

1 Faster than NERFs – 3D Models from 2D Images

2 Aniket Rajnish 

3 Indian Institute of Technology, Gandhinagar, India

4 <http://makra.wtf>

5 aniket.r@iitgn.ac.in

6 Progyan Das 

7 Indian Institute of Technology, Gandhinagar, India

8 <http://progyan.me>

9 progyan.das@iitgn.ac.in

10 — Abstract —

11 Slow inference and training times have always been an issue with Neural Radiance Fields, and
12 voxel representations, often used in these papers, lead to prohibitively large memory requirements
13 and very long waiting times. While the accuracy of NERFs are very high, we may be willing to
14 sacrifice representation accuracy in the favour of time. We note that game-development, in particular,
15 suffers from a large latency from ideation to the actual creation of assets for deployment.

16 In addition, although faster implementations of NERFs, like Instant-NGP [4] and Plenoxels [7]
17 have been famously published in the past year, we wanted to devise an end-to-end tool for going
18 from an image to a quickly deployable 3D model with the least number of steps. Traditional ML
19 techniques use voxel, mesh, or point-cloud based rendering techniques – these are volumetric, and
20 often are bound by high time complexities (voxel-based rendering, for example, runs at $O(N^3)$
21 where N is the number of *voxels*, or 3-dimensional pixels, we are rendering). Instead, we use Signed
22 Distance Functions, a (somewhat) more complicated but overall less algorithmically complex solution
23 for rendering, that sacrifices some benefits of volumetric rendering for speed.

24 To make this possible, we wrote our own raymarching rendering engine on *C#* and *HLSL* (short
25 for High-Level Shader Language, a C-like language for use in Direct3D applications[3]), implemented
26 it in Unity, and combined it with a Convolutional Neural Network trained on a custom data-set
27 made on Blender. Our functioning assumption is that most complex shapes that we may want to
28 approximate may be built with a set of 35 primitive shapes and an accompanying boolean expression
29 combining a few operations (union, intersection, subtraction). This is the foundational principle of
30 Constructive Solid Geometry[8], and we have found that it works with great effect in our project.

31 Our end-product, SDFNet, is a tool built by game-developers, for game-developers, that takes a
32 simple 2D image, and generates a 3D, surface-rendered model that is immediately deployable for
33 development.

34

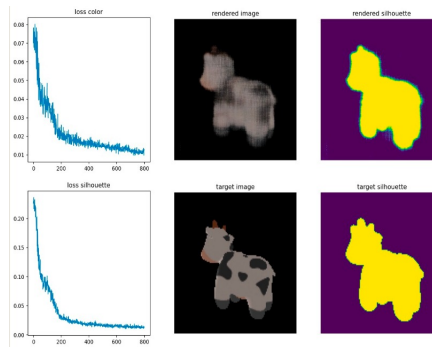
35 **1** Implementing the State of the Art

36 Before we started writing our own rendering engine, we wanted to see how Neural Radiance
37 Fields are implemented in *python*. To that end, we implemented a truncated version of the
38 classical NERF paper from Nvidia Research, based on *pytorch3D* documentation[1].

39 **1.1 Neural Radiance Fields (NERFs) through raymarching**

40 The input is a number of images of the target, from different angles and their corresponding
41 cameras, and our network attempts to build a scalar field that allows us to generate views
42 from any other angle. To fit the radiance field, we render it from the viewpoints of the target
43 cameras, and we compare the results with the observed target images and target silhouettes.





■ **Figure 1** Side-by-side: the drop in the huber loss, and our trained model.

44 Loss is calculated as the mean huber loss (related to smooth-L1 loss) between rendered
 45 colours and the sampled target images, predicted masks and sampled target silhouettes.

