. There exist several other variations of Mandelbulb formula, each balancing performance with fractal detail.

The most notable performance improvement for this algorithm was discovered by David Mankin [18], and uses a distance estimation function, which for any point in 3D space, returns an estimation of the distance to the surface of the geometry. Distance estimation functions exist for many other 3D fractals as well, meaning that this is a valid generalised method for rendering 3D fractals.

In addition to the Mandelbulb, there exist many other 3D fractals. One interesting example are the Julia sets, which come from the same $f(z) = z^2 + c$ equation as the Mandelbrot set, but the value of z is varied instead of the value of c .

2.1.4 3D Fractals Summary

In summary, there exist two main ways of creating 3D fractals. The first approach is used for creating the platonic solid fractals and relies on taking primitive shapes and applying

transformations, rotations, and scaling operations to them. The main bottleneck of rendering this type of fractal is the number of objects in the scene that must be rendered. The second approach is used for rendering more natural looking fractals and relies on plotting the convergence of equations and arbitrarily colouring them depending on how quickly the values converge. The main bottleneck of this approach is the number of iterations required for the values to converge. However, this can be improved, and the most common optimisation is the use of a distance estimation function.

Both approaches discussed create completely different looking fractals, and the application will contain examples of both. However, to view these fractals, a suitable rendering approach which supports both primitive object transformation, rotation, and scaling operations, and supports the rendering of a surface defined by a series of points in 3D space. This will allow both fractal rendering approaches to be supported.

2.2 Fractal Rendering Methods

2.2.1 Rasterization

In computer graphics, there are two commonly used methods of rendering an image of a 3D scene. The first is called rasterization, which renders an image of a scene by iterating through all objects in that scene, determining which pixels on the screen are affected by that object, and modifying them accordingly. This approach has many benefits [19], most notably its speed when rendering an image. Additionally, when doubling the number of pixels in an image, the time taken to rasterize the image generally increases by less than double. Because of these benefits, rasterization is the most common rendering method, and all graphics cards are designed to efficiently render images using this approach.

Rasterization is very well suited for rendering 3D objects that are stored as meshes, which contain vertices, edges, and faces. Unfortunately, a 3D fractal could not be converted into a mesh without losing detail, as a mesh only contains a finite number of points, and a 3D fractal must contain close to infinite detail.

2.2.2 Ray Tracing

Ray tracing is another method of rendering an image of a 3D scene. When rendering an image using ray tracing, for each pixel on the screen, a ray (simply a line in 3D space) is extended or traced forwards from the camera position until it intersects with the surface of an object. From there, the ray can be absorbed or reflected by the surface and more rays can be sent out recursively. Ray tracing is ideal for photorealistic rendering as it takes into consideration many of the properties of light, through simulating reflections, light refraction, and reflections of reflections [20]. Often, ray tracers do not render images in real-time as the process is computationally expensive. To make a ray tracer capable of rendering in real-time, many approximations must be made, or hybrid approaches used.

Fractals cannot be efficiently rendered using pure ray tracing. Early attempts chose to sample points uniformly along each ray to check at which point the ray had collided with the fractal. While this approach did work, it is a crude and inefficient solution. With a little modification, the ray tracing algorithm can be modified to render an image represented by a distance estimation function instead. This replaces the uniform sampling of points along the ray with points distributed using the distance estimation, which is both a faster and more detailed approach, since empty space will not be sampled, and detail will not be missed as points can be sampled close to each other. This approach is called ray marching, and it can be used to render 3D fractals since fractal geometry can be represented using a distance estimation function. The optimisation for the Mandelbulb fractal which uses distance estimation [18] gives such a massive performance increase when compared to the crude ray tracing approach, that this approach can be used to render fractals in real time.

2.2.3 Ray Marching

Ray marching (sometimes known as sphere tracing) is a variation of ray tracing, which only differs in the method of detecting intersections between the ray and geometry. Instead of a using a ray-surface intersection function which returns the position of intersection, ray marching uses a distance estimation (DE) function, which simply returns the distance from any

given position in the scene, to the surface of the closest geometry. Instead of shooting the ray in one go, ray marching uses an iterative approach, where the current position is moved/marched along the ray in small increments until it lands on the surface of an object. For each point on the ray that is sampled, the DE function is called and marched forward by that distance, and the process is repeated until the ray lands on the surface of an object. If the distance function returns zero at any point (or is close enough to an arbitrary epsilon value), then then the ray has collided with the surface of the geometry. Figure 2.2.1 shows a ray being marched from position p_0 in the direction to the right. The distance estimation for each point is marked using the circle centred on that point.

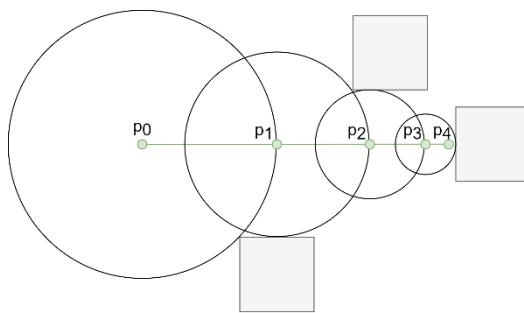


Figure 2.2.1 Ray marching diagram

Technically, the DE does not have to return the exact distance to an object, as for some objects this may not be computable, but it must never be larger than the actual value. However, if the value is too small, then the ray marching algorithm becomes inefficient, so a fine balance must be found between accuracy and efficiency.

2.2.3.1 Benefits of Ray Marching

Ray marching may sound more computationally complex than ray tracing since it must complete multiple iterations of an algorithm do what ray tracing does in a single ray-surface intersection function, however, it does provide several benefits. As mentioned above, ray marching does not require a surface intersection function like ray tracing does, which means it can be used to render geometry for which these functions do not exist, such as many 3D fractals.

While many effects such as reflections, hard shadows and depth of field can be implemented almost identically to how they are in ray tracing, there are several optical effects that the ray marching algorithm can compute very cheaply.

Ambient occlusion is a technique used to approximate how exposed each point in a scene is to ambient lighting [20]. This means that the more complex the surface of the geometry is (with creases, holes etc), the less ways ambient light can get into it those places and so the darker they should be. With ray marching, the surface complexity of geometry is usually proportional to the number of steps taken by the algorithm [21]. This property can be used to implement ambient occlusion and comes with no extra computational cost at all.

Soft shadows can also be implemented very cheaply, by keeping track of the minimum angle from the distance estimator to the point of intersection, when marching from the point of intersection towards the light source [22]. This second round of marching must be done anyway if any type of lighting is to be taken into consideration, so minimum check required for soft shadows is practically free.

A glow can also be applied to geometry very cheaply, by keeping track of the minimum distance to the geometry for each ray. Then, if the ray never actually collided with the geometry, a glow can be applied using the minimum distance the ray was from the object, a strength value, and colour specified [21].

2.2.3.2 *Signed Distance Functions*

A signed distance function (SDF) for a geometry, is a function which given any position in 3D space, will return the distance to the surface of that geometry. The distance contains a positive sign if the position is outside of the object, and a negative sign if the position is inside of the object. If a distance function returns zero for any position, then the position must be exactly on the surface of an object. Every single geometry in a scene must have its own SDF. The scenes distance estimation (DE) function will loop through all the SDF values for geometry in the scene and will return the distance to the closest geometry surface.

The sign returned by the SDF is useful as it allows the ray marcher to determine if a camera ray is inside of a geometry or not, and from there it can use that information to render the objects differently. We may want to render geometry either solid or hollow, or potentially add transparency.

Signed distance functions are already known for most primitive 3D shapes, such as spheres, boxes, and planes. A full list of these functions can be found on Inigo Quilez's web page [23]. A full example of the SDF for a sphere is included in appendix 6.2. Since distance estimation can be used to represent primitive shapes, the feature still required to render the platonic solid fractals is the ability to transform, rotate and scale primitives.

2.2.3.3 Transforming SDFs

Translating and rotating objects is trivial when using distance estimation. All that is required, is to transform the point being sampled with the inverse of the position and rotation used to place an object in the scene [23]. A function to transform a point by a translation and rotation is given in Equation 2.2.i.

$$\text{transform}(p, t) = \text{invert}(t) \times p$$

where $p = \begin{bmatrix} x \\ y \\ z \end{bmatrix}$, t is a 3×4 transformation matrix storing only translation and rotation

Equation 2.2.i Equation for transforming a point in space

Using matrices to store transformations and rotations is common practice in computer graphics, and is fully explained in this article [24].

Uniform scaling of an SDF can be completed easily by first decreasing the scale back to one, sampling the point, and then increasing the scale back up again. Functions for decreasing and increasing the scale of a point in 3D space are given in Equation 2.2.ii.

$$\text{scaleDown}(p, s) = \frac{p}{s}, \quad \text{scaleUp}(p, s) = p \times s, \quad \text{where } p = \begin{bmatrix} x \\ y \\ z \end{bmatrix}, s \in \mathbb{R} \text{ is the scale}$$

Equation 2.2.ii Equations for scaling a point down (left) and up (right)

2.2.3.4 Combining SDFs

Another feature that would be nice to include in the renderer is the ability to combine SDFs, and to have multiple objects in the same scene. SDFs can be combined using the union, subtraction, and intersection operations [25], as given in Equation 2.2.iii.

$$\text{union}(a, b) = \min(a, b), \text{ subtraction}(a, b) = \max(-a, b), \text{ intersection}(a, b) = \max(a, b)$$

where $a, b \in \mathbb{R}$ are the values returned from object a and b's SDF

Equation 2.2.iii Equations for union (left), subtraction (middle) and intersection (right) of two SDFs

There also exist variations of the formulas given in Equation 2.2.iii which can be used to apply a smoothing value to the operation. These formulas have been included in the appendix 6.3. The images in Figure 2.2.2 were rendered using the application.

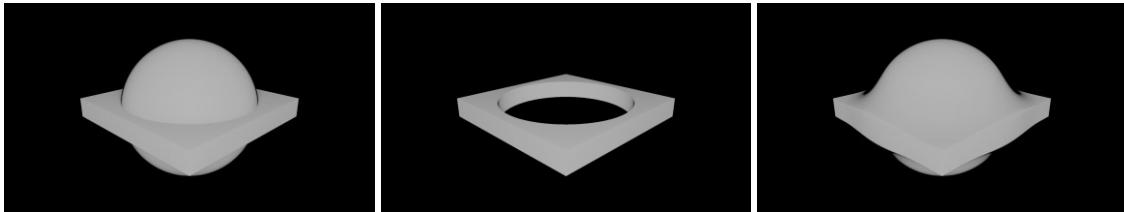


Figure 2.2.2 Ray marched sphere and box scene experiment union (left), subtraction (middle) and smooth union (right)

There are several additional alterations that can be applied to primitives once we have their SDF [23]. A primitive can be elongated along any axis, its edges can be rounded, it can be extruded, and it can be “onioned” – a process of adding concentric layers to a shape. All these operations are relatively cheap. Signed distance functions can also be repeated, twisted, bent, and surfaces displaced using an equation such as a noise function or sin wave, though these alterations are more expensive. All these techniques mentioned will be essential when creating more complex geometry.

2.2.3.5 Surface Normal

The surface normal of a position on the surface of a geometry, is a normalised vector that is perpendicular to the that surface. This information is essential for most lighting calculations and surface shading techniques, such as phong shading, which is discussed in section 2.3.1. When using distance estimation, the surface normal of any point on the geometry in a scene can be determined by probing the SDF function on the x, y and z axis, using an arbitrary epsilon

value. The formula used to calculate the surface normal of any point using distance estimation can be found in appendix 6.4.

2.2.3.6 Ray Marching Summary

In summary, ray marching is a variation of ray tracing which contains all the building blocks required for rendering both types of 3D fractals. It uses distance estimation and signed distance functions for calculating intersections with geometry, primitives can be modelled using SDFs, and SDFs can be translated, rotated, and scaled. In addition, there are many mathematical operations that can be used for combining SDFs to create more complex geometry and the surface normal of geometry can be calculated which means that it is possible to create advanced rendering features like surface shaders.

2.3 Surface Shading

The following section gives a brief introduction to the common surface shading technique Phong shading. There exist many surface shading techniques, but Phong is one of the most common as it provides an excellent trade-off between visual quality and performance. Other Surface shading techniques have not been discussed as they are out of scope of the project.

2.3.1 Phong Reflection Model

The Phong reflection model is a surface shading technique used in computer graphics [20], which describes the way that a surface appears visually when light is shining on and reflecting off it. It does this by combining ambient light, diffuse light, and specular highlights. The Phong reflection model attempts to mimic the real-world characteristics of shiny and rough objects, which create small and intense specular highlights, and large and dull highlights respectively.

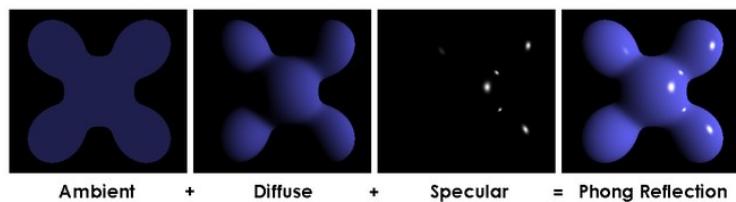


Figure 2.3.1 Phong reflection model [26]

The Blinn-Phong reflection model is a modification made to the Phong reflection model which reduces the computation cost when the light source and object are very far away from the origin in space. This modification does not reduce visual quality in any way and should be used in place of the Phong reflection model due to its better performance in certain scenarios.

The following section gives a brief introduction to GPU computing and how it will be used in the project to render scenes in parallel on the GPU.

2.4 GPU Computing

GPU computing is when a GPU is used in combination with a CPU to execute some code [27]. The correct use of GPU computing improves the overall performance of the program by offloading some computation from the CPU to GPU. A CPU is designed for executing a sequence of operations, called a thread. A CPU can execute a few tens of threads in parallel and is designed to execute them as fast as possible [28]. On the other hand, a GPU is designed to execute a few thousand of these threads in parallel but does it significantly slower than the CPU would, achieving a higher overall throughput. This difference in architecture between CPUs and GPUs means that CPUs are better suited for computations requiring data caching and flow control, while GPUs are better suited for highly parallel arithmetic computations.

Implementing a graphics renderer using GPU computing is the perfect choice, as this is a massively parallel task requiring a computation to be executed for every single pixel on the screen. Making use of this architecture is the only feasible way to implement a real-time renderer. When implementing a real-time fractal renderer, the extra features that libraries like OpenGL [29] provide, such as compute shaders, aren't needed. Instead, it makes more sense to make use of a general-purpose GPU computing library, as the overhead of this software will be significantly less, giving us better overall performance.

Currently, there are two suitable general-purpose GPU computing interfaces, CUDA (Compute Unified Device Architecture) and OpenCL (Open Computing Language). Both interfaces make

use of a kernel language, a specialised programming language in which the code to be executed in parallel is written. This code is then compiled and executed on the devices GPU. Sections 2.4.1 and 2.4.2 dicuss CUDA and OpenCL.

2.4.1 CUDA

CUDA [28] is a closed source, cross platform development environment for creating high performance GPU accelerated programs, created by Nvidia, and only supported on Nvidia GPUs. This means that programs written using CUDA can only be run on Nvidia brand graphics cards. Despite this downside, CUDA is incredibly powerful and provides the CUDA C++ kernel language, which allows the use of many modern C++ features when writing the kernel. Additionally, the CUDA runtime system provides excellent performance as it makes full use of all available architecture dependant features, such as the new ray tracing cores on the current generation of Nvidia GPUs, which are designed for efficiently computing ray tracing calculations.

One downside of CUDA is that the Nvidia CUDA compiler `nvcc` is required when compiling CUDA programs at runtime. Ideally, the fractal rendering application would want to support this, as it would allow users to create new scenes and modify existing scenes. However, if CUDA was used, this would add a large dependency to the program, which users would have to install alongside the application. The other option would be to only allow developers to create scenes, and to ship a selection of examples alongside the application, however, this would significantly reduce the appeal of the application for most users. This downside is a major drawback for the suitability of CUDA when implementing a 3D fractal rendering program, where users can create their own scenes.

2.4.2 OpenCL

OpenCL [30] is an open source, cross platform, and cross device GPU computing interface, created by the Khronos Group. This means that OpenCL programs can run on most hardware [31], some common vendors include Nvidia, AMD, Intel, Samsung, Google, Nintendo, ARM, Raspberry PI, and QUALCOMM. This makes OpenCL the superior choice when aiming to deploy

an application to any type of graphics card. However, this flexibility comes with a cost, as each vendor implements their own OpenCL runtime system, meaning that the behaviour of implementations will vary between devices, which could result in inconsistent behaviour. Additionally, this makes it unlikely that OpenCL programs achieve the best possible performance, as the OpenCL API must cater to all vendors, making architecture dependant features of certain devices unusable.

OpenCL uses the OpenCL C kernel language, based on the C99 standard. The OpenCL C++ (C++17) kernel language is starting to be adopted but is not supported on all devices. This limits the kernel language to OpenCL C, which is very primitive when compared to C++, and makes development more difficult. However, this language has been supported for decades and all devices can compile kernels written in this language at runtime, making it an excellent choice for creating scenes for the application.

2.4.3 GPU Computing Summary

The main differences between CUDA and OpenCL, is that CUDA can only run on NVIDIA branded GPUs, and requires a large dependency for compiling kernels at runtime. While this technically does allow CUDA to make full use of architecture dependant features, potentially giving better performance, it also makes deployment of the application far less portable as the code will not run on other brand of GPUs. OpenCL is, therefore, the superior choice when aiming to write portable GPU code.

2.5 Review of Existing Applications

A review of relevant and popular existing 3D fractal renderers was completed, and the key features of each application were recorded. Appendix 6.5 contains screenshots from all applications discussed in the following section.

2.5.1 Fragmentarium

Fragmentarium [32], [33] was originally released in 2011, and was one of the first open-source programs capable of rendering 3D fractals. It has received many new features over the years and is still being updated to this day. The Fragmentarium interface contains three main panels: a GLSL shader language IDE with syntax highlighting, used for defining geometry in the scene, a main window displaying the visual output for the camera, and a panel containing a selection of sliders for controlling various parameters of the scene or visual output. This interface allows users to update parameters of the scene in real time and view the output. Users can also manually edit and recompile the GLSL shader code for the scene and view its output. This is a very nice feature as it allows users to develop scenes and view their changes in real time, without having to close the application or open the scene file in a different program.

Fragmentarium contains many example scenes other than just the Mandelbulb and Sierpiński fractals, and the rendering quality of the output image is good but does vary between scenes.

2.5.2 Synthclipse

Synthclipse [34] is shader prototyping tool, which extends the Eclipse IDE. It builds upon features from an early version of Fragmentarium, by adding support for JavaScript scripting, OpenGL, music playback, key frame animations, and Shadertoy integration. Synthclipse isn't just a tool for modelling fractals and can edit and run shaders written in most languages. Unfortunately, Synthclipse doesn't come with any example scenes which makes it very hard to learn how to use the application. Additionally, Synthclipse has not been updated since 2019.

