

# Unraveling Imagination Synthesis

## constructing Conceptual Artistry into a 3D scene through Neural Radiance Fields

(David) Zichen Liu<sup>a</sup>

<sup>a</sup>*Shenzhen college of international education,*

---

### Abstract

In this research, I propose an architecture that can reconstruct images into implicit 3D representation through Neural radiance field (NeRF). With a consistent combination between the world graphics and computer vision, I dive deep into the detailed implementation of such model and discover the constraint in time usage as well as computational power, I strongly believe that in the future, this area will become the pioneering research in meta verse, virtual reality and augmented reality.

**Keywords:** Graphics, 3D reconstruction, computer vision, NeRF

---

### 1. Introduction

Realistic or Imaginary 3D scene construction have always been a difficult problem to achieve, however, highly demanded in the fields of gaming, virtual reality, and robotic receptions. Imagine that we can obtain the novel 3D structure just with the 2D images whether it's photo from reality or conceptual art made by artist? In order to give the computer the power of understanding 3D constitution from 2D images, we can empower it with the advanced methods – Neural Radiance Field (NeRF) (Mildenhall et al. 2021), (Gao