46 We use two losses –

- 47 1. The color loss – as it is traditionally used in NERFs, the color loss essentially compares
 48 a snapshot of the model from the same camera position and angle with the data-point,
 49 which in this case is our corresponding image.

```
50
51 ray_color = sampled_images(
52     target_images[batch_idx],
53     sampled_rays.xys
54 )
55 color_error = jnp.mean(
56     jnp.abs(
57         huber_loss(
58             rendered_images,
59             colors_at_rays,
60         )))
61
```

- 62 2. The silhouette loss – the silhouette loss forces the model to absorb the rays where necessary,
 63 and not pass through it, and vice-versa. This is possible because our dataset comes with
 64 segmentation-masks – otherwise, we could have used an architecture like Mask-RCNN for
 65 image segmentation and obtained a good approximate for the same.

```
66
67 silhouettes_at_rays =
68     sampled_images(
69         target_silhouettes[batch_idx, None],
70         sampled_rays.xys
71     )
72 silhouette_err = jnp.mean(
73     jnp.abs(
74         huber(
75             rendered_silhouettes,
76             silhouettes_at_rays,
77         )))
78
```

79 1.2 Drawbacks of NERFs

80 As we can see, current state-of-the-art techniques produce very accurate results. However,
 81 as the paper mentions, they are slow, and often prone to taking hours to train. While it is

82 possible to do away with neural networks and use classical machine learning for a speedup
 83 [7], we wanted to minimize the time-consuming process of training a radiance field altogether,
 84 and therefore, we decided to shift to neural networks only for *detecting* shapes and structures
 85 in our images, and not for *reconstructing* our models from the images.

86 **2 Signed Distance Functions and Surface Rendering**

87 Computationally, geometry is often stored explicitly as a list of points, triangles, or other
 88 geometric fragments [5]; however, these methods are computationally expensive, and we
 89 can devise both parametric and non-parametric methods for expressing these geometries
 90 implicitly. Signed Distance Functions, therefore, are a method for parametric implicit surface
 91 representation. [6]

92 These signed distance functions, or SDFs for short, are defined as continuous functions
 93 that, when passed the coordinates of a point in space, will return the shortest distance
 94 between that point and some arbitrary surface corresponding to that specific function. The
 95 sign of the return value indicates whether the point is inside that surface or outside (hence
 96 *signed* distance function).

For example, for a sphere centered at the origin, the standard SDF is mentioned below.
 [6]

$$f(p) = \vec{p} - \vec{r}$$

97 **2.1 Rendering Shapes**

98 We wrote an Image Effect shader to render objects directly in the screen space instead of
 99 creating instances of individual objects. We wrote a raymarching loop in the shader to render
 100 these shapes using their individual signed distance functions. All the parameters were taken
 101 from our model in a CSV file and were communicated to the shader from a C# script using
 102 Compute Buffers. The dimensional parameters were stored in a custom class of `Vector12`
 103 with 12 fields (maximum dimensional inputs that any shape can take) for floats, as the
 104 Shader language doesn't support dynamic arrays. So these parameters were communicated
 105 in the following way:

```
106 dimensions[0] = new vector12(cyl.r, cyl.h, 0,0,0,0,0,0,0,0,0,0);
107 dimensions[1] = new vector12(cap.r1, cap.r2, cap.h, 0, 0, 0, 0,
108                                0, 0, 0, 0,0);
```

111 Computer buffers were also used to communicate other information like the number of shapes
 112 to be rendered and the blend factor for the operations for each shape. All the shapes are
 113 rendered on render texture in front of the camera, the dimensions of which are communicated
 114 to the shader.

115 **2.2 Building a custom-editor in Unity**

116 A custom editor was developed in Unity to aid the need to fine-tune the objects rendered on
 117 the screen The editor could be accessed using the Unity's inspector. The following parameters
 118 were governed using the Custom Editor –

Shape Operation	Color	Blend Factor	Dimension Factor
-----------------	-------	--------------	------------------

:4 Faster than NERFs – 3D Models from 2D Images

120 Similarly, following spatial parameters were governed by the Unity’s inspector component
121 –

122	Shape position	Shape Orientation (Quaternion/Euler Angles)
-----	----------------	---

123 **3 Predicting Shapes – designing the Pipeline.**

124 We used a CNN architecture similar to *Alex Net* and built on top of that. The input to
125 the model was an image of the object, in our example case, the bottle with a constructive
126 geometry made of primary shapes. A total of 17 metadata entries, along with every image,
127 were stored to reconstruct the model. These entries correspond to the presence, shape, and
128 color of the sub-parts of the bottles. These were the output of the model.