2.5.3 Mandelbulb3D

Mandelbulb3D [35] is a modern fractal rendering application, featuring a clean and stylish user interface. Like Fragmentarium, Mandelbulb3D contains many windows for controlling various parameters in the scene, allowing the user to manipulate the fractal and see the effects in real time. However, unlike Fragmentarium, Mandelbulb3D contains many more parameters which makes it much more powerful. This also makes the interface feel clunky as it is sometimes hard

to find features. Additionally, Mandelbulb3D uses its own language as an interface to define scenes which makes the learning curve of the application very steep.

Mandelbulb3D contains powerful animation and video editing tools, making it perfect for creating highly detailed screenshots and videos of 3D fractals. While this application provides real-time feedback when changing parameters in the scene, it is less suited for exploring the scene in real time as the visual quality is so high. To make real-time exploration possible, the developers have reduced the output resolution and visual quality when moving the camera in preview mode. Even with this change, the performance of the real-time preview is sub optimal which makes it hard to navigate the scene. Output images created by Mandelbulb3D are very detailed, coloured beautifully and contain many advanced optical effects.

2.5.4 Mandelbulber

Mandelbulber [36] is a cross platform 3D fractal rendering application, which contains similar features to Mandelbulb3D but with a considerably simpler and more intuitive user interface, which reduces the learning curve of the application. However, the visual output of Mandelbulber is worse than that of Mandelbulb3D. Additionally, the real time preview generated by Mandelbulber is very low quality. Mandelbulber uses OpenCL C to define its scenes, but it also uses some custom syntax which can make it harder to learn. Output images created by Mandelbulber are highly detailed contain advanced optical effects.

2.5.5 FractalLab

FractalLab [37] is a web-based 3D fractal renderer, which uses the GLSL renderer in WebGL (web-based OpenGL) to compile and run scenes. This is a hobby project and has remained in the proof-of-concept stage since 2011. Despite this, FractalLab is powerful enough to update and render 3D fractals in real time, all from inside a web browser. Admittedly the visual quality is not as high as some of the other applications discussed, but it does contain materials, lighting, and ambient occlusion. FractalLab uses a simple but well-designed user interface for controlling

scene parameters at runtime. Out of all the applications investigated, this is by far the most impressive in terms of creativity and real time rendering performance.

2.5.6 Existing Applications Summary

After reviewing existing 3D fractal rendering applications, there are several key points that should be considered when designing the application. To summarise, of all the applications discussed, FractalLab is the only one designed for real-time rendering of 3D fractals, while the others are designed for offline rendering of images or videos, and the real-time preview is low resolution and low quality. The applications also use different programming languages to define their scenes, Fragmentarium and FractalLab both use the GLSL shader language, Mandelbulber uses OpenCL, Mandelbulb3D uses its own language, and Synthclipse supports all common shader languages. Additionally, all applications discussed use their own parser on the scene file to create the interactive GUI elements for editing the scene parameters at runtime. This adds a large amount of complexity and development time and isn't required for rendering 3D fractals in real time but would be nice feature to implement.

Additionally, several applications suffered from clunky interface design which made it hard to access features and increased the learning curve of the application. Some of the applications made it hard to find example scenes.

2.6 Literature Review Summary

In summary, there are two main methods of creating fractal patterns - recursively applying transformations to primitive shapes and plotting the convergence of equations. To be able to render both types of fractals, a variation of ray tracing called ray marching must be used, which uses a distance estimation function to calculate the distance to the closest piece of geometry in the scene. In addition, to be able to render 3D fractals in real time, GPU computing techniques must be utilised to execute code in parallel. Several existing applications already use similar approaches to great success, however, these applications do contain some flaws and a note of these have been made so that this project will not make the same mistakes.

3 Development

The following section describes the development environment and chosen design for the application, the technologies used, and development strategy employed to achieve the project aims and objectives. Additionally, the application documentation is referenced and future improvements to the application are discussed.

3.1 Development Environment

One of the most efficient ways of developing a C++ application on Windows is using Visual Studio [38], which provides powerful development tools such as debugging support, a unit testing interface, refactoring options, code snippets, and IntelliSense, Microsoft's automatic code completion software. Visual Studio was chosen as the development IDE for these reasons, and Visual Studio 2019 was installed on a desktop Windows 10 machine to be used for development of the application. Early on, a [project GitHub repository](#) [5] was created ensure that all work carried out for the project was backed up using a version control system. GitHub Desktop was used as the version control interface as it is intuitive to use.

Development of the application was completed using a Windows 10 computer with a Nvidia RTX 3060 Ti graphics card, a mid-tier graphics card from the current new generation of Nvidia GPUs.

3.2 Application Design

3.2.1 High Level Overview

Figure 3.2.1 displays a high-level overview of how the `FractalGeometryRenderer` interacts with the CPU and GPU of the device it is currently running on. For the most part, the application is single threaded and runs on the CPU. When the application starts, the scene kernel file is loaded to the devices main memory, then it is passed to the OpenCL runtime system which compiles the scene and then loads it onto the devices GPU. Once this is done, the application enters a loop where it will instruct the kernel about *how* to render the scene, specifying

information such as the resolution to output, the position of the camera, the direction that the camera is facing, and the current time that the application has been running for. The kernel then computes the pixel data and writes it to a buffer in main memory, and the application reads that pixel data and updates the display accordingly.

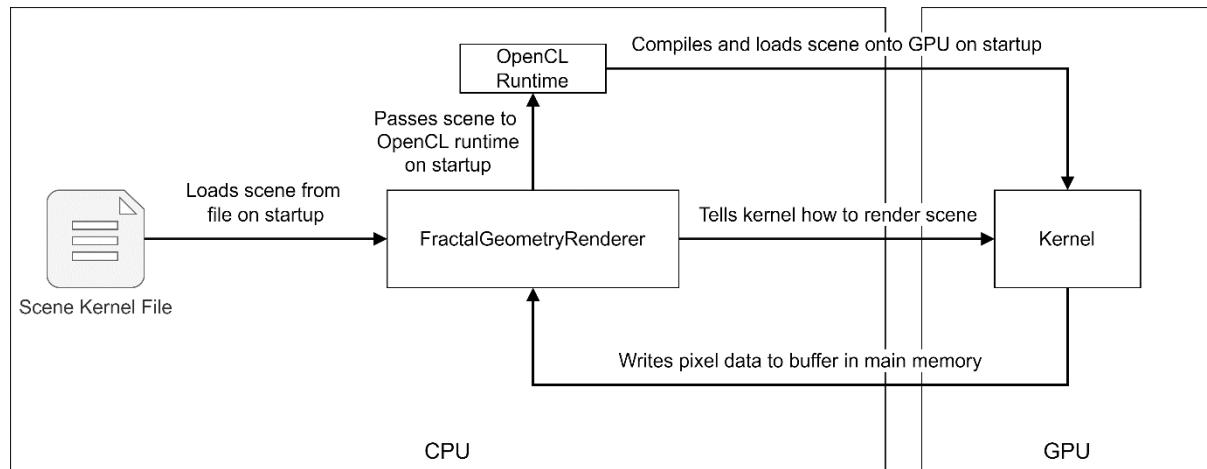


Figure 3.2.1 Application high level overview

Most of the time taken to render each frame comes from computing the colour of each pixel, a process which is executed on the GPU. During this time, the CPU is idle as it must wait for this process to complete before continuing. This means that the application is GPU bound, as the limiting factor for the application performance is the performance of the graphics card itself. This is important because it allows the specifications of the CPU to be ignored when benchmarking the application as they have little to no impact on performance, which allows a direct comparison of the performance of the application over various GPUs to be made. Additionally, over the coming years as the performance of GPUs increase, so will the performance of the application.

3.2.2 Application Design Principles

The application has been developed following modern C++ practices [39]. This encourages the use of objects and emphasises the principle of *resource acquisition is initialisation* (RAII), which states that resources such as heap memory, file handles, and sockets should be *owned* by the object that creates that resource in its constructor and, therefore, is also responsible for deleting that resource in its destructor once it is no longer required. The principle of RAII ensures that

all resources are correctly returned to the operating system once an object goes out of scope. This is especially important for this application, which must allocate and deallocate memory buffers on the GPU. There would be severe performance consequences for the device if GPU memory was not correctly deallocated.

The application has been structured using several core classes, which are loosely coupled and highly abstracted, which encourage the object-oriented principle of encapsulation. An object-oriented programming style was chosen to ensure the RAII principle was followed. Table 3.2.1 lists the core classes of the application and their responsibilities.

Table 3.2.1 Class responsibilities

Class name	Responsibilities
FractalGeometryRenderer	<ul style="list-style-type: none"> • Drives the application using a game loop • Contains instances of Renderer and Window
Renderer	<ul style="list-style-type: none"> • Calculates the pixel data for the current frame, and writes that data to main memory
Window	<ul style="list-style-type: none"> • Reads pixel data from main memory and updates the window accordingly • Reads window events and user input events

The colour format chosen by the application is RGBA8888, which means that each red, green, blue, and alpha channel of a colour uses eight bits, representing a value from 0 to 255. Theoretically, this allows for 4,294,967,296 different colours, but in practice only 16,777,216 of these colours will be used by the application as the alpha channel is always set to 255 which makes the colour opaque. Leaving this channel unused is not ideal, but it is the simplest solution, and many GPUs are optimised for rendering textures using this colour format which gives us good performance at the cost of a small amount of wasted memory. While this format contains more than enough colours, the main benefits of this format are memory usage and performance. Each pixel requires one byte of memory for storage, meaning a resolution of 1920 by 1080 takes up only 2.1 MB of memory. In the application, the contents of the pixel data buffer is copied to another buffer so that the graphics library SDL2 can efficiently stream the data to a texture and display it on the window. If the device supports it, SDL2 will stream this pixel data very efficiently using the GPU. The pixel data must be duplicated in memory as OpenCL and

SDL2 should not have access to the same memory location as this can cause conflicts and unexpected behaviour, as SDL2 may alter the pixel data when streaming it to the window. The duplication of pixel data is not ideal, but it is the simplest solution.

3.2.3 Kernel Design Principles

The kernel code is structured into several key OpenCL C files. The file `main.cl` contains all the algorithmic ray marching logic and contains the main kernel method which the C++ side of the application interacts with. The file `types.cl` contains `struct` definitions for `Material`, `Light` and `Ray`. The file `defines.cl` contains all pre-processor definitions used by `main.cl` and is used for assigning default values if a definition is not specified in the scene. A dependency graph for a typical scene is included in appendix 6.6.

The kernel code of the application has been designed so that a custom parser is not required by the application, as is done by all applications reviewed in section 2.5. This reduces the complexity of the application overall, at the cost of making the kernel code verbose.

Compiler directives are used throughout the file `main.cl` to enable and disable specific branches of code, depending on the settings defined in the scene. This means that features that are not wanted in the current scene aren't even compiled, which gives the application better performance as there is no redundant code. Each scene kernel file must implement several key methods which will be called by `main.cl`, these methods are listed in Table 1.1.1.

Table 3.2.2 Key scene methods to implement

Method signature	Description
<code>float DE (float3 p, float t)</code>	The signed distance estimation between point <code>p</code> and geometry in the scene, sampled at time <code>t</code>
<code>Material getMaterial (float3 p, float t)</code>	The material used to render the geometry at point <code>p</code> and time <code>t</code>
<code>Light getLight (float t)</code>	Properties of the scene light at time <code>t</code>
<code>float boundingVolumeDE (float3 p, float t)</code>	The signed distance estimation between point <code>p</code> and bounding volume used by all geometry, sampled at time <code>t</code> Only used if <code>USE_BOUNDING_VOLUME</code> is defined as <code>true</code>

The design choice to separate the scene signed distance estimation function and the scene material sampling function has significant consequences. By separating the two functions, it gives significantly better performance for scenes for which the material is not calculated as part of the distance estimation function. For example, the scene `mandelbrot.cl` which maps the 2D Mandelbrot set onto a 3D plane, benefits significantly from this design choice as the graphics rendered onto the plain can be calculated independently to the distance estimation to the plane. This means that the graphics are only calculated once per frame when the ray has collided with the geometry. While this design choice gives significantly better performance for some scenes, it does provide slightly worse performance for others. For example, the scene `mandelbulb.cl` will perform slightly worse as one extra call will be made to the scene geometry function, since the material of this geometry is calculated as a by-product of calculating the distance to the geometry. In most cases, any negative performance is negligible as the computation cost of calling these functions once should be low, considering these functions are called hundreds of times per frame for each pixel on the screen.

Another important design choice was the use of the OpenCL data type `float3`, which contains three float values and is used to represent vectors and colours. Behind the scenes, OpenCL uses a `float4` to store the `x`, `y`, and `z` values of the `float3`, and it initialises the fourth component `w` to zero. This means that for every `float3` type that is used in the kernel code, four floats are allocated, and one is not used. This is wasteful of memory, though the benefits of using this type are significant. Most GPU architectures support SIMD operations (Single Instruction Multiple Data) which makes use of the parallel nature of the GPU to combine multiplication or addition operations on vector types into a single instruction. This means that the addition or multiplication of two `float4` types is completed in one operation, instead of taking four operations (one for each of the `x`, `y`, `z`, and `w` components). The performance benefits gained through using this type are significant, in our case giving around a 3x speedup assuming the GPU architecture supports this technique, which in most cases is worth wasting the memory of a 32-bit float type. If a GPU did not support this technique, then it would fallback to SISD

operations (Single Instruction Single Data) meaning that addition or multiplication would be completed in four operations. This is very unlikely as most GPUs have supported the basic addition and multiplication SIMD operations for many years now, though it is worth noting as this could contribute to poor performance on old devices.

When using GPU programming, realistically the precision of data types is restricted to 32-bit precision, as GPU architecture is designed for computations on these types. 64-bit data types are supported, however, GPUs can compute orders of magnitude more operations on 32-bit types than 64-bit types, making the use of 64-bit types unfeasible for a high-performance application.

A coalesced memory transaction is when multiple accesses to memory are combined into a single transaction. When accessing the main GPU memory, it is essential that this principle is followed as the main GPU memory is orders of magnitude slower to access than the local GPU memory. Memory coalescence for a work item is achieved by ensuring that each work item only accesses a single block of memory, which is consecutive to the memory accessed by previous work items. This principle has been followed when writing data to the output colour buffer.

It was initially planned to use the new OpenCL C++ kernel language [40], which allows kernels to be written using many C++17 features, most notably allowing classes (instead of C structs) and method/operator overloading, as these features would have allowed the kernel code to follow similar coding principles to those used by the main application. However, support for this new language was made optional in OpenCL 3.0 and was not supported on the device used for development, so unfortunately it could not be used. Additionally, if this new language had been used, scenes written in this language might not have been backwards compatible with devices using older OpenCL runtimes.

3.3 Technologies

The application uses OpenCL, a parallel programming framework that allows code to be executed on the GPU. OpenCL has bindings in several languages, most notably C++ which is a

high-performance system language. CMake is a cross platform language designed for automating the build process for large scale applications. CMake provides a simple API for including files and linking libraries and works well with package managers on most operating systems. In the Windows environment, the vcpkg package manager was used for installing development packages.

Table 3.3.1 lists the packages used by the application. All packages chosen are cross platform, as this provides an abstraction layer over platform specific libraries, which allows the program implementation to remain decoupled from the deployment platform. The application has only been built and tested on Windows x64, however, the use of CMake and cross platform libraries means that building the application on other architectures should be relatively straight forward.

SDL2 was chosen as the graphics library as it is low level, lightweight, and gives good performance for uploading 8-bit colour values directly to the window. A more complex graphics library was not required as the application does not contain a GUI.

Table 3.3.1 Packages used by the application

Technology	Description	Justification
OpenCL	GPU parallel programming framework	<ul style="list-style-type: none"> Standardised parallel programming framework for GPUs Version 1.2 is the mandatory minimum for all graphics cards Allows code to be deployed to any brand of GPU
SDL2	Graphics library	<ul style="list-style-type: none"> Provides window manipulation, user input event polling, and high-performance timers Has extensive documentation and examples
Eigen3	Maths library	<ul style="list-style-type: none"> Provides useful vector and quaternion maths functions, used when calculating the camera view direction
GoogleTest	Unit testing framework	<ul style="list-style-type: none"> Powerful unit testing framework Comes bundled with the Visual Studio 2019 C++ development package

3.4 Development Strategy

An Agile [41] development strategy has been used throughout the duration of the project to ensure that work was completed on time and that the project aims and objectives were fully achieved. The Agile approach uses small sprints of work to complete specific and defined tasks,

which allows teams to respond to change quickly as requirements and plans are updated. While this is an individual project and the Agile approach is normally used in teams, it was beneficial in this case as it encouraged continuous reflection and improvement of the application features and design. To ensure that requirements were implemented in time, a Gantt chart was created which broke down development into eight key sprints, each two weeks in length. The highest priority requirements were implemented in the earlier sprints which left the less important functionality and stretch goals to be implemented later in development. The initial and updated project Gantt charts are located in appendix 6.11.3.

Development of the application began in sprint zero, which aimed to set up the project environment and refactor existing code that was created when experimenting with a simple ray tracing program. This ray tracing code was refactored and translated to the OpenCL C kernel language so that it could be executed on the GPU. Additionally, SDL2 was set up which allows the manipulation of windows so that the output of the program can be viewed by the user.

The first functionality of the application was implemented in sprint one, which aimed to add a game loop, basic benchmarking tools, and user input controls to move the camera around the scene at runtime. This ended up being straightforward to implement as the progress made in sprint zero allowed these tasks to be completed quickly. The dynamic scenes task from the next sprint was brought forward and completed in this sprint instead.

Sprint two aimed to add the first visual features of the application, including ambient occlusion, lighting, and shadows. Lighting and shadows were implemented easily, as these features don't differ much from traditional ray tracing, however, there were problems finding a suitable ambient occlusion algorithm to implement, as the one chosen in the initial research report suffered from a significant flaw, which makes it unusable in some situations. This task was put on the back burner and materials were implemented instead using the Blinn-Phong reflection model. Additionally, the benchmark scene was not defined during this sprint as was originally planned due to lack of time, so this task was pushed back to the next sprint.

Sprint three aimed to complete the final release of the application and to create a full user guide explaining how users and developers should use the application. It was during this sprint that progress began to significantly fall behind the project plan, as there were performance optimisations that had to be made before the application could be released, and there were tasks from the previous sprint still to implement. These optimisations were implemented along with several key fractal scenes and the first draft of the user guide was created. More research into ambient occlusion algorithms was completed, and it was decided to move this task to the stretch goals and to instead focus on defining the benchmark scenes and producing a stable version of the application.

Sprint four aimed to complete the first draft of the final report and record all results for the project. Since progress had fallen behind the original estimation these tasks were not completed and instead the benchmark scenes were defined, and final build of the application created. This required lots of fine tuning of parameters in example scenes, ensuring that all features were accessible from the command line interface, and that there were no bugs in the final build.

In sprint five progress was made towards the first draft of the report and the application benchmarks were prepared and run on several different computers.

In sprint six the first draft of the report was completed, and in sprint seven changes were made to the draft based on feedback from readers, before completing the final draft.

Overall, an Agile development strategy worked well for this project as the continuous reflection on the project ensured that tasks were reprioritised when required and the project scope was adjusted if needed. The free time allocated in the original project plan was required as tasks were delayed and had to be pushed back. Several risks identified in the original risk analysis did occur, specifically *R-2 Change in Requirements* and *R-6 Delays due to Bugs*, but the negative impact of these risks was minimal, thanks to the effective mitigation plans created. No unidentified risks occurred during the duration of the project.

3.5 Documentation

Extensive documentation for the application was created using Doxygen, an automated documentation generator, and is [hosted using GitHub pages](#) [6]. The *Home* page of the documentation contains system requirements, an installation guide, and a basic user guide with instructions on how to use the application.

The *Manual Builds* page contains developer information explaining how to build the application manually on your own device, should you wish to modify the application or build it for your own operating system. Documentation of all classes used by the application is also included.

The *Scene Development Guide* page contains information explaining how one would create their own scenes for the application. This page contains links to the pages for some core kernel files, which list all methods and compiler directives that can be used to modify the contents of the scene and quality of the render. A hello world scene is included as a starting point.

3.6 Interface Design

The application features a simple command line user interface. A full list of commands can be found on the *Home* documentation page, discussed in section 3.5.

A graphical based user interface was not chosen as it simply was not necessary. All scene configuration information is located within the scene kernel file itself, the only parameters that need to be passed into the application at runtime describe which scene should be rendered, what resolution to render the scene at, and the mouse sensitivity for controlling the camera.

4 Evaluation

To accurately evaluate the successfulness of the project in respect to the original aims and objectives, several different evaluation strategies must be employed. First, the unit testing strategy will be discussed as this form of testing was vital during development in ensuring that the application behaviour was consistent with the expected behaviour. Then, the application benchmarking framework must be discussed, and results gathered from this framework will be analysed and used to assess the value of the application in terms of performance, functionality, and scalability. The overall success of the project in terms of requirements implemented and aims and objectives achieved will then be evaluated using a goal-based evaluation strategy. Finally, possible future developments for the application will be discussed.

4.1 Unit Testing

Unit testing was a small but vital part of ensuring that the application meets the requirements specified and behaves as expected. Unit tests were created for all key classes used by the application, including profiling tools such as the high-performance timer and benchmark data classes, and classes containing complex behaviour such as the scene class which is responsible for moving the camera along a camera path during benchmarks. In addition, unit tests were created for all key scenes in the application, and these tests ensured that each scene can be loaded by the application without crashing.

Due to the nature of the application, automated tests cannot be used to determine whether functionality has been implemented or if the functionality behaves as expected, as most of the features of the application have a visual output which cannot be detected using automated tools. Instead, a goal-based evaluation strategy has been employed to keep track of which requirements have been implemented and the corresponding objectives that each requirement falls under. This is discussed in section 4.5.

4.2 Benchmarking Framework

The benchmarking framework in the application is responsible for storing performance data of the application and outputting that data when the application closes. The benchmark provides two main pieces of functionality, firstly it compares the duration taken to render each frame to the current minimum and maximum and updates their values respectively. Secondly, it is responsible for calculating the time taken to execute key pieces of code. The performance benchmarking tools were added to the application to enable comparisons to be made between the performance of each new feature and the value brought by each new feature, in terms of visual quality or functionality.

The performance benchmark is interfaced solely through the scene kernel file, as this is more flexible than using a command line argument, and it keeps all scene parameters within one file. For the benchmark scenes, either a stationary camera or camera path was used as this ensures that the geometry viewed by the camera is consistent across runs. It would not be appropriate to allow the user to control the camera as this could significantly increase the variance of results. The camera path is defined using the `CAMERA_POSITIONS_ARRAY` and `CAMERA_FACING_DIRECTIONS_ARRAY` compiler directives, which both contain an array of `float4` types, for a position/facing direction (x, y, and z components) and a time (w component) at which this value should occur. Additionally, `BENCHMARK_START_STOP_TIME` is used to specify the duration of time that the benchmark should be active and `CAMERA_DO_LOOP` is used to move the camera back to the starting position when it reaches the end. This information is loaded from the scene file on start-up and data is passed to the `Scene` class, which is responsible for linearly interpolating between camera position/facing values, and for closing the application once the duration is exceeded.

While the benchmark is running, the time taken to render each frame is compared to the current minimum and maximum that has occurred which is updated accordingly. Additionally, this frame time is also added to the total benchmark duration and the number of frames that

have occurred is incremented. This total duration value is important as the benchmark will not run for the exact amount of time specified by the `BENCHMARK_START_STOP_TIME` directive, as it will always have to finish rendering the current frame before exiting. This means that the total duration result is likely to exceed the benchmark duration time by at most the duration of one frame, which varies between systems as this is dependent on system performance.

Once the benchmark has concluded, the results are appended to the file “`results.txt`” which can then be copied and pasted into a spreadsheet. Additionally, the results are also displayed on the terminal window so that the user can read them. Table 4.2.1 describes all benchmark results that are recorded by the application.

Table 4.2.1 Benchmark framework results

Description	Units	Origin
Scene name		Runtime Parameter
OpenCL build options		Runtime Parameter
Resolution width	Pixels	Runtime Parameter
Resolution height	Pixels	Runtime Parameter
Device name		OpenCL Query
Device version		OpenCL Query
Work group size		OpenCL Query
Clock frequency	MHz	OpenCL Query
Number of parallel compute units		OpenCL Query
Global GPU memory	Bytes	OpenCL Query
Local GPU memory	Bytes	OpenCL Query
Constant GPU memory	Bytes	OpenCL Query
Total benchmark duration	Seconds	Benchmark
Total number of frames rendered		Benchmark
Maximum frame time	Seconds	Benchmark
Minimum frame time	Seconds	Benchmark

Table 4.2.2 describes the results calculated from the raw data specified in Table 4.2.2.

Table 4.2.2 Results calculated from benchmark

Description	Units	Calculation
Total number of pixels	Pixels	$resolution\ width \times resolution\ height$
Global GPU memory	Gigabytes	$\frac{global\ GPU\ memory\ (bytes)}{1000000000}$
Local GPU memory	Kilobytes	$\frac{local\ GPU\ memory\ (bytes)}{1000}$

Constant GPU memory	Kilobytes	$\frac{\text{constant GPU memory (bytes)}}{1000}$
Mean frame time	Seconds	$\frac{\text{benchmark duration (seconds)}}{\text{total number of frames rendered}}$
Mean frames per second	Frames per second	$\frac{1}{\text{mean frame time (seconds)}}$
Maximum frames per second	Frames per second	$\frac{1}{\text{minimum frame time (seconds)}}$
Minimum frames per second	Frames per second	$\frac{1}{\text{maximum frame time (seconds)}}$
Upper difference	Percent	$\frac{\text{maximum FPS} - \text{mean FPS}}{\text{maximum FPS} + \text{mean FPS}} \cdot 2$
Lower difference	Percent	$\frac{\text{mean FPS} - \text{min FPS}}{\text{min FPS} + \text{mean FPS}} \cdot 2$

4.3 Fractals Implemented

The application features several scenes containing 3D fractal geometry, including the Sierpiński cube and tetrahedron fractals, and the Mandelbulb fractal. Two core benchmark scenes were created to analyse the performance of the application when rendering these fractals. The scene `sierpinski_collection.cl` contains both the Sierpiński cube and tetrahedron fractals, while the scene `mandelbulb.cl` contains the Mandelbulb fractal.

Several other non-trivial scenes have been implemented which don't feature 3D fractal geometry. The scene `mandelbrot.cl` contains the 2D Mandelbrot fractal mapped to a 3D plane, which the user can explore. The scene `mandelbrot_zoom.cl` contains a camera path zooming into the Mandelbrot fractal. The scene `terrain.cl` contains mountainous terrain generated using 2D Perlin noise and the scene `planet.cl` contains an Earth like planet generated using 3D Perlin noise. Additionally, the scene `infinite_spheres.cl` contains an infinite number of spheres which are calculated by distorting the space at which the ray position is sampled from. The application also contains several trivial scenes such as `hello_world.cl`, `trivial.cl`, and `sphere_box.cl`, which are all useful as examples for helping new users. The scenes `planet.cl` and `trivial.cl` were also benchmarked.

4.4 Evaluation of Application

4.4.1 Optimisations

Ensuring that the performance of the application was good enough to render 3D fractals in real time was a challenge, and there were several key optimisations implemented to achieve this.

4.4.1.1 *Linear Epsilon*

The first optimisation, called the linear epsilon optimisation, involves linearly increasing the epsilon value used for ray intersections as the ray gets further away from the camera. This means that the further a ray gets from the camera, the less precise the collision between the ray and geometry needs to be. This optimisation has a nice side effect of reducing the amount of noise generated on geometry far away, as distant geometry becomes less detailed. This makes scenes feel more natural and realistic, while also increasing performance significantly. The `LINEAR_INTERSECTION_EPSILON_MULTIPLIER` directive controls the rate at which the epsilon value linearly increases, which is useful for fine tuning the performance boost with respect to visual quality. This optimisation performs best when the camera is far away from geometry and has little effect on anything close.

4.4.1.2 *Bounding Volumes*

The second key optimisation involved adding bounding volumes around geometry in the scene. A bounding volume is any shape which contains the main geometry inside of its volume. This can be used when ray tracing or ray marching as a collision with geometry only needs to be calculated if the ray collides with the bounding volume first. Since the collision function for a primitive sphere or box is considerably cheaper to compute than that of the Mandelbulb fractal for example, this can result in a significant performance boost as it saves unnecessary computation when the camera is far away from the geometry or facing in a different direction.

Bounding volumes only work well when the bounding volume function is simple and cheap to compute, while also being as close to the original geometry shape as possible. This means that bounding volumes are less suitable for geometry which contains lots of empty space or does not

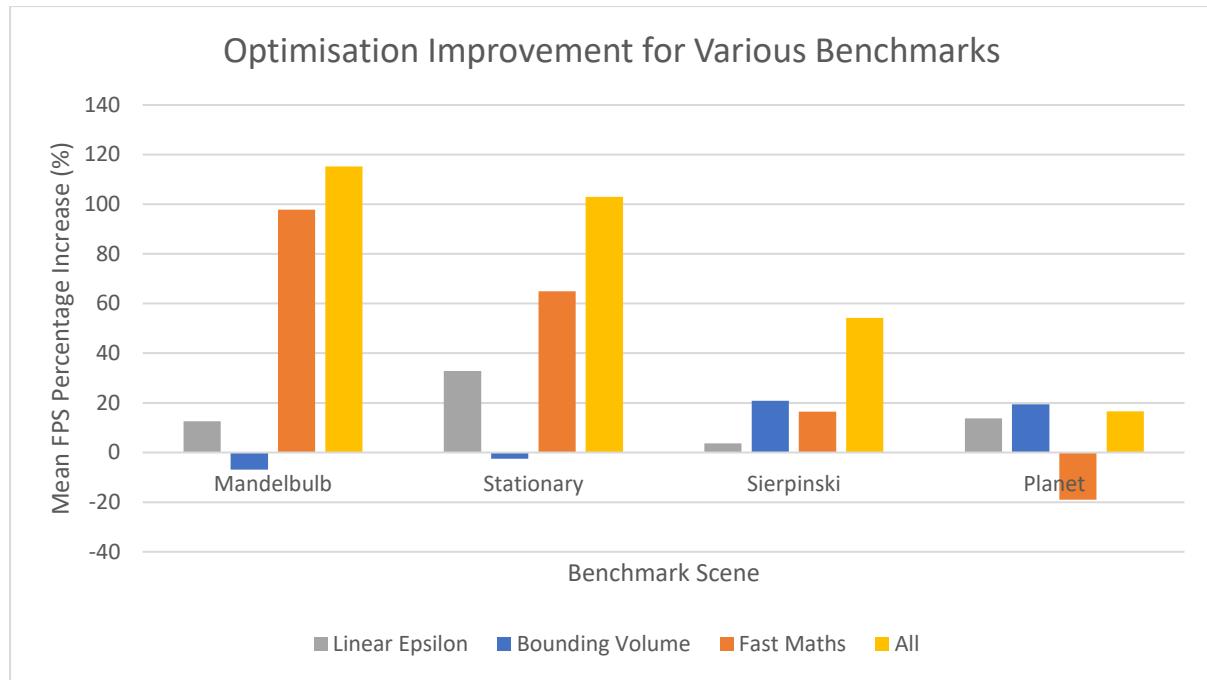
fit easily within a primitive shape. A single geometry could be subdivided into multiple bounding volumes, but this would add more complexity and computation cost and may not be worth it depending on the scene. Each scene must implement bounding volumes into a separate distance estimation function which the application uses when it doesn't require a precise distance estimation. Since bounding volumes may not be required for all scenes, this optimisation is disabled by default but can be enabled using the `USE_BOUNDING_VOLUME` directive. Some scenes may not see any performance boost when using bounding volumes or may even see worse performance if the bounding volume is too complex.

4.4.1.3 Fast Maths Build Flag

The final key optimisation results from using the `-cl-fast-relaxed-math` OpenCL build flag when compiling the scene at runtime. This build flag instructs the compiler to assume that the input will never be `NaN` or infinity, to ignore the signedness of zero, and to reduce the precision of certain arithmetic operations (such as $a \times b + c$) [42]. These optimisations provide a trade-off between performance and result correctness and subsequently violate several IEEE 754 conventions, which define floating point arithmetic. This optimisation only effects the compiled OpenCL kernel code, not the C++ side, meaning that only the visual output might be affected and not the accuracy of the profiling timers. Several screenshots comparing the output of the application when this optimisation is enabled and disabled are included in appendix 6.7. No noticeable visual difference can be seen when comparing the two images so this optimisation is enabled by default, but can be disabled using the application command line argument `--force-high-precision`.

4.4.1.4 Evaluation

The performance of the three optimisations discussed above were benchmarked using the scenes defined in section 4.2 and the results are displayed in Graph 4.4.1. Each benchmark was run three times and an average of the raw data was taken.



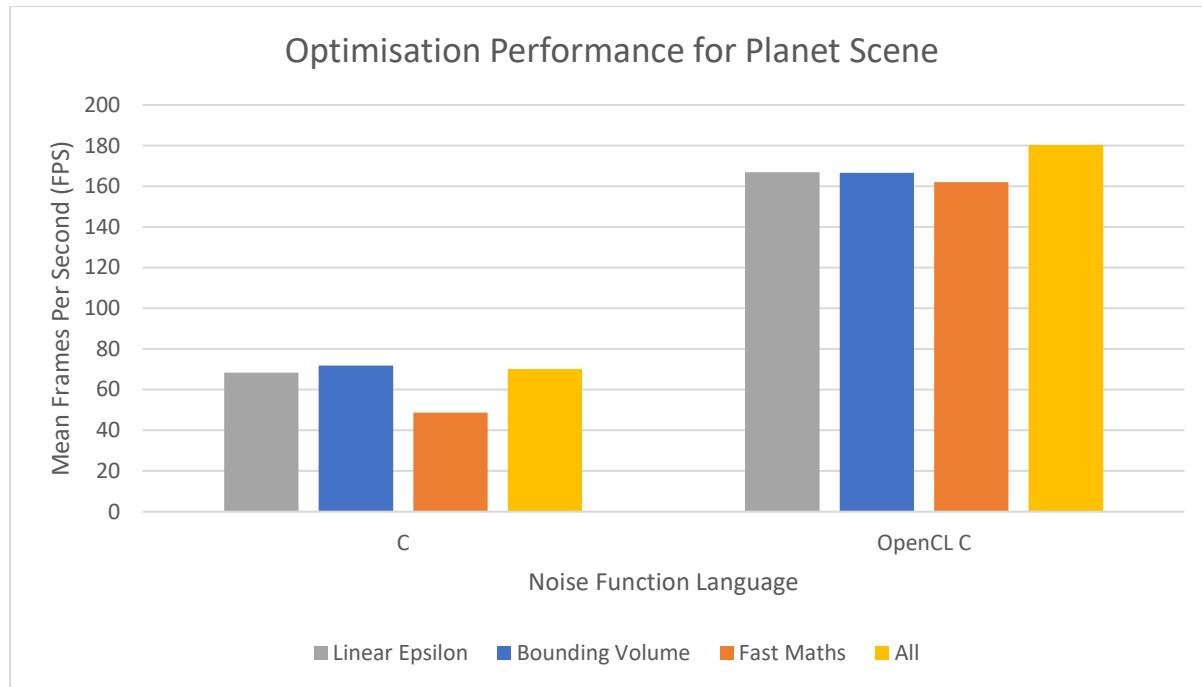
Graph 4.4.1 Performance of optimisations in the Mandelbulb scene

The linear epsilon optimisation resulted in between a 3% and 32% increase in mean frames per second when compared to the base version. This optimisation behaves similarly to the bounding volume optimisation in the sense that a significant performance boost is only seen when the camera is further away from the geometry, though in this case it performed significantly better than the bounding volume optimisation as the computational cost of using this optimisation is trivial, and so it has no negative performance effects when in suboptimal conditions.