129 **3.1 Architecture and Parameters of Neural Network.**

130 The architecture of the model consisted of CNNs, Pooling layers, activation functions such as
131 ReLU, and dense linear layers towards the end. Finally, we received the 17 labels as output
132 from the model, and these were read by the renderer. The output is received from the model
133 in a csv sheet which includes –

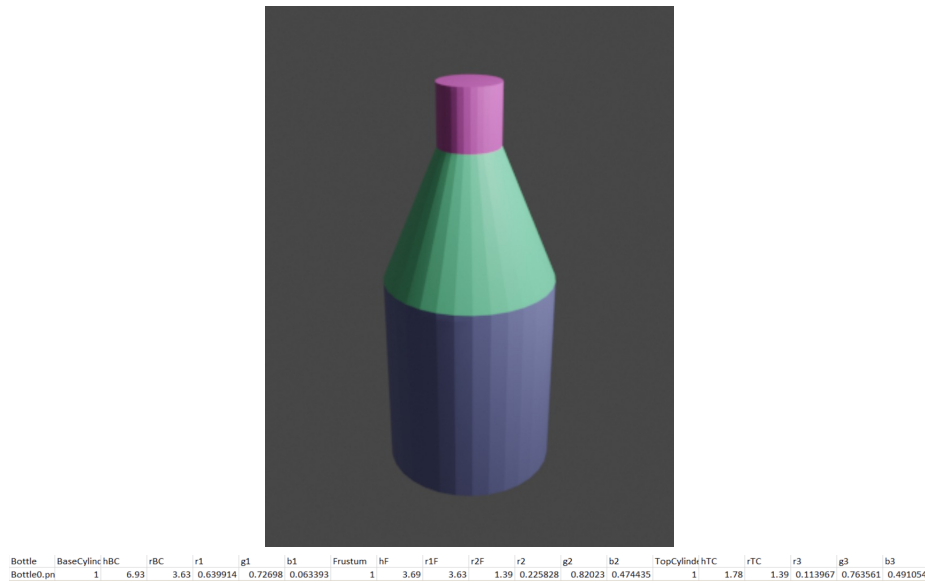
Parameter	Data-type	Description
Shape Index	int	Describes which out of a predefined list of primitives the given shape is.
Shape Position	Vector3	Describes the position of the shape in 3D space
Shape Rotation	Vector3	Describes the orientation of the shape in 3D space
134 RGB Values	Vector3	Denotes the RGB values of the color of the shape in the form of a 3-dimensional vector.
Shape Dimensions	Vector12	As a 12-dimensional vector, describes important absolute and relative dimensions for attributes like radius, height, et cetera, for the object. This is a sparse vector, and depending on the shape, many or none of the components may remain 0.

135 **3.2 Creating our dataset in Blender**

136 We used Blender to prepare the dataset of bottles of different shapes, sizes, and colours.
137 Random gaussians were used to generate the dimensions of the bottles. The dimensions
138 of the bottles along with the information of color were normalized before being fed to the
139 model in the subsequent steps. The images and the information regarding the dimensions of
140 the bottles were saved. Each bottle comprised two cylinders and one frustum. A total of
141 17 metadata entries, along with every image, were stored to reconstruct the model. These
142 entries correspond to the presence, shape, and color of the sub-parts of the bottles. In the
143 end, 400 distinct data points were generated and used.

144 **3.3 In more detail: Parameters in use.**

145 The *shape-index* is an int used to determine the shape we’re trying to render. For instance,
146 2 denotes a cylinder, and 5 denotes a Frustum. The *Shape Position* basically represents
147 the represents predicted center of mass of the segmented shapes. The coordinates of which
148 are scaled values of the coordinates of the pixel. Thus, these dimensions are not absolute



■ **Figure 2** A sample datapoint from the data-set, and the corresponding photograph.

149 but rather relative. The *Shape Rotation* values are the quaternion of the segmented shapes,
 150 again, it is calculated relative to the vertical axis. The *RGB* values are the average *RGB*
 151 values of each pixel on the segmented shapes. The *Shape Dimensions* are the individual
 152 dimensions required to define a particular shape. Take, for instance, a cylinder needs radius
 153 and height, a sphere just needs a radius, and a frustum/capped cone needs height, top, and
 154 bottom radius. Again, these are scaled values of the pixels covered.