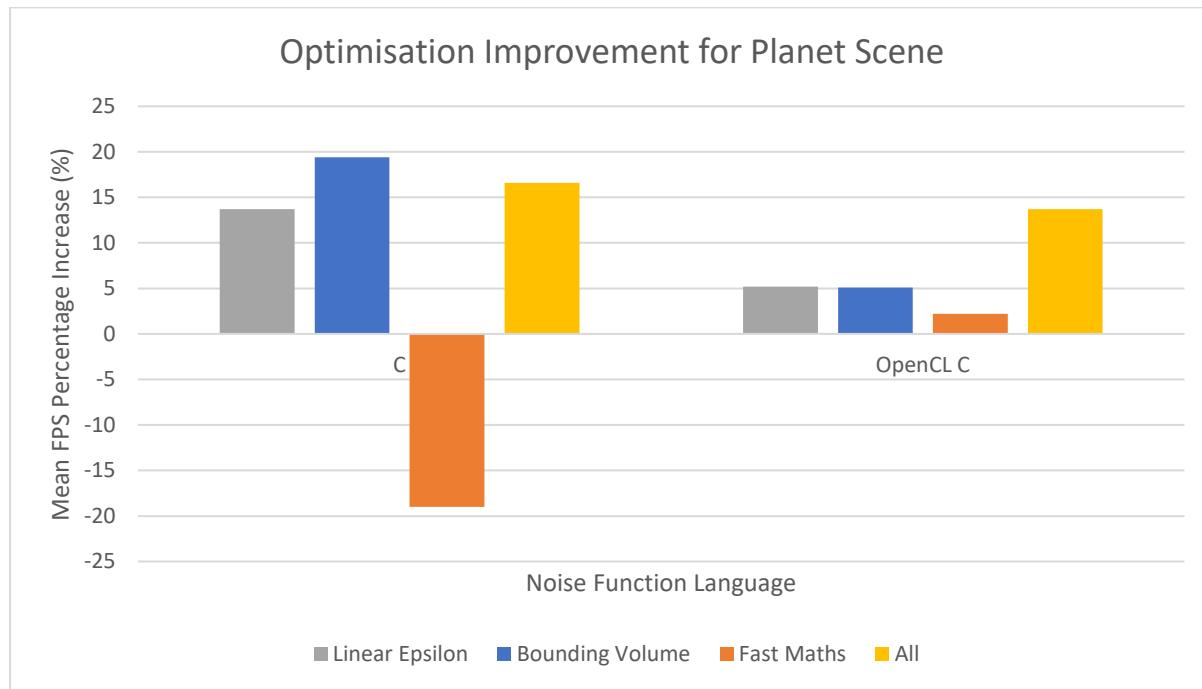
Interestingly, the Mandelbulb scene performed worse when using the bounding volume optimisation. This is likely due to a combination of the camera being close to the geometry and therefore not getting the benefit of using a bounding volume, and because the bounding volume for the Mandelbulb is substantially larger than the geometry itself, to ensure that the whole fractal is within the bounding volume since it changes shape over time. The efficiency of a bounding volume optimisation is highly dependent on the scene and camera position, and in this case, it didn't provide a performance increase. However, if the benchmark had spent more time further away from the geometry, then the performance could be significantly better. The bounding volume optimisation performed better when viewing the Mandelbulb from a distance in the Stationary scene, though performance of this optimisation was still worse than not using

it at all. For the Sierpinski and planet scenes, the bounding volume optimisation performed significantly better.

The fast maths optimisation performed significantly better for the Mandelbulb and stationary Mandelbulb scenes and performed acceptably for the Sierpinski scene. However, this optimisation performed significantly worse for the Planet benchmark scene. This is due to the large arrays of constant random noise data used by the C functions. When the C file was changed to a .cl file, and the arrays of data were marked with the OpenCL `__constant` identifier, the performance of this scene improved significantly, and the fast maths optimisation no longer resulted in worse performance. Presumably, when the OpenCL compiler converts a C code to OpenCL C code, it does not automatically identify constant data, meaning that each work item was instantiating its own copy of the data which would require large amounts of memory. Each work item only has a small amount of private memory allocated to it for temporary variables. It is likely that these arrays far exceed the amount of memory allocated to each work item, meaning that this data would spill over to the main GPU memory, which has orders of magnitude higher latency than the private memory of each work item, resulting in substantially higher computation times. These results can clearly be seen in Graph 4.4.2 which displays the mean FPS for the original C noise function and updated OpenCL C noise function, which results in an overall FPS increase over two times the original value. Additionally, Graph 4.4.3 displays the updated percentage increase for each optimisation, note that while the percentage increase in FPS is less than before for each optimisation, the overall FPS is considerably higher.



Graph 4.4.2 Comparison of FPS for Planet scene with C and OpenCL C noise function



Graph 4.4.3 Comparison of percentage FPS increase for Planet scene with C and OpenCL C noise function

For all benchmark scenes in Graph 4.4.1, the combination of all optimisations resulted in a greater FPS increase than the sum of their component parts. This is due to the synergistic nature of the linear epsilon optimisation and bounding volume optimisation, as when combined this reduces the accuracy of rays that are far away from the camera but also interacting with the bounding volume. Additionally, the use of the fast maths optimisation results in better

performance in most cases and when combined with the other two optimisations, this results in even better performance as the fast maths can optimise the computations performed by the other optimisations.

4.4.2 Efficiency

The GPU efficiency of the application has been profiled regularly throughout development to ensure that the implementation makes full use of all resources available to it. Nvidia Nsight Systems [43] is a GPU profiling tool for Nvidia graphics cards, which can be used to record kernel usage, command queue requests, and throughput of data packets. The profiling results for the Mandelbulb benchmark are displayed in Figure 4.4.1, with the green highlighted section representing the computation of one frame. Note the GPU usage during this time is close to 100%. The empty space after each frame is the time taken to display the frame on the window, a duration of between two and three milliseconds for a 1920×1080 image.

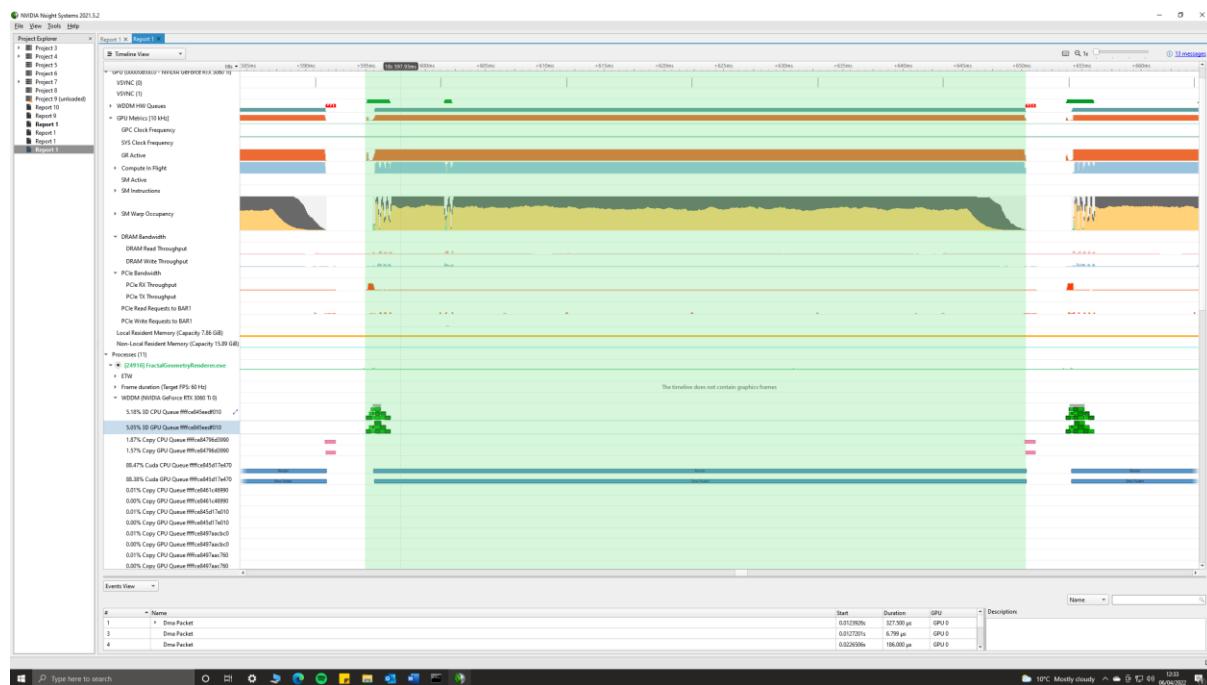


Figure 4.4.1 Nvidia Nsight Systems output for Mandelbulb benchmark

In addition, the Windows 10 task manager displays the GPU compute usage as very close to 100% for the Mandelbulb benchmark, as shown in Figure 4.4.2.

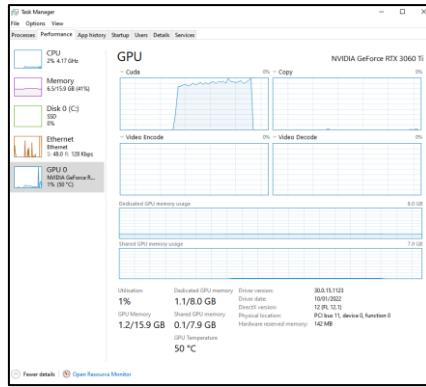


Figure 4.4.2 Windows GPU usage for Mandelbulb benchmark

While the application makes use of almost 100% of the available CUDA computation cores, it uses a very small amount of GPU memory, as it only requires memory for the output pixel data buffer, and any data that overflows from a work item's private memory. The application doesn't make use of any device specific features, such as Nvidia's new ray tracing cores, as these are not accessible through the OpenCL API. OpenCL is designed to be cross platform and cross device, meaning that only features that are available on all GPUs can be accessed through the API. It could be worth considering rewriting the application to use the Nvidia CUDA API, as this would allow the ray tracing API to be accessed, potentially giving better performance as more device features would be utilised.

4.4.3 Unmeasurable Features

There are several features of note provided by the application which aren't measurable but are still significant reasons to use the application. The first and most significant, is that every point that is sampled using the distance estimation function can have its own material. This means that every pixel rendered on the screen can be rendered using a different material, allowing geometries to blend between colours seamlessly. This gives a massive amount of flexibility and is incredibly powerful as it allows geometry to behave more like real world objects, which contain impurities and aren't made from one solid material. This feature is not common among rasterization rendering, which requires that geometry be grouped using the material by which it should be rendered.

The second feature of note is that there is no performance cost for moving geometries or lights around the scene. All lighting and pixel data calculated by the application is computed in real time and is computed every frame, even if the camera, light, or geometry hasn't moved. This does leave room for future improvements, but real time lighting is a benefit of using ray tracing over rasterization rendering, which commonly "bakes" lighting data into a texture for later use. This makes lighting calculations very cheap for rasterization rendering as all data is precomputed, but also means that it is difficult for lighting to be updated in real time. Ray tracing completely negates this problem by making all lighting update in real time, which requires more computation but gives more realistic results. While the two features discussed in this section are not measurable, they are significant reasons to use this application.

4.4.4 Measurable Features

This section evaluates the features of the application that can measured, by calculating their performance cost and comparing this against the visual benefits of that feature.

4.4.4.1 *Surface Shading*

The Blinn-Phong surface shading model is a physically based rendering method which combines ambient colour, diffuse colour and specular highlights to efficiently approximate how real-world materials behave under light. This is a common surface shading model that is computationally cheap while giving a good visual result.

4.4.4.2 *Glow*

A glow can be applied to the space outside of an object. This is calculated for each pixel by keeping track of the minimum distance that this ray was from the scene geometry, and if this ray didn't collide with geometry, then the glow strength should be calculated using $\frac{\text{minimum distance to geometry}}{\text{maximum glow distance}}$.

4.4.4.3 *Shadows*

Both hard shadows and soft shadows have been implemented and either can be added to a scene. To calculate if a point on a geometry is in shadow, a ray must be cast from the point

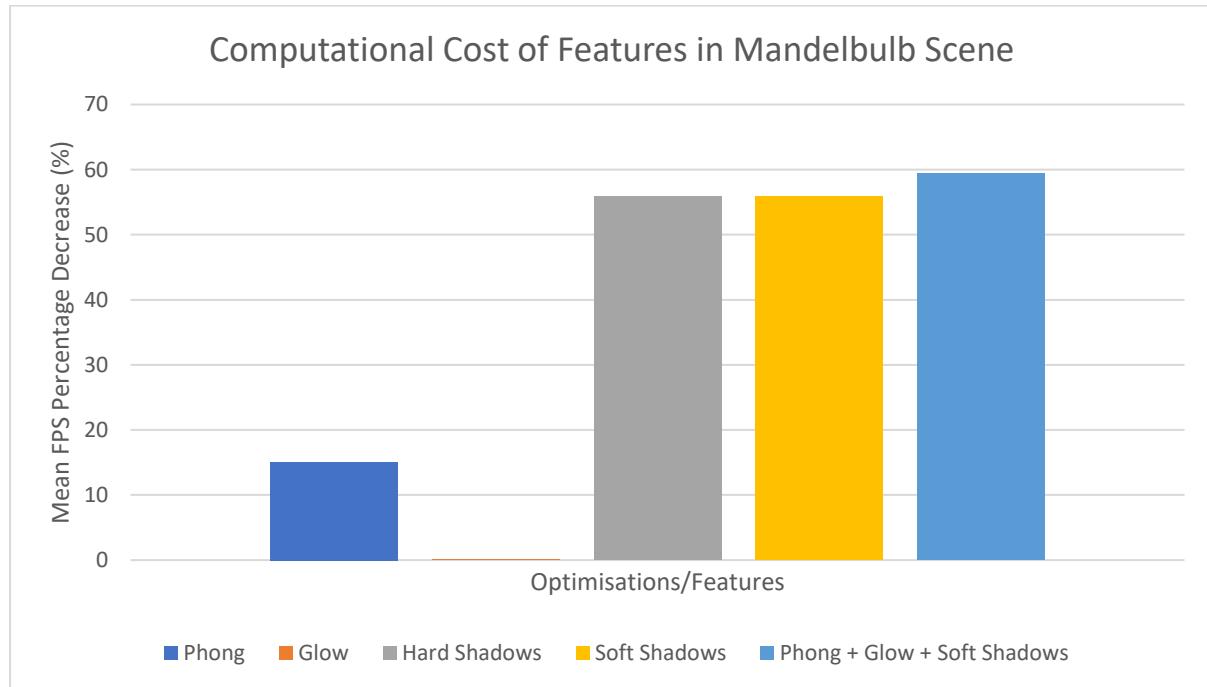
towards the main light in the scene, and if the ray collides with any other geometry, then the original point should be in shadow. Hard shadows mean that a point is either in shadow or not, which makes the shadows have very harsh edges and gives a stylised, cartoon look. Soft shadows use a falloff value to determine the distance used to blend between shadow and no shadow, which gives more realistic shadows at a higher computation cost.

To generate accurate shadows, the main scene distance estimation function is called for each position on the ray, however, this could be optimised to use the bounding volume distance estimation. This was not implemented as incorrect shadows could be introduced accidentally by allowing the bounding volume to cast shadows on geometry.

4.4.4.4 Evaluation

The performance of the four features discussed above were benchmarked using the Mandelbulb benchmark scene, where scene was run three times and an average of the raw data was taken.

The results are displayed in Graph 4.4.4.



Graph 4.4.4 Computational Cost of Features in Mandelbulb Scene

The Blinn-Phong feature saw a performance decrease of only 15.1%, which in most cases is well worth it considering how substantial the visual payoff is. Figure 4.4.3 compares basic ambient shading to the Blinn-Phong surface shading model, note the darker areas where no diffuse light

reaches, and the specular highlights created by the light source. In most cases, the Blinn-Phong surface shading technique is worth the performance cost considering how substantially it increases the visual quality of the render.

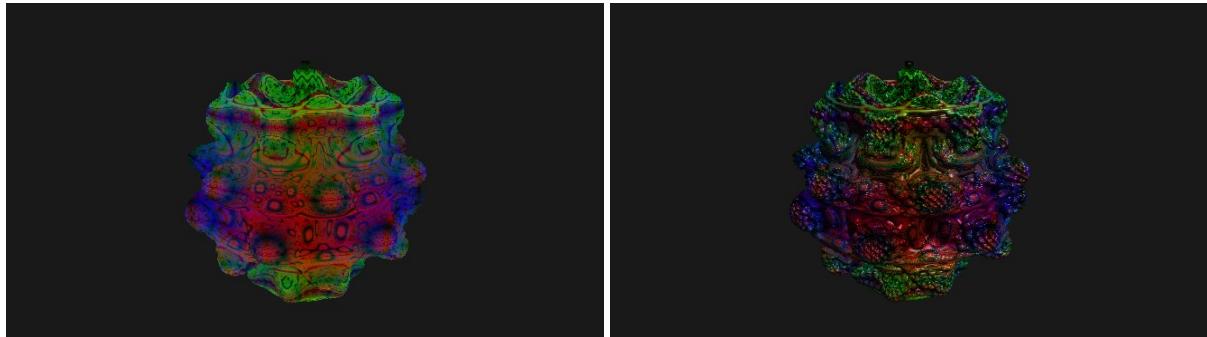


Figure 4.4.3 Basic shading (left) and Blinn-Phong shading (right)

The glow feature resulted in a 0.03% performance decrease. The computation cost of this feature is negligible and is essentially free for the ray marching algorithm and can produce a nice effect in some scenes. While this feature is cheap, it doesn't look good in every scene which is why it is disabled by default. Figure 4.4.4 displays the geometry glow effect on the Mandelbulb scene.

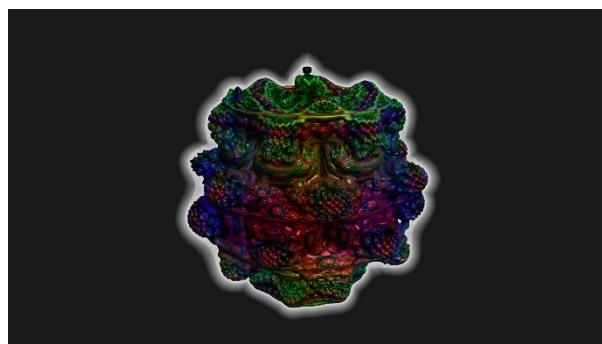


Figure 4.4.4 Geometry glow applied to Mandelbulb

The use of hard shadows decreased the average frames per second by 55.8%, which is a substantial performance loss. This performance decrease comes from the additional ray that must be cast for each pixel of a geometry. The performance of this feature is affected by both the computational cost of the scene distance estimation function, and by the distance between the geometry surface and the scene light. If the light is far away from the geometry, then it will take more iterations of calling the distance estimation function for the ray to reach the light, and therefore will take more computation time.

Soft shadows performed similarly to hard shadows, resulting in the same mean FPS decrease to hard shadows, at 55.8%. The computation required to calculate soft shadows is almost identical hard shadows but does require more operations and temporary variables. It would be expected for soft shadows to perform slightly worse than hard shadows, but in this case inconsistencies between runs makes up this difference. Figure 4.4.5 compares hard shadows and soft shadows, note the smoother blending between shadow and light areas for soft shadows.

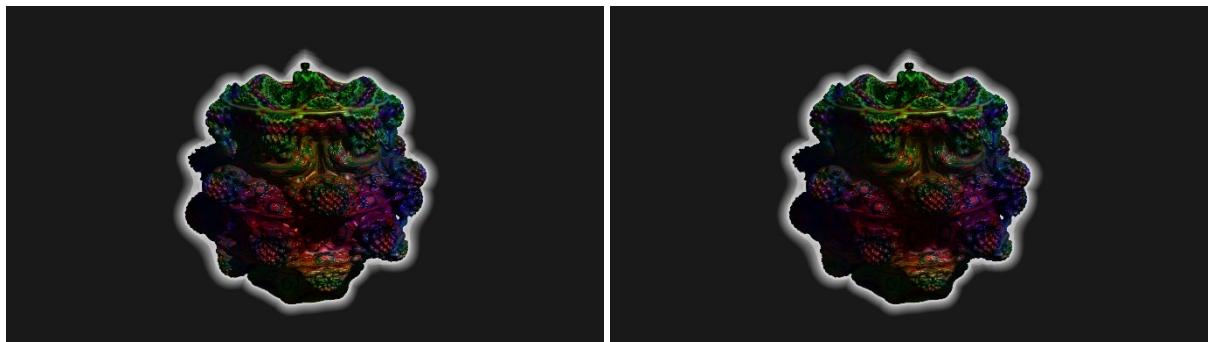


Figure 4.4.5 Hard shadows (left) and soft shadows (right)

Interestingly, when combining Blinn-Phong, soft shadows, and glow, the performance decrease is significantly less than the sum of their component parts, which is an unexpected result. It is expected for optimisations to perform better when combined due to their synergistic nature but combining features should not have the same effect. This inconsistency is likely due to caching, which occurs every time the value of a variable is required. The calculation of shadows requires several variables that are also used when calculating the material, most notably a struct containing the scene light information, and the current position of the ray. Since the calculation of shadows occurs directly after the computation of the material value, these shared values should already be in the work item cache and so can be accessed very quickly.

While the visual output of the application is aesthetically pleasing, it doesn't look anywhere near as good as an offline render of a fractal as many compromises have been made to ensure that the application can render the image in real time. However, in the future, this may not always be the case as Moore's law observes that the number of transistors in a dense integrated circuit (such as a CPU or GPU) doubles approximately every two years, and therefore performance also approximately doubles. While the upper limit for the number of transistors in CPUs and GPUs

has likely been reached, the relationship of doubling performance roughly every two years is almost still being achieved, through increasing the parallel degree of chips and, therefore, allowing more operations to be completed per second [44]. If this trend keeps up, higher visual quality real time rendering will be achievable in the next couple years, though it is likely that offline renders will always create a superior visual output.

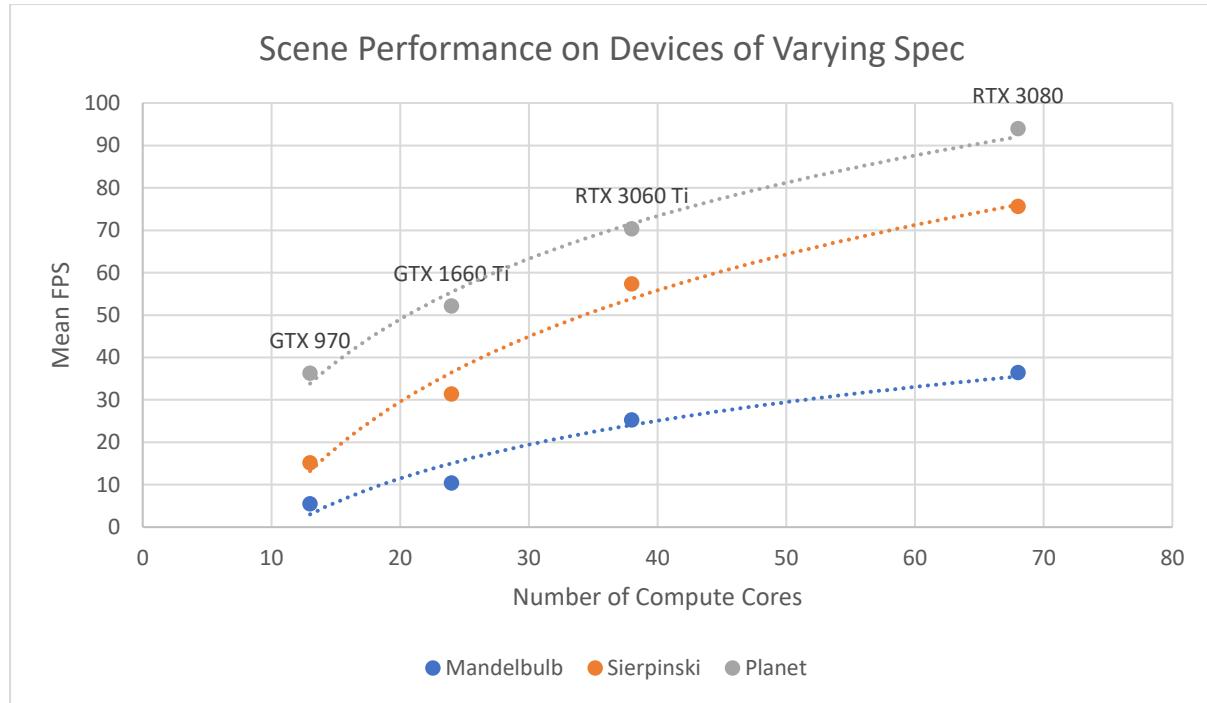
4.4.5 Application Scalability

The benchmark scenes were profiled across several different systems to determine how well the application scaled over systems of varying performance and age. Table 4.4.1 lists the graphics cards used by the benchmarking systems and the age of their respective architectures. The devices tested against contains variety of current generation and older generation Nvidia GPUs. Unfortunately, no devices containing AMD or Intel graphics were available to benchmark the application on. Windows was the only operating system that the application was tested on, and Windows 11 is still in early access and is not yet fully supported, so performance of devices using this operating system may not be accurate.

Table 4.4.1 Graphics cards used by the benchmarking systems

Graphics Card	Date Released	Price Range	Operating System
NVIDIA GeForce GTX 970	2013/2014	Mid-tier	Windows 10
NVIDIA GeForce GTX 1660 Ti	2018/2019	Mid-tier	Windows 11
NVIDIA GeForce RTX 3060 Ti	2020 (current generation)	Mid-tier	Windows 10
NVIDIA GeForce RTX 3080	2020 (current generation)	High-end	Windows 11

The core scenes were benchmarked on the systems listed in Table 4.4.1 and results plotted in Graph 4.4.5. The axis “Number of Compute Cores” refers to the total number of work groups (a collection of work items) that each graphics card contains. This is the parallel degree of the graphics card, which has increased substantially over the last decade.

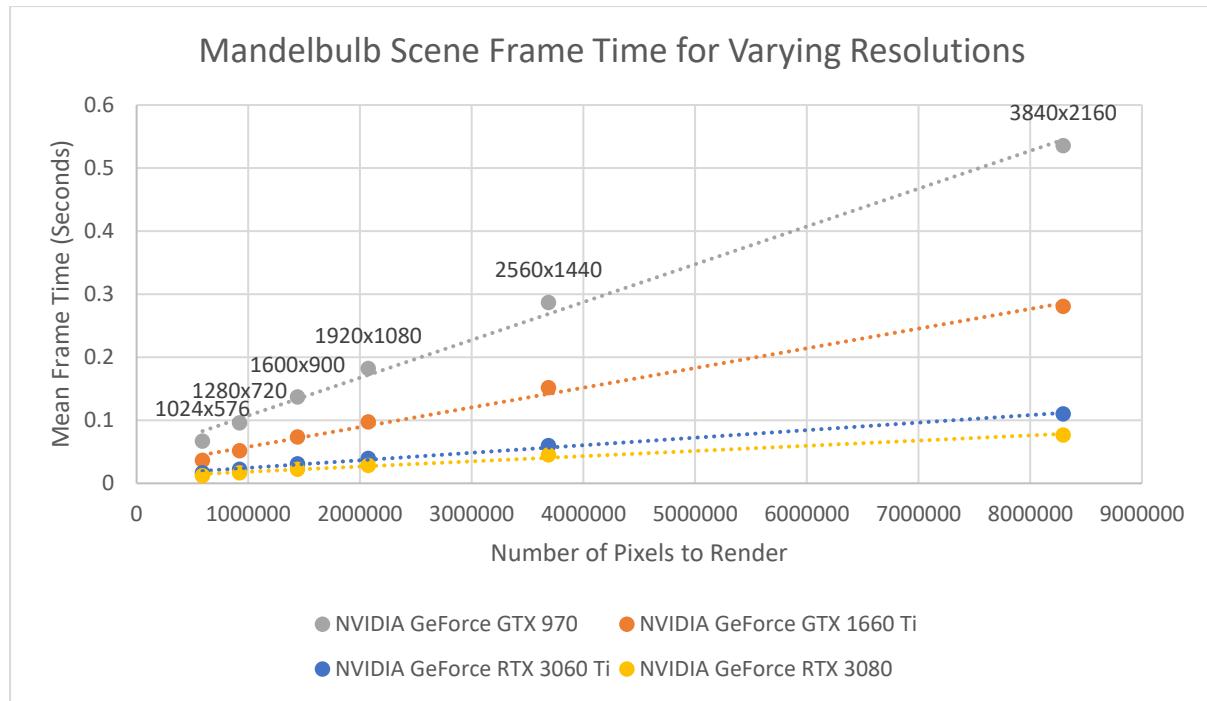


Graph 4.4.5 Scene Performance on Devices of Varying Spec

The mean frames per second increases logarithmically as the number of compute cores increases linearly. This result is expected as per Amdahl's law, which states that all performance improvements are limited by the amount of time that the improved section is spent in execution. This means that because it will always take at least 3ms per frame to draw pixels onto a window for a resolution of 1920×1080, the total computation time for one frame can never fall below 3ms. This is a hard limit for performance and is the theoretical best performance achievable.

To determine the point at which this limit is reached, further benchmarking would be required on much higher power systems. Benchmarking of the application using a high performance system such as the Heriot-Watt Robotarium cluster [45] or Google Colaboratory [46] was considered, but was decided to be out of scope for this project and instead is discussed in section 5.3 Future Work.

The performance of the application over various resolutions was also benchmarked, and results displayed in Graph 4.4.6.



Graph 4.4.6 Mandelbulb Scene Frame Time for Varying Resolutions

As expected, the average time taken to render a frame increased linearly as the number of pixels in the frame increased linearly. This result is important as it means that upscaling is a valid optimisation improvement that can be made. The resolution can easily be halved and image upscaled back to the original resolution resulting in significantly better performance, at the cost of reduced detail. In most cases this would work very well as ray tracing creates images which have more than enough detail, meaning they can be upscaled and still retain good quality.

Interestingly, when doubling the number of pixels to render on the screen, the computation time increases to less than double the original time. Intuitively, the opposite would be expected, however, when taking into consideration that for larger resolutions, work groups will consist of more pixels taking a similar computation time (since the pixels will be sampling points much closer to each other than before), this means that the amount of work item stalling is much lower, allowing the GPU to achieve higher a throughput than before. Additionally, the performance benefits of caching will be more significant since more items will be accessing data.

To summarise, the application scales very well when increasing the parallel degree of the graphics card and is consistent with the expected result. The application also scales well when

increasing the resolution of the application, which makes reaching high resolutions feasible once graphics cards have become more powerful.

4.5 Evaluation of Requirements Specification

The full project requirement specification can be found in appendix 6.9. This specification only received small changes since the initial research report as thorough project planning ensured that all requirements were realistically achievable and within a reasonable scope. Requirements were prioritised using the following strategy:

- MUST – a requirement that is of the highest priority to the project
- SHOULD – a requirement that is not essential, but would be valued
- COULD – a requirement that is optional

All thirteen of the MUST priority requirements, and the one and only SHOULD priority requirement were fully implemented. Of the five COULD priority requirements included as stretch goals to be implemented if good progress was made with the project, two were fully implemented, one was partially implemented and two were not implemented.

In terms of requirements implemented, this project was definitely a success as 100% of the MUST and SHOULD priority requirements were implemented, and 40% of the optional COULD requirements were completed.

4.6 Evaluation of Aims and Objectives

A goal-based evaluation strategy has been used to determine if the projects aim and objectives have been achieved. For an objective to be considered achieved, all requirements prioritised MUST and SHOULD related to that objective must have been implemented. For the project aim to be considered achieved, all objectives must be achieved.

The purpose of objective one was to complete background research into the project area to gain a better understanding of the chosen topic and the scope of the project. This objective was

mostly completed when creating the initial research report, though it did run for the entire duration of the project to ensure that the project stayed up to date with current developments.

The purpose of objective two was to investigate existing solutions and to analyse their strengths and flaws. This objective was also mostly completed when creating the initial research report, though popular fractal communities were monitored so that any new fractal rendering software could be investigated. No new fractal rendering software of note was released during the duration of this project.

The purpose of objective three was to implement the core functionality of a non-real-time 3D fractal renderer, which formed the safe core of the project. This objective was linked directly to the MUST priority requirements FR-1.1, FR-1.2, FR-1.3, FR-1.4 and FR-1.5, all of which were fully implemented. Therefore, objective three has been achieved as all requirements associated with the objective have been implemented.

The purpose of objective four was to add additional functionality to the renderer, making it capable of real-time rendering by adding a game-loop, adding a controllable camera, making scenes dynamic, and adding optical effects and lighting. The scope of this objective was left flexible so that it could be increased or decreased as necessary, and several stretch goals were included in the requirements specification in appendix 6.9. The requirements associated with this objective were FR-2.1, FR-2.2, FR-2.3 and FR-2.4, all of which were MUST priority and fully implemented. Additionally, two of the COULD priority requirements representing the project stretch goals were fully implemented, FR-2.7 and FR-2.8, though the remaining three stretch goals FR-2.5, FR-2.6, and FR-2.9 were not fully implemented, though some progress had been made towards one of them. Since all MUST priority requirements associated with this objective were fully implemented, and good progress was made implementing additional stretch goals, objective four has been achieved to a large degree.

The purpose of objective five was to benchmark the application performance across various systems. This objective was tied to the MUST priority requirement FR-3.1 which was achieved.

Additionally, the application performance is evaluated in section 4.4, making this objective achieved without a doubt.

The purpose of objective five was to create documentation to aid users and developers of the application. This objective was not tied to any requirements but was still achieved and is located on the [online documentation page](#) [6].

The project objectives were created to help guide the project in the correct direction by splitting up the work to be completed into several key chunks. All six of the project objectives have been achieved without a doubt.

The project aim was to develop an application which can update and render 3D fractal geometry in real-time. This was undoubtedly achieved. Additionally, the aim specified that the render should use common optical effects including the Blinn-Phong surface shading method, lighting, hard and soft shadows, and geometry glow, which as discussed in section 4.4.4 have been fully implemented and evaluated. Finally, the aim also specified that it should be possible to create scenes and view them using the application, and this process should be as straight forward as possible. This part of the aim is more subjective than the previously discussed parts, as some users will find it more straight forward to create and view a scene than others, though for the most part, this process is as simple as possible for the current setup. In the future it may be worth conducting a usability study to quantitatively ascertain if this part of the aim was achieved.

In summary, all six of the project objectives have been completed and several stretch goals were also completed. The project aim has been achieved and additional constraints mostly met, though in the future a usability study of the user interface would be valuable. Overall, the project has undoubtedly completed the aims and objectives that it set out to achieve.

4.7 Comparison with Literature Results

Table 4.7.1 displays a matrix comparing the work completed by this project to results from existing pieces of literature. The data should be compared with caution, as there are many

differences between the projects in terms of work completed and evaluation strategy, which make a side-by-side comparison incorrect. Some values are missing from the matrix, most notably the render times for two of the real time rendering solutions.

The most recent significant piece of work completed in this field is Realtime Rendering of Complex Fractals [14], a book chapter from Ray Tracing Gems II, published by Nvidia and Apress. This project implemented the Mandelbulb and Juliabulb fractals using HLSL (High-Level Shader Language) which is the language used by DXR, the DirectX Raytracing API. Unfortunately, this project didn't report any performance results.

Rendering Methods for 3D Fractals [16] is a bachelor thesis from 2010, which investigated two main ways of rendering 3D fractals – a real-time voxel-based volume renderer, and an offline super sampler solution. Of the two solutions, the real-time volume renderer is most similar to this project, as it could render the Mandelbulb fractal in real time. However, this solution didn't render a fractal with large amounts of detail, instead a voxel grid was computed for the fractal and this data was rendered instead. The voxel grid was recomputed as the user zoomed into the fractal giving the illusion of infinite detail. This solution used GLSL (OpenGL Shading Language) with OpenGL. Unfortunately, this thesis did not report runtime results for the real-time solution.

Finally, Ray Tracing Deterministic 3-D Fractals [47] is a journal article published 1989 and was the first piece of work to propose how to render 3D fractals efficiently using distance estimation and signed distance fields. This was an astonishing piece of work at the time which significantly furthered the field of fractals and computer graphics, and without it, real time rendering of 3D fractals would not be possible.

To summarise, there have been many projects which set out to render images of 3D fractals, but only a few which aimed to render them in real time. Unfortunately, these projects did not report their performance runtimes which makes a side-by-side comparison impossible. Installation of these programs on the development device should be considered for future work.

Table 4.7.1 Comparison matrix for literature results

Project	Project Type	Year	Implementation Language	Real time or offline rendering	Processor Used	Device	Render Time
Realtime Rendering of 3D Fractal Geometry	Honours project	2022	OpenCL	Real-time rendering	GPU	Nvidia RTX 3080	Average of 28ms per 1920×1080 frame, for Mandelbulb benchmark scene
Realtime Rendering of Complex Fractals [14]	Book chapter	2021	HLSL for DXR	Real-time rendering	GPU		
Rendering Methods for 3D Fractals [16]	Bachelor thesis	2010	Voxel renderer using GLSL with OpenGL	Real-time rendering and voxel calculation	GPU	AMD Turion X2 Ultra-Dual-Core Mobile ZM-82 2.80 GHz processor 3 GB RAM ATI Radeon HD 3200	Approximately 30 minutes to render a one-megapixel image
			Super sampler using C++ and Devil image library	Offline rendering	CPU		
Ray Tracing Deterministic 3-D Fractals [47]	Journal article	1989		Offline rendering	CPU	AT&T Pixel Machine containing 64 parallel processors, each running at about 10 MFLOPS	Approximately one hour to render a 1280×1024 image

4.8 Future Work

4.8.1 Features

There are many visual features that could be added to the application. Additional surface shading methods like Lambert, Oren-Nayar, or a subsurface scattering algorithm could be implemented. Additionally, ambient occlusion would be a valuable addition and the scene background could be converted from a static colour to a gradient colour easily.

The application would benefit visually from the addition of transparency and reflections, however, these features are costly as they both require additional rays to be cast. It is likely both features would have a performance cost comparable to shadows, which resulted in a roughly 50% drop in mean frames per second. For any of these visual features to be viable, the application performance would first have to be significantly improved.

A graphical user interface could also be valuable for the application. Currently, a GUI is not needed as there are not any features that would make use of it, however, if more features are added in the future, then a GUI may could be considered.

4.8.2 Algorithmic Optimisations

There are several small and easy to implement optimisations that would result in a moderate performance boost, such as optimising the shadow ray cast to use the geometry bounding volume, and to use an alternate method for calculating the surface normal of a geometry, which uses four calls to the scene distance estimation function instead of six [48].