155 3.4 Justification: Why Unity?

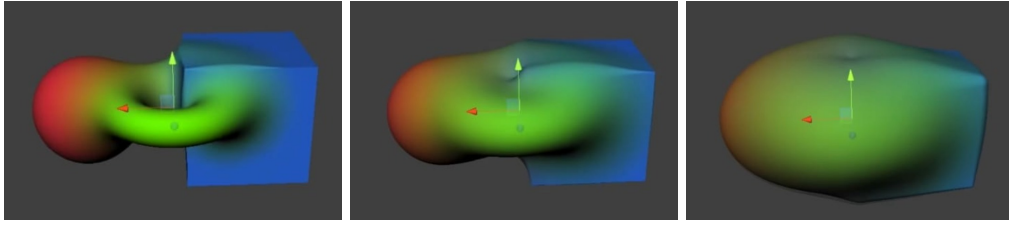
156 The Unity Engine was primarily used for the following reasons –

- 157 1. To aid in rendering the shadows using Unity’s Directional Light Object. It takes its
 158 quaternion in Euler form (`Vector3`) into account.
- 159 2. To map the 2D image onto the screen space using Camera Frustum (Matrix 4×4) and
 160 Camera to World Matrix (Matrix 4×4).

161 We wrote a custom Raymarcher class in *C#* to provide the following data manually –

Parameter	Data-type	Description
Shape Count	<code>int</code>	The length of the rows of the <code>csv</code> sheet, i.e, the number of shapes to be rendered on-screen. We later refactored this to be part of the output file.
Operation Index	<code>int</code>	An integer value used to denote which operation (union, intersection, subtraction) to perform with each shape.
Blend Factor	<code>float</code> , $0 \leq x \leq 1$	Whether or not to smoothen out the edges of different shapes, and finely blend them with each other.

163 Ideally, the model should have predicted these data-points as well, but we couldn’t train it
 164 to do so at the moment and would be working on it further. We have found architectures like
 165 *CSG-Net*, that infer boolean expressions for Constructive Solid Geometry from 3D Models,
 166 that have piqued our interest, and we look forward to using them in our work.



■ **Figure 3** Left to right: How the blend factor increases from 0.36, to 0.59, to 0.93

167 3.5 Compute Buffer Conversion – C# to HLSL

168 All this data is passed to an HLSL-based shader to be rendered. This communication between
 169 the shader and C is done using a Compute Buffer of stride (size) 96 bytes.

Parameter	C#	HLSL (Compute Buffer Conversion)
Shape Index	int	int
Shape Position	Vector3	float3
Shape Rotation	Vector3	float3
RGB Values	Vector3	float3
Light Direction	Vector3	float3
Camera Frustrum	Matrix 4 × 4	uniform float 4×4
Camera to World Matrix	Matrix 4 × 4	uniform float 4 × 4
Shape Count	int	int
Operation Index	int	int
Blend factor	float	float