A more substantial optimisation would be to make use of an acceleration structure for storing the geometry in the scene. An acceleration structure is a data structure (usually some kind of tree) which stores all objects in a scene. This means that when calculating ray intersections, the tree can be iterated over to find the collision, instead of all objects in the scene. In the best case, this decreases the complexity of searching for the object to $O(\log N)$, down from $O(N)$. However, in practice, it is hard to reach this theoretical performance as the scene must be partitioned efficiently. Additionally, a performance gain is only seen when the scene contains

many objects, which our application does not. Finally, the OpenCL C kernel language would be unpractical for implementing this optimisation as it does not feature classes, and work items cannot access a shared state. This optimisation may be beneficial in some cases but would require rewriting most of the application code.

Progressive rendering could also be considered. This is the process of splitting the rendering of a frame into multiple parts, each part reliant on the previous part. This has the benefit of allowing the renderer to fall back to a lower quality render if the performance is sub optimal, and to increase quality if performance is good.

Additionally, level of detail optimisations could be implemented which swap out highly detailed geometries with less detailed geometries when the camera is far away from the object, as this extra detail is not seen.

4.8.3 Design Optimisations

An important optimisation to make would be to reduce the duplication of pixel data as this is an unnecessary waste use of memory. Additionally, it could be worth changing colour format as the alpha channel is completely unused by the application.

It may also be worth investigating other graphics libraries and finding one that would allow the pixels to be displayed asynchronously using the CPU, which could be completed while the application is using the GPU to compute the next frame. This would result in quite a significant performance gain as it takes approximately 3ms to render a 1920×1080 frame to the window using the development device.

Additionally, some of the sequential code in the application could be replaced with concurrency, such as when reading the user input. However, this would require a significant rework as threads would now be required to synchronise before computing a frame, and the performance gain would be minimal, but possibly worth it as it would reduce the amount of sequential code in the application, therefore increasing the theoretical maximum FPS as per Amdahl's law.

Finally, it could be worth implementing upscaling in the application which would give a significant performance gain as the application could render a lower quality image, and then upscale it to the desired resolution.

4.8.4 Fractals

There are many more fractals that could be implemented as scenes. Some of which include the Julia bulb family, kaleidoscopic fractal family, and IFS (Iterated Function System) fractals.

4.8.5 Evaluation

More sophisticated evaluation of the application could be performed. This could include performance benchmarks on different commercial GPUs such as AMD or Intel, or even on high-performance systems such as the Heriot-Watt Robotarium cluster [45] or Google Colaboratory [46]. The performance of the application on different architectures may vary significantly, due to differences in data type and memory pointer sizes and differences in the implementation of the OpenCL API.

Additionally, it could be worth benchmarking the application using some ray tracing hardware, though the application renderer might have to be rewritten using a different language to make use of platform specific API features, such as the CUDA ray tracing API. This would likely result in significantly better performance but would require thorough research and planning and a significant rework of the application.

Finally, a more thorough evaluation of the project in respect to existing work could be completed. This would be a complex task, which would involve installing existing solutions from literature on the development machine, creating a standardised scene which can be rendered in all applications, rewriting that scene in every language required so that it can be run on all programs, and benchmarking the performance of the scene for all the applications. With these results, the performance of this application could be compared directly to the performance of existing solutions, which would allow a definitive analysis of the value of this application.

5 Conclusion

5.1 Achievements

The aim of this project was to develop an application which can update and render 3D fractal geometry in real-time. This has undoubtedly been achieved as the application can render the Sierpiński cube and tetrahedron, or the Mandelbulb fractal in real time, running at a resolution of 1920×1080 . The render uses the Blinn-Phong lighting and surface shading technique, hard or soft shadows, and geometry glow. Additionally, users can create their own scenes using the library functions provided by the application, and scenes can be viewed in the application by passing the scene path as a command line argument.

All six of the project objectives have been achieved. Objectives one and two were met in section 2, which evaluated relevant background literature surrounding the project area, accompanied by a critical analysis of existing solutions. Objectives three and four specify the functionality that the application should provide in terms of software requirements, which are evaluated in section 4.5. Objective five specifies the evaluation to be performed on the application and project, which is completed in section 4. Finally, objective six specifies the documentation to be created to aid users and developers, which is located on the [online documentation page](#) [6].

What is the value of the application, where would it be used?

5.2 Limitations

At its core, the application is only a basic implementation, providing limited functionality. Most currently available fractal viewing applications provide more visual features than just surface shading, shadows, and glow, and contain many more example scenes. Additionally, the application was only tested using several fractals - the Mandelbulb and Sierpiński cube and tetrahedron, and it was only tested on Windows 10 and 11 using Nvidia GPUs.

5.3 Future Work

A thorough discussion of possible future improvements and features for the application is included in section 4.8, which covers five main categories, features, algorithmic optimisations, design optimisations, fractals, and evaluation.

There are many visual features that could benefit the application, such as additional surface shading techniques, ambient occlusion, transparency, and reflections. However, these features are costly, and the application would need to be optimised further before these could be implemented. Algorithmically, there are several simple optimisations that could be made, such as modifying the shadow calculation to use bounding volumes and to use a cheaper alternative method for calculating the geometry surface normal. The application design could also be improved by removing the duplication of pixel data in main memory and by using an alternative library for updating the display window asynchronously. Additionally, some sequential code could be made concurrent and image upscaling could be considered. There are many more fractals that could also be implemented, and more sophisticated evaluation of the application could be completed using different hardware and operating environments.

5.4 Final Thoughts

This project was a proof of concept which proved that it is possible to create a real-time renderer for 3D fractals. The project was completed using low-level languages, in the hopes that this would give ideal performance, however, this highlighted the effort required to create a high-performance application using low level programming languages. Looking to the future, this field seems to be moving towards using higher level languages, which make use of compilers to efficiently optimise code, probably doing a better job than a developer would by hand. This is a very exciting time, as real-time rendering for 3D fractals is becoming easier and hence more common, and I expect many new pieces of work to be completed in the coming years.

6 Appendix

6.1 Experimentation Renders

The images in Figure 6.1.1 were rendered using a very early version of the application, during the research and experimentation stage.

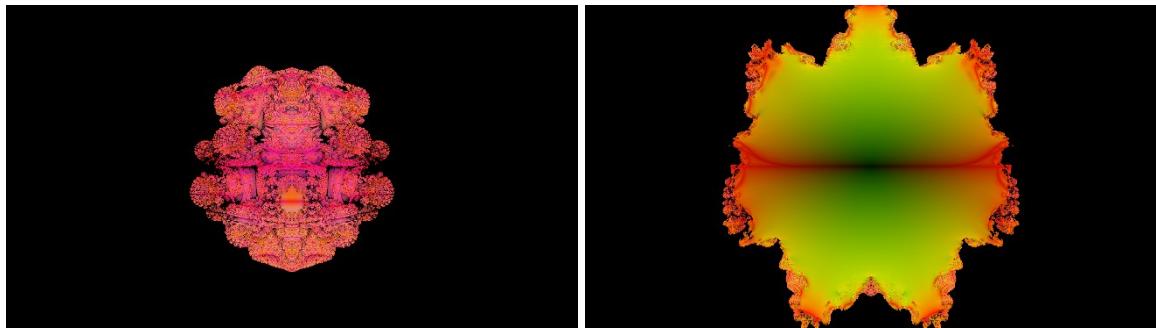


Figure 6.1.1 Render of the Mandel bulb fractal (left) and cross section (right) using equation from [49]

The value of a surface normal can be converted to a colour by mapping the x, y, and z values to r, g, and b, which can be useful when debugging. The images in Figure 6.1.2 show the surface normals of the sphere and box experiment scene, and an example phong shading, a surface shading technique.

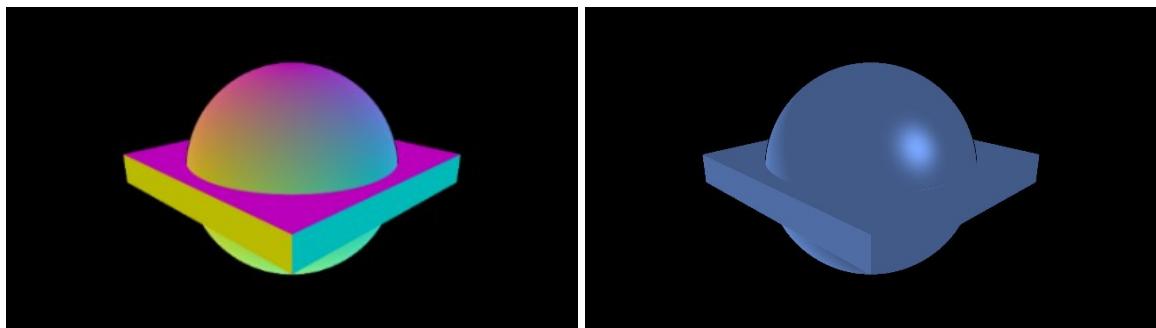


Figure 6.1.2 Surface normal visualised (left) and phong shading experiment (right)

6.2 Sphere SDF Example

Included is the signed distance function (SDF) for a sphere with radius R , positioned on the origin.

$$\text{sphereSDF}(p) = |p| - R$$

where $p = \begin{bmatrix} x \\ y \\ z \end{bmatrix}$, $|p|$ is the magnitude of the vector p , R is the circle radius in world units

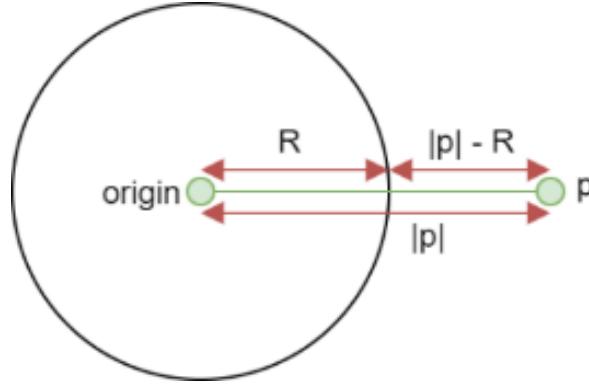


Figure 6.2.1 Sphere SDF diagram

6.3 Smooth SDF Combinations

The following equations can be used to combine two SDF values using a smoothing value.

$$\text{smoothUnion}(a, b, s) = \min(a, b) - h^2 \times \frac{0.25}{k}$$

where $a, b \in \mathbb{R}, s \in \mathbb{R}$ is the smoothing value, $h = \max(s - \text{abs}(a - b), 0)$

$$\text{smoothSubtraction}(a, b, s) = \max(-a, b) + h^2 \times \frac{0.25}{k}$$

where $a, b \in \mathbb{R}, s \in \mathbb{R}$ is the smoothing value, $h = \max(s - \text{abs}(-a - b), 0)$

$$\text{smoothIntersection}(a, b, s) = \max(a, b) + h^2 \times \frac{0.25}{k}$$

where $a, b \in \mathbb{R}, s \in \mathbb{R}$ is the smoothing value, $h = \max(s - \text{abs}(a - b), 0)$

6.4 SDF Surface Normal Calculations

$$\text{normal} = \text{normalise} \left(\begin{bmatrix} DE(p + \begin{bmatrix} e \\ 0 \\ 0 \end{bmatrix}) - DE(p - \begin{bmatrix} e \\ 0 \\ 0 \end{bmatrix}) \\ DE(p + \begin{bmatrix} 0 \\ e \\ 0 \end{bmatrix}) - DE(p - \begin{bmatrix} 0 \\ e \\ 0 \end{bmatrix}) \\ DE(p + \begin{bmatrix} 0 \\ 0 \\ e \end{bmatrix}) - DE(p - \begin{bmatrix} 0 \\ 0 \\ e \end{bmatrix}) \end{bmatrix} \right)$$

where $p = \begin{bmatrix} x \\ y \\ z \end{bmatrix}$, e is the epsilon value, DE is the scene distance estimation function

6.5 Existing Solution Renders

6.5.1 Fragmentarium Renders

The following screenshots were taken when experimenting with Fragmentarium [32], [33].

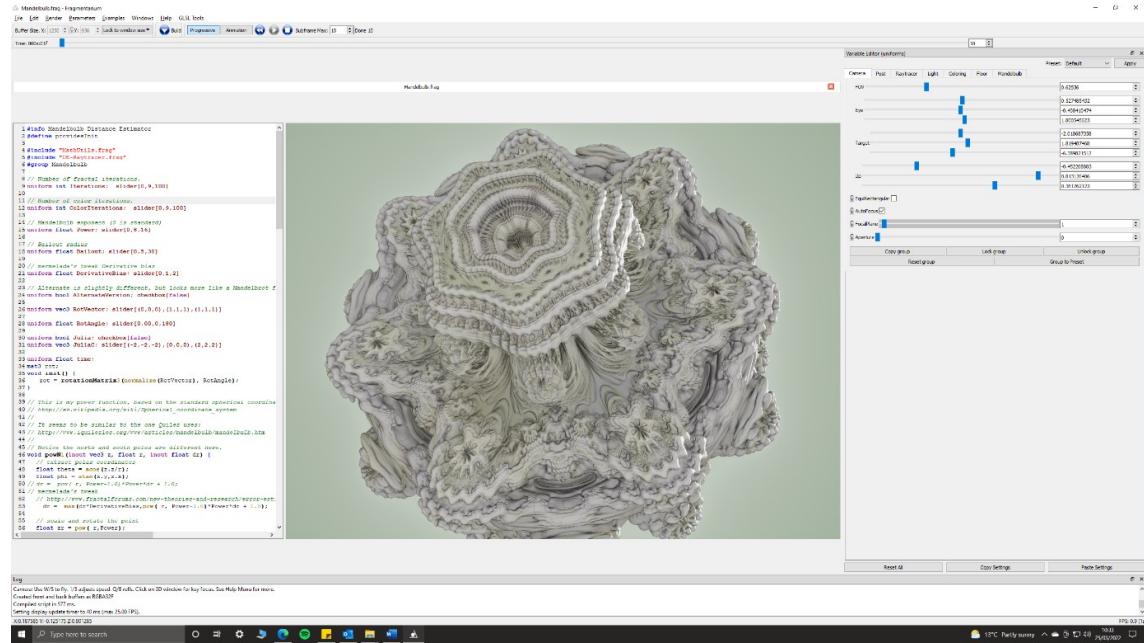


Figure 6.5.1 Fragmentarium interface when viewing the Mandelbulb fractal



Figure 6.5.2 A recursive scene (left) and Mandelbulb fractal (right)



Figure 6.5.3 Tree fractal

6.5.2 Synthclipse Renders

The following screenshot was taken when experimenting with Synthclipse [34].

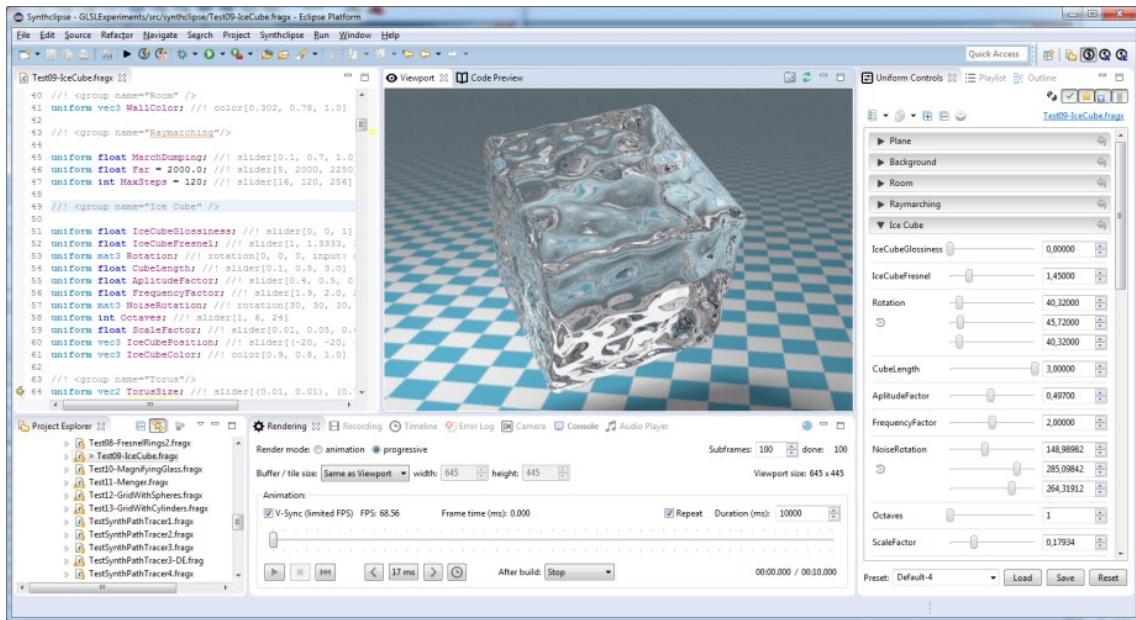


Figure 6.5.4 Synthclipse interface when viewing an ice cube shader

6.5.3 Mandelbulb3D Renders

The following screenshots were taken when experimenting with Mandelbulb3D [35].

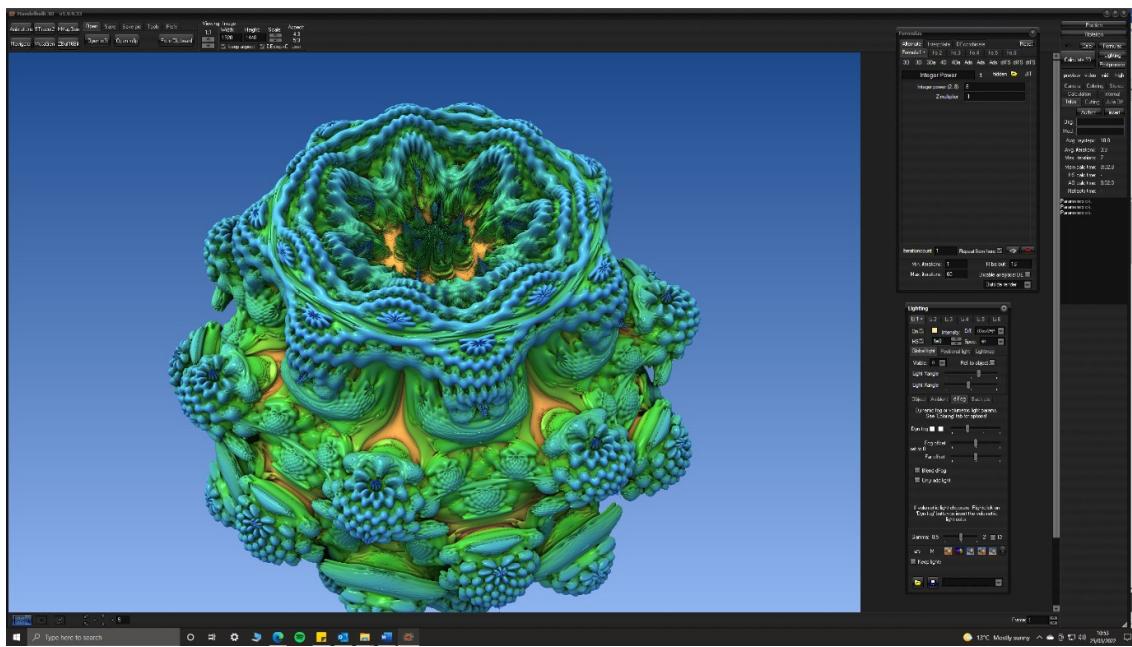


Figure 6.5.5 Mandelbulb3D interface when viewing the Mandelbulb fractal



Figure 6.5.6 Mandelbulb fractal with natural looking colouring

6.5.4 Mandelbulber Renders

The following screenshot was taken when experimenting with Mandelbulber [36].