171 Apart from this the Image Effect shader used additional parameters –

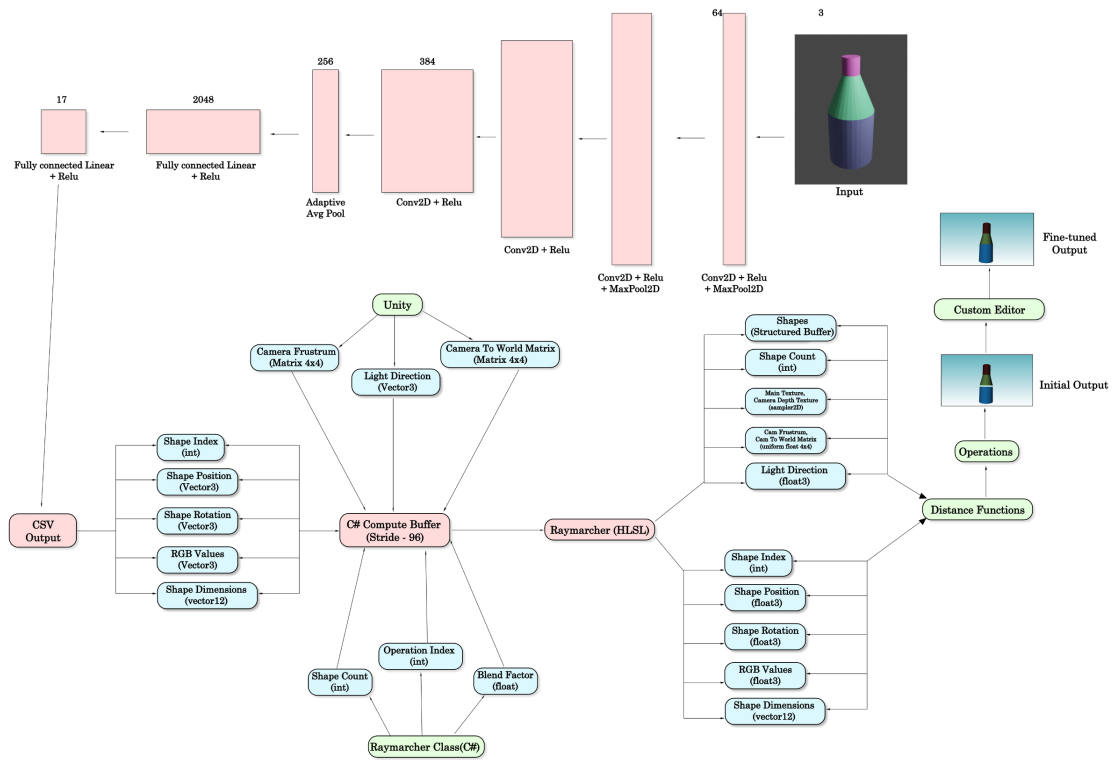
- 172 1. Main Texture (`sampler2D`)
- 173 2. Camera Depth Texture (`sampler2D`)
- 174 3. Shapes (`Structured Buffer`)

175 The Main Texture and Camera Depth Texture is used to render multiple objects in the
 176 screen space without needing to create an instance for each. These are Image Effect shaders,
 177 that work like a post-processing effect over the screen-space, as opposed to vertex-shaders
 178 that work in world-space, which makes them more efficient and light for rendering crowded
 179 scenes. Please note that the entire SDF renderer is written in HLSL.

180 4 Putting it all together with raymarching.

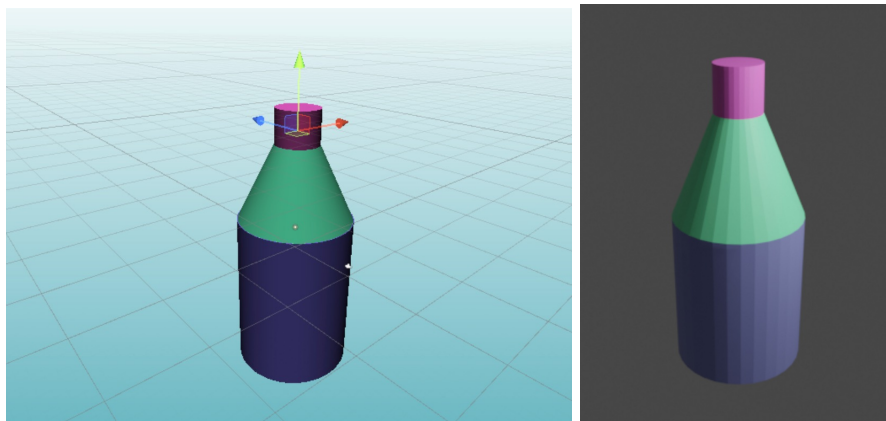
181 We use raymarching for rendering – here, all attributes of the scene are implicitly defined in
 182 terms of some signed distance function. To find the intersection between the view ray and
 183 the scene, we start at the camera, and move a point along the view ray. At each step, we
 184 check if the SDF evaluates to a negative number at that point. If it does, we consider this a
 185 collision and initialize a surface at the point. [6][2]

186 This data is then used by the raymarching loop to decide the distance functions and
 187 operations for each shape and the parameters that these functions would use to render every



■ **Figure 4** A flowchart of the entire pipeline for going from 2D image to 3D model

188 shape as perceived from the 2D image. The predicted model is surface-rendered and, as
 189 anticipated, comes with some flaws, which can later be fine-tuned using the custom editor
 190 that we wrote to reconstruct a fairly accurate model.