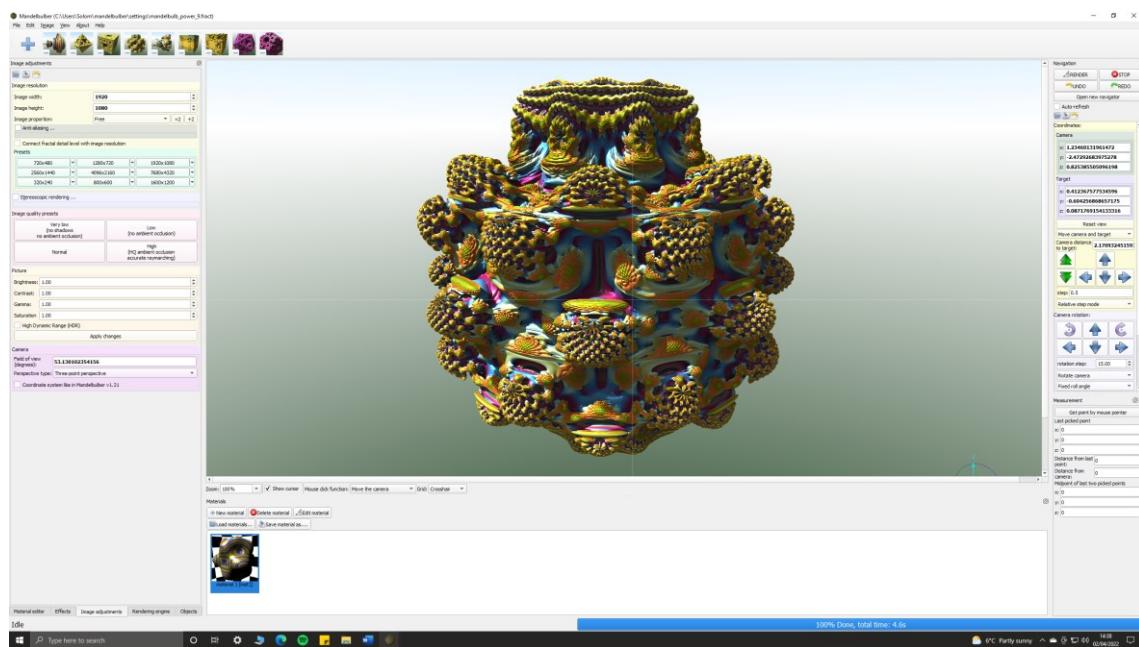


Figure 6.5.7 Mandelbulber interface when viewing the Mandelbulb fractal

6.5.5 FractalLab Renders

The following screenshot was taken when experimenting with FractalLab [37].

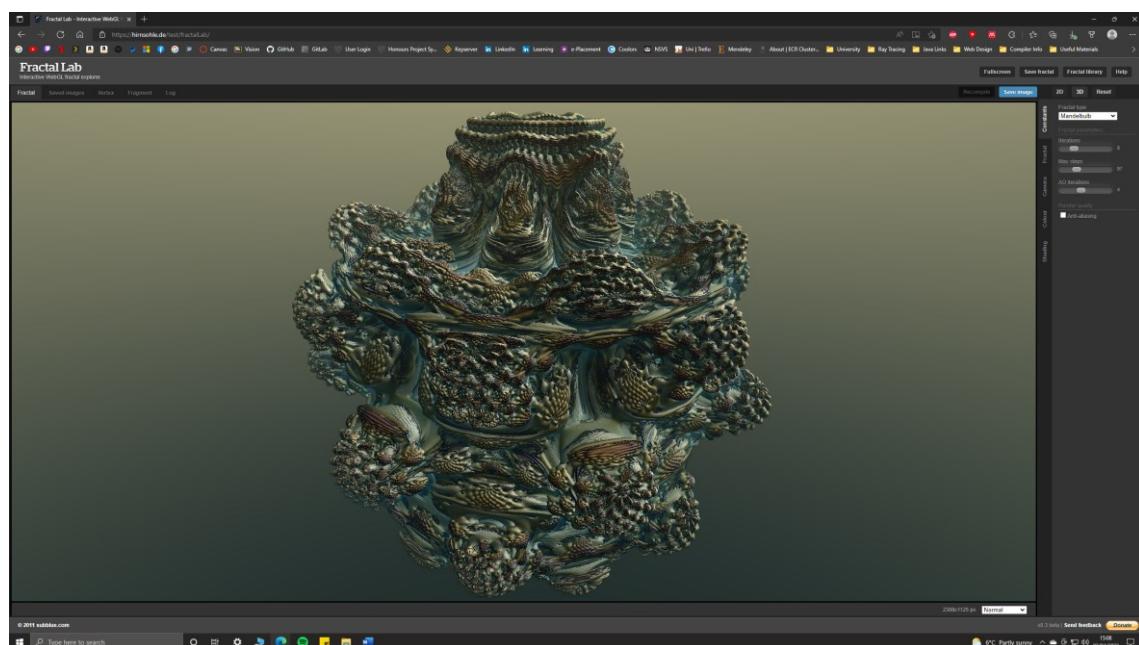


Figure 6.5.8 FractalLab interface when viewing the Mandelbulb fractal

6.6 Application Class Diagrams

Full documentation for the application including class diagrams, dependency graphs, directory structure and namespace structure can be found in the [project documentation page](#) [6].

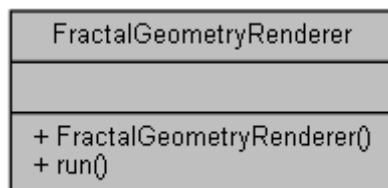


Figure 6.6.1 Application public methods

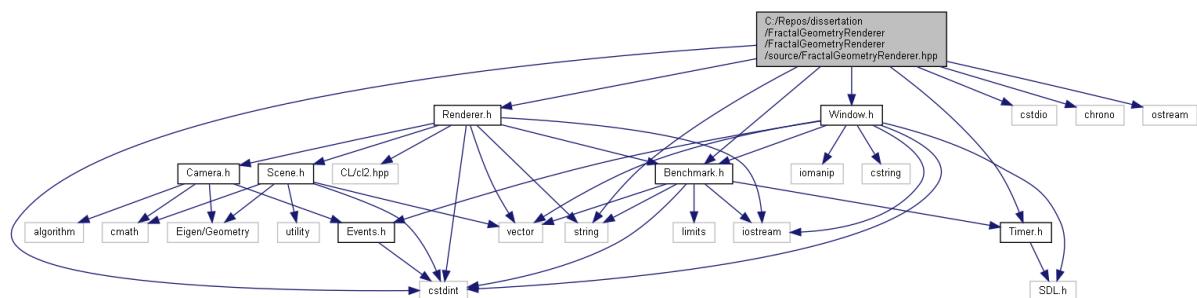


Figure 6.6.2 Dependency graph of the application

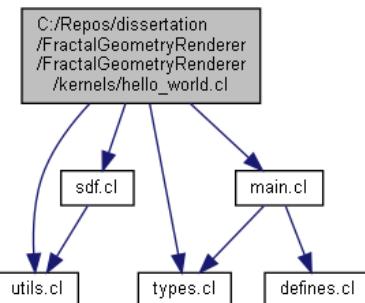


Figure 6.6.3 Dependency graph for example scene hello_world.cl

6.7 Comparison of -cl-fast-relaxed-math Optimisation

No visual differences between -cl-fast-relaxed-math enabled or disabled.

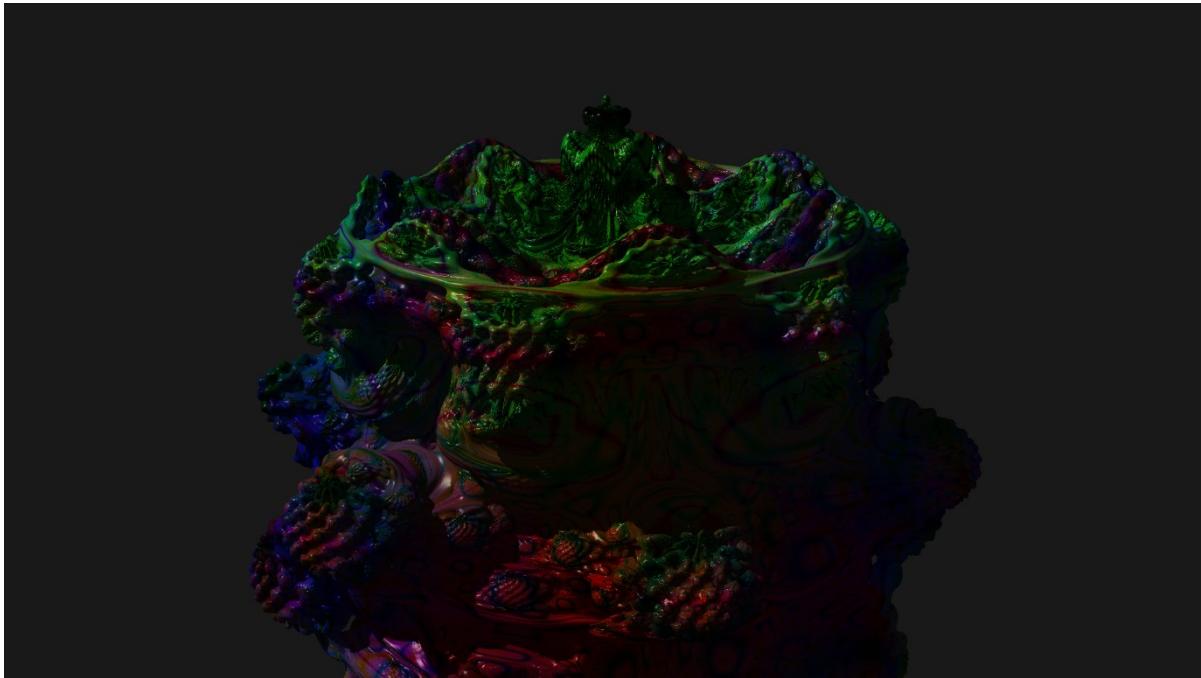


Figure 6.7.1 Mandelbulb.cl scene with -cl-fast-relaxed-math **enabled**

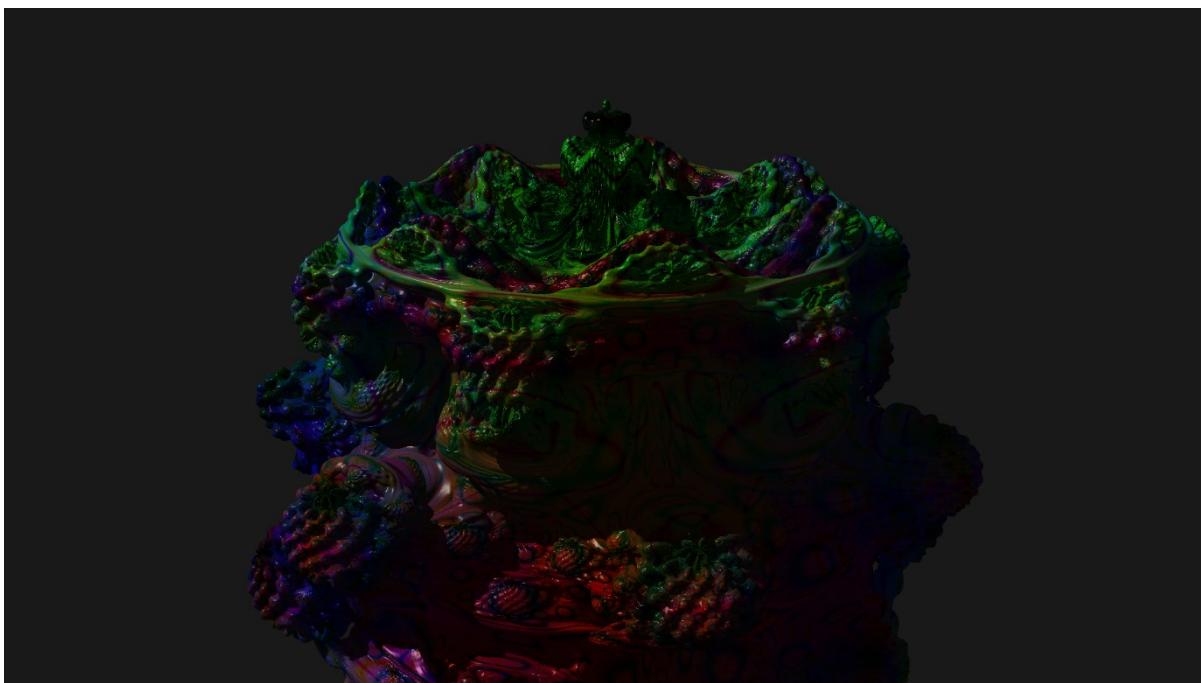


Figure 6.7.2 Mandelbulb.cl with -cl-fast-relaxed-math **disabled**

6.8 Comparison of Glow with Bounding Volume

Figure 6.8.1 displays the Mandelbulb scene with glow effect and compares the visual differences between enabling and disabling bounding volumes **before** the design choice discussed in section 3.2.3 was implemented.

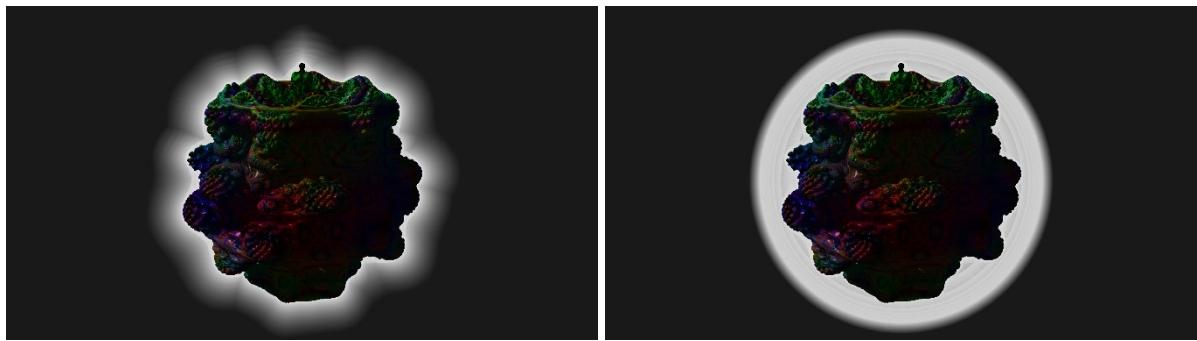


Figure 6.8.1 Mandelbulb scene with glow without bounding volume (left) and with bounding volume (right)

6.9 Requirements Specification

The tables below display the key functionality to be implemented to achieve the projects aim and objectives, specified in section 1.2. Requirements have been grouped by the project objective that they fall under. Requirements were prioritised using the following strategy:

- MUST – a requirement that is of the highest priority to the project
- SHOULD – a requirement that is not essential, but it would be good if the project had it
- COULD – a requirement that is optional
- WON'T – a requirement that would be implemented if the project could run for more time

Table 6.9.1 Functional requirement specification

ID	Name	Description	Priority	Status
FR-1	Objective 3	Core functionality		
FR-1.1	Scene requirements	A scene must contain: <ul style="list-style-type: none">• Geometry• A light• Camera	MUST	Achieved
FR-1.2	Example scenes	The application must contain multiple example scenes, some of which should include: <ul style="list-style-type: none">• Sierpinski tetrahedron fractal• Menger sponge fractal• Julia set fractal• Mandel bulb fractal	SHOULD	Achieved
FR-1.3	Mandatory optical effects	The application must support the following optical effects: <ul style="list-style-type: none">• Ambient occlusion• Hard and soft shadows• Glow	MUST	Achieved
FR-1.4	Load scene	Scenes must be defined each in their own kernel file and the user must be able to load them into the application at runtime	MUST	Achieved
FR-1.5	Settings	Basic application settings must be editable. This includes: <ul style="list-style-type: none">• Output resolution• Optical effects enable/disable• Mouse sensitivity	MUST	Achieved

FR-2	Objective 4	Additional functionality		
FR-2.1	Real-time	The application must be capable of rendering the view of the scene in real time	MUST	Achieved
FR-2.2	Game Loop	The application must make use of a game loop to update the scene. Methods should be called in the following order: <ol style="list-style-type: none"> 1. Poll for user input 2. Update scene 3. Render scene 	MUST	Achieved
FR-2.3	Controllable camera	The user must be able to control a camera, using a keyboard and mouse to move it around the scene	MUST	Achieved
FR-2.4	Dynamic scenes	The contents of a scene (requirement FR-1.1) must be able to move around at runtime	MUST	Achieved
FR-2.5	Optional optical effects	The application could additionally support the following optical effects: <ul style="list-style-type: none"> • Reflections • Depth of field • Transparency 	COULD	Not Achieved
FR-2.6	Optional Surface Shading	The application could support the following surface shading algorithms: <ul style="list-style-type: none"> • Lambert • Oren-Nayar • Phong or Blinn-Phong • Subsurface scattering 	COULD	Partly Achieved
FR-2.7	Fixed camera paths	The application camera could additionally move using a specified camera path	COULD	Achieved
FR-2.8	Image output	It could be possible for the application to output an image of the current camera view	COULD	Achieved
FR-2.9	Video output	It could be possible for the application to output a video sequence for the current scene	COULD	Not Achieved
FR-3	Objective 5	Evaluation		
FR-3.1	Performance benchmark	The application must contain a performance benchmark scene	MUST	Achieved

Table 6.9.2 Non-functional requirement specification

ID	Name	Description	Priority	Testing strategy
NF-1	Executable	The application must run from a compiled executable	MUST	Achieved
NF-2	Operating system	The application must run on Windows 10. Where possible, cross platform libraries should be used, though other operating systems will not be officially supported	MUST	Achieved
NF-3	Display resolutions	The application must support the following common display resolutions: 1366×768, 1920×1080, 2560×1440 and 3840×2160	MUST	Achieved
NF-4	GPU Parallel Computing	The application must run in parallel on the GPU	MUST	Achieved

6.10 Use Cases

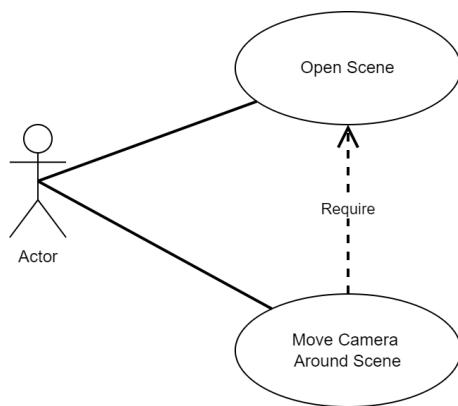


Figure 6.10.1 Use case diagram for the application

6.11 Project Plan

6.11.1 Professional, Legal, Ethical & Social Issues

The main principles of the open-source definition [50] state that a product must have publicly available source code, must allow modifications and derived works of the product, and allow free redistribution of the product. This project is licensed under the GNU General Public License v3.0 [51] which complies with the open-source definition and grants permission for modification, distribution, and commercial use of this product. However, all changes made to the licensed material must be documented, the modified source code must be made public, and the modified work must be distributed under the same license as this product. The license can be viewed on the [project GitHub repository](#) [5], within the LICENSE file.

The British Computer Society (BCS) codes of conduct [52] specify the professional standards expected to be displayed by a member of the BCS, which includes working for the interest of the public, displaying professional integrity, accepting relevant authority, and showing a commitment to the profession. These standards will be respected and complied with for the duration of the project.

This project does not require any ethical approval as no studies requiring participants will be completed. Instead, the application's performance will be benchmarked over several available PCs, and results will be calculated to determine the how successful the project was.

6.11.2 Risk Analysis

A risk analysis for the project has been completed to identify any potential risks that may affect the successfulness of the project, and a mitigation plan has been drawn up for each risk. The mitigation plan has been designed to reduce the chances of each risk happening and to reduce the negative impact of the risk if it were to happen. Each risk was then assigned a rating using the table below.

Table 6.11.1 Risk rating matrix

		Severity		
		LOW	MEDIUM	HIGH
Probability	LOW	LOW	LOW	MEDIUM
	MEDIUM	LOW	MEDIUM	HIGH
	HIGH	MEDIUM	HIGH	HIGH

Table 6.11.2 Risk analysis matrix

ID	Description	Probability	Severity	Mitigation plan	Rating
R-1	Loss of work	LOW	HIGH	All work will be backed up regularly using online version control	MEDIUM
R-2	Change in requirements	LOW	MEDIUM	A thorough requirements specification has been prepared to reduce the probability of this happening In addition, an Agile development approach will be used, which by design minimises the negative effects of changes in requirements	LOW
R-3	Change of deadlines	LOW	HIGH	Communication from the course leaders will be regularly monitored This risk is very unlikely to happen as assurances have been made that the course deadlines will not change	MEDIUM
R-4	Delays due to learning new software	HIGH	LOW	Time was assigned during the planning stage to experiment with new software like OpenCL and GLSL shaders Additionally, some free time has been allocated at the end of the project timetable to allow for delays	MEDIUM
R-5	Delays due to illness	LOW	MEDIUM	Free time has been allocated at the end of the timetable to allow for delays	LOW
R-6	Delays due to bugs	MEDIUM	MEDIUM	Free time has been allocated at the end of the timetable to allow for delays	MEDIUM

R-7	Renderer can't be made real-time	LOW	MEDIUM	Thorough research has been completed with the conclusion that this project is feasible. There is a safe core to the project and room for scaling back or adding extra functionality to the application if this were to occur. In addition, performance analysis of a non-real-time renderer could be performed instead	LOW
-----	----------------------------------	-----	--------	--	------------

The project risks will be regularly monitored, and mitigation of higher priority risks will be prioritised. Each sprint of work will conclude with an analysis of the current state of the project, and from there any risks likely to occur will be identified and corresponding mitigation plans will be consulted.

6.11.3 Project Timeline

The project timeline for the second deliverable has been split into a total of eight sprints of work, each sprint being two weeks in duration. Sprint zero will occur in December during the Christmas holidays, and its purpose is to set up the project working environment and to refactor the existing experimentation code, created when researching and investigating existing solutions. Once this has been completed, Objective 3, the safe core of the project, will almost have been achieved. From there, sprints one and two will be used to achieve Objective 3, to add additional functionality to the application. Then, during sprint three and four, the project evaluation will be completed, and any remaining documentation will be created. Sprints five and six will be used to write up the project evaluation and complete the report, and sprint seven has been allocated as free time, in full expectation that there will be delays when completing tasks and the timeline will have to be pushed back.

A Gantt chart for the project was created using an online tool [53] and has been included below. Tasks have been coloured using the following strategy: documentation (blue), implementation (red), and unallocated (green).

6.11.3.1 Initial Gantt Chart



Figure 6.11.1 Initial Gantt chart for deliverable one (top) and deliverable two (bottom)

6.11.3.2 Revised Gantt Chart

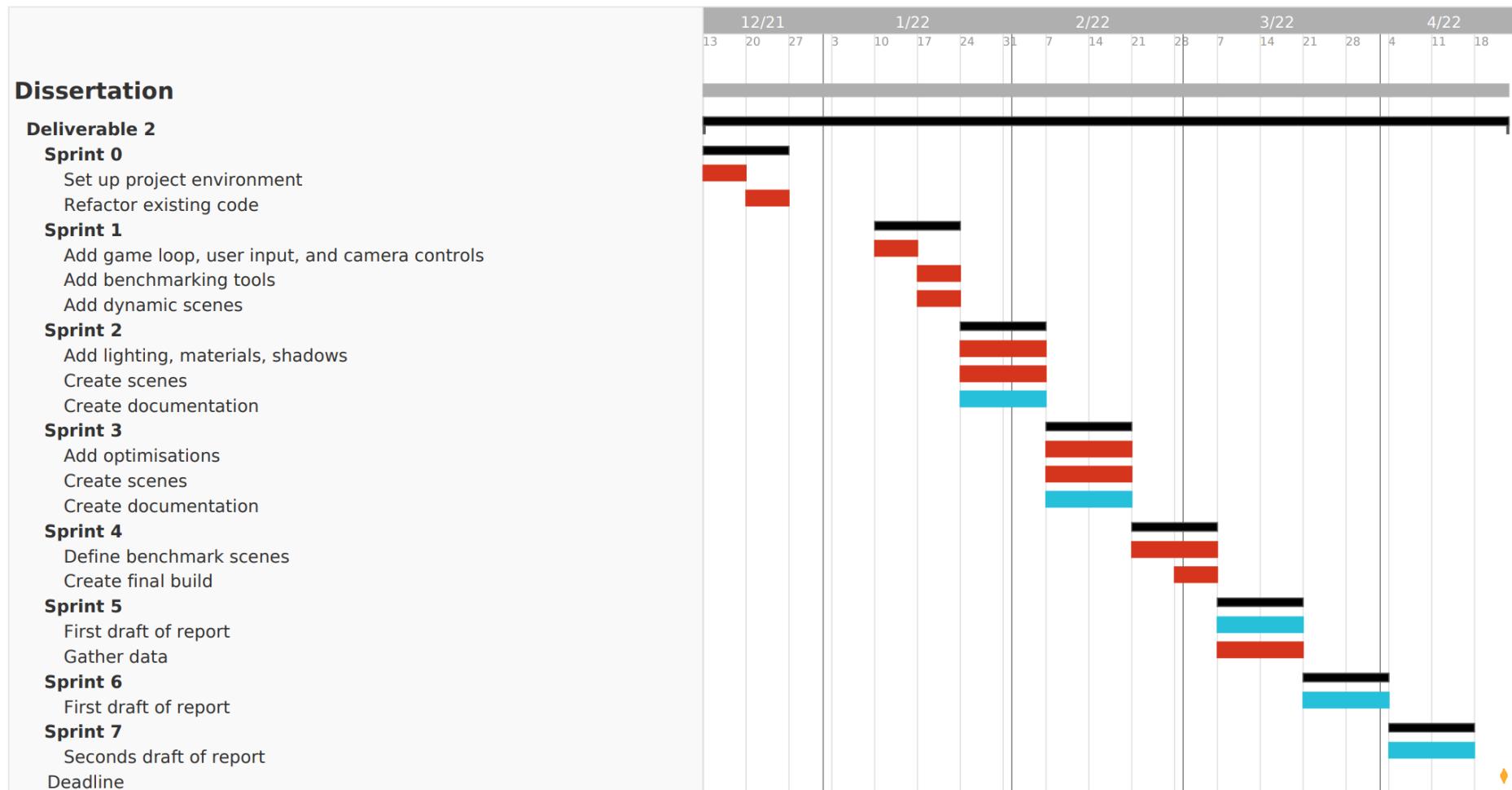


Figure 6.11.2 Revised Gantt chart for deliverable two

7 References

- [1] Jack Challoner, "How Mandelbrot's fractals changed the world - BBC News," 2010. <https://www.bbc.co.uk/news/magazine-11564766> (accessed Nov. 10, 2021).
- [2] University of Waterloo, "Top 5 applications of fractals." <https://uwaterloo.ca/math/news/top-5-applications-fractals> (accessed Nov. 10, 2021).
- [3] T. Kluge, "Fractals in nature and applications," 2000. <https://kluge.in-chemnitz.de/documents/fractal/node2.html> (accessed Nov. 10, 2021).
- [4] The Fractal Foundation, "Chapter 12 - FRACTAL APPLICATION." <http://fractalfoundation.org/OFC/OFC-12-2.html> (accessed Nov. 10, 2021).
- [5] S. Baarda, "3D Fractal Geometry Renderer," 2022. <https://github.com/SolomonBaarda/fractal-geometry-renderer> (accessed Sep. 24, 2021).
- [6] S. Baarda, "Realtime Fractal Renderer Documentation," 2022. <https://solomonbaarda.github.io/fractal-geometry-renderer/FractalGeometryRenderer/documentation/html/index.html> (accessed Mar. 16, 2022).
- [7] Wikipedia, "Wacław Sierpiński." https://en.wikipedia.org/wiki/Wac%C5%82aw_Sierpi%C5%84ski (accessed Nov. 11, 2021).
- [8] H. Segerman, "Fractals and how to make a Sierpinski Tetrahedron", Accessed: Nov. 11, 2021. [Online]. Available: <http://www.segerman.org>
- [9] L. Riddle, "Sierpinski Carpet - Agnes Scott College." <https://larryriddle.agnesscott.org/ifs/carpet/carpet.htm> (accessed Nov. 11, 2021).
- [10] J. Baez, "Menger Sponge | Visual Insight," 2014. <https://blogs.ams.org/visualinsight/2014/03/01/menger-sponge/> (accessed Nov. 11, 2021).
- [11] Wikipedia, "n-flake." <https://en.wikipedia.org/wiki/N-flake> (accessed Nov. 16, 2021).
- [12] A. Douady, "Julia Sets and the Mandelbrot Set," *The Beauty of Fractals*, pp. 161–174, 1986, doi: 10.1007/978-3-642-61717-1_13.
- [13] Wikipedia, "Mandelbrot set." https://en.wikipedia.org/wiki/Mandelbrot_set (accessed Nov. 13, 2021).
- [14] V. da Silva, T. Novello, H. Lopes, and L. Velho, "Real-time Rendering of Complex Fractals," in *Ray Tracing Gems II*, 2021, pp. 529–544. doi: https://doi.org/10.1007/978-1-4842-7185-8_33.
- [15] D. White, "The Unravelling of the Real 3D Mandelbrot Fractal," 2009. <https://www.skytopia.com/project/fractal/mandelbulb.html> (accessed Nov. 12, 2021).
- [16] R. Englund, S. Seipel, and A. Hast, "Rendering Methods for 3D Fractals, Bachelor Thesis," 2010.

- [17] D. White, "The Unravelling of the Real 3D Mandelbrot Fractal page 2," 2009.
<https://www.skytopia.com/project/fractal/2mandelbulb.html> (accessed Nov. 20, 2021).
- [18] D. Mankin, "True 3D mandelbrot type fractal, distance estimation optimisation," 2009.
<http://www.fractalforguums.com/3d-fractal-generation/true-3d-mandlebrot-type-fractal/msg8073/#msg8073> (accessed Nov. 20, 2021).
- [19] M. McGuire, E. Enderton, P. Shirley, and D. Luebke, "Real-time Stochastic Rasterization on Conventional GPU Architectures," *Proceedings of the conference on high performance graphics*, pp. 173–182, 2010, Accessed: Nov. 20, 2021. [Online]. Available: <http://research.nvidia.com>
- [20] J. Peddie, *Ray Tracing: A Tool for All*. 2019. doi: 10.1007/978-3-030-17490-3.
- [21] Mikael Hvidtfeldt Christensen, "Distance Estimated 3D Fractals," 2011.
<http://blog.hvidtfeldts.net/index.php/2011/08/distance-estimated-3d-fractals-ii-lighting-and-coloring/> (accessed Nov. 04, 2021).
- [22] I. Quilez, "Soft Shadows in Raymarched SDFs," 2010.
<https://iquilezles.org/www/articles/rmshadows/rmshadows.htm> (accessed Nov. 04, 2021).
- [23] I. Quilez, "Distance Functions," 2013.
<https://www.iquilezles.org/www/articles/distfunctions/distfunctions.htm> (accessed Oct. 28, 2021).
- [24] Khan Academy, "Matrices as transformations."
<https://www.khanacademy.org/math/precalculus/x9e81a4f98389efdf:matrices/x9e81a4f98389efdf:matrices-as-transformations/a/matrices-as-transformations> (accessed Apr. 10, 2022).
- [25] I. Quilez, "Distance to Fractals," 2004.
<https://www.iquilezles.org/www/articles/distancefractals/distancefractals.htm> (accessed Nov. 04, 2021).
- [26] Wikipedia, "Phong reflection model."
https://en.wikipedia.org/wiki/Phong_reflection_model (accessed Nov. 18, 2021).
- [27] Boston, "What Is GPU Computing?" <https://www.boston.co.uk/info/nvidia-kepler/what-is-gpu-computing.aspx> (accessed Nov. 16, 2021).
- [28] NVIDIA, "Programming Guide :: CUDA Toolkit Documentation."
<https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html> (accessed Nov. 17, 2021).
- [29] The Khronos Group Inc, "OpenGL Overview." <https://www.khronos.org/opengl/> (accessed Nov. 17, 2021).
- [30] The Khronos Group Inc, "OpenCL Overview." <https://www.khronos.org/opencl/> (accessed Nov. 17, 2021).
- [31] The Khronos Group Inc, "Conformant Products," 2022.
<https://www.khronos.org/conformance/adopters/conformant-products/> (accessed Apr. 09, 2022).

- [32] “Fragmentarium (original version, now unmaintained),” 2017.
<https://github.com/Syntopia/Fragmentarium> (accessed Nov. 17, 2021).
- [33] “Fragmentarium,” 2022. <https://github.com/3Dickulus/FragM> (accessed Nov. 17, 2021).
- [34] “Synthclipse,” 2019. <http://synthclipse.sourceforge.net/> (accessed Apr. 02, 2022).
- [35] “Mandelbulb3D,” 2020. <https://github.com/thargor6/mb3d#Coloring> (accessed Nov. 17, 2021).
- [36] “Mandelbulber,” 2022. <https://www.mandelbulber.com/> (accessed Apr. 02, 2022).
- [37] “Fractal Lab - Interactive WebGL Fractal Explorer,” 2011.
<https://hirnsohle.de/test/fractalLab/> (accessed Mar. 24, 2022).
- [38] Microsoft, “Visual Studio Development Features,” 2022.
<https://visualstudio.microsoft.com/vs/features/> (accessed Mar. 16, 2022).
- [39] Microsoft, “Welcome back to C++ - Modern C++,” 2022. <https://docs.microsoft.com/en-us/cpp/cpp/welcome-back-to-cpp-modern-cpp?view=msvc-170> (accessed Mar. 16, 2022).
- [40] Khronos®, “The C++ for OpenCL 1.0 Programming Language Documentation,” 2021.
https://www.khronos.org/opencl/assets/CXX_for_OpenCL.html#_the_c_for_opencl_programming_language (accessed Nov. 04, 2021).
- [41] Atlassian, “What is Agile?” <https://www.atlassian.com/agile> (accessed Nov. 13, 2021).
- [42] The Khronos Group Inc, “clBuildProgram.”
<https://www.khronos.org/registry/OpenCL/sdk/1.0/docs/man/xhtml/clBuildProgram.html> (accessed Mar. 09, 2022).
- [43] Nvidia, “NVIDIA Nsight Systems,” 2022. <https://developer.nvidia.com/nsight-systems> (accessed Apr. 06, 2022).
- [44] D. Vellante and D. Foyer, “A new era of innovation: Moore’s Law is not dead and AI is ready to explode - SiliconANGLE,” 2021. <https://siliconangle.com/2021/04/10/new-era-innovation-moores-law-not-dead-ai-ready-explode/> (accessed Mar. 30, 2022).
- [45] Heriot-Watt University, “Internal documentation site for the ECR Robotarium Cluster at Heriot-Watt University.” <https://ecr-cluster.github.io/> (accessed Mar. 30, 2022).
- [46] Google, “Google Colaboratory FAQ.” <https://research.google.com/colaboratory/faq.html> (accessed Mar. 30, 2022).
- [47] J. C. Hart, D. J. Sandin, and L. H. Kauffman, “Ray Tracing Deterministic 3-D Fractals,” *ACM SIGGRAPH Computer Graphics*, vol. 23, no. 3, pp. 289–296, 1989, doi:
<https://doi.org/10.1145/74334.74363>.
- [48] I. Quilez, “Normals for an SDF,” 2015.
<https://www.iquilezles.org/www/articles/normalsSDF/normalsSDF.htm> (accessed Mar. 31, 2022).
- [49] I. Quilez, “Mandelbulb,” 2009.
<https://www.iquilezles.org/www/articles/mandelbulb/mandelbulb.htm> (accessed Nov. 04, 2021).

- [50] Open Source Initiative, “The Open Source Definition,” 2007. <https://opensource.org/osd> (accessed Nov. 18, 2021).
- [51] Open Source Initiative, “GNU General Public License version 3,” 2007. <https://opensource.org/licenses/GPL-3.0> (accessed Nov. 18, 2021).
- [52] British Computer Society, “Code of Conduct for BCS Members,” 2021. Accessed: Nov. 18, 2021. [Online]. Available: <https://www.bcs.org/media/2211/bcs-code-of-conduct.pdf>
- [53] TeamGantt, “Online Gantt Chart Maker.” <https://www.teamgantt.com/> (accessed Nov. 18, 2021).