■ **Figure 5** Side-by-side: reconstructed, SDF-rendered model, and the input image

191 **5 Scope for improvement.**

192 The pertinent areas of improvement for our tool have been listed below. These are areas of
 193 rapid development, and we foresee them being implemented very soon.

194 **1. Expanding for radially asymmetric geometries.**

195 Since our tool only takes into account *one* image, we understand that for asymmetric
196 geometries, there shall exist attributes that are occluded in any one camera angle. In
197 addition, even for radially symmetric geometries, our model requires an angle that properly
198 exposes all primitives present in the shape.

199 **2. Expanding for complex foregrounds and crowded backgrounds.**

200 The model often fails for images with complex foreground geometries or crowded back-
201 grounds. We wish to forego that limitation with a combination of foreground-background
202 image segmentation, and some flavour of the *CSG-Net* model, to break images down to
203 corresponding boolean expressions for constructive solid geometry.

204 **3. Raising number of primitives.**

205 While we believe the 30 primitives that have been integrated into the Raymarching engine
206 should be enough to express most geometries out there, there might be some esoteric,
207 complex shapes that our model may not be able to approximate. We hope to soon get
208 from 30 primitives to a planned 47 primitives.

209 **4. Predict Blend factor and Operations.**

210 At the moment, our model does *not* predict the boolean operations and the blend factor.
211 We wish that both can be implemented soon.

212 **5. Increase dataset variety.**

213 Our model was only trained on a dataset consisting of frustrums and cylinders, out of
214 the 30 primitives available. We wish to increase that variety with a much wider number
215 of primitives.

216 **6 Conclusion and acknowledgements**

217 We have been able to produce a functioning prototype that can take a simple 2D image of a
218 radially symmetric geometry and reconstruct a 3D representation through signed distance
219 functions, with a combination of constructive solid geometry and neural networks. There is
220 huge scope for improvement, and we are excited to keep working on this project, and also
221 branch out into other domains.

222 We are thankful to Prof. Shanmuganathan Raman, for his guidance across the duration
223 of the project. We are also grateful to his student, Ashish Tiwari, for his timely help in
224 understanding hard concepts whenever we required it, and we are indebted to Shruhid
225 Bantia, a third-year undergraduate student at IIT Gandhinagar, who helped us build a
226 significant portion of the tool.

227 **References**

- 228 **1** Pytorch3d · a library for deep learning with 3d data. URL: [https://pytorch3d.org/
229 tutorials/fit_simple_neural_radiance_field](https://pytorch3d.org/tutorials/fit_simple_neural_radiance_field).
- 230 **2** Ray marching and signed distance functions. URL: [https://jamie-wong.com/2016/07/15/
231 ray-marching-signed-distance-functions/](https://jamie-wong.com/2016/07/15/ray-marching-signed-distance-functions/).

- 232 **3** Yong He, Kayvon Fatahalian, and Tim Foley. Slang: Language mechanisms for extensible real-
233 time shading systems. *ACM Trans. Graph.*, 37(4), jul 2018. doi:10.1145/3197517.3201380.
- 234 **4** Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics
235 primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15,
236 July 2022. doi:10.1145/3528223.3530127.
- 237 **5** Jeong Joon Park, Peter Florence, Julian Straub, Richard A. Newcombe, and Steven Lovegrove.
238 DeepSDF: Learning continuous signed distance functions for shape representation. *CoRR*,
239 abs/1901.05103, 2019. URL: <http://arxiv.org/abs/1901.05103>, arXiv:1901.05103.
- 240 **6** Inigo Quilez. Signed distance functions. URL: [https://iquilezles.org/articles/
241 distfunctions/](https://iquilezles.org/articles/distfunctions/).
- 242 **7** Sara Fridovich-Keil and Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and
243 Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022.
- 244 **8** Gopal Sharma, Rishabh Goyal, Difan Liu, Evangelos Kalogerakis, and Subhransu Maji. Csgnet:
245 Neural shape parser for constructive solid geometry. *CoRR*, abs/1712.08290, 2017. URL:
246 <http://arxiv.org/abs/1712.08290>, arXiv:1712.08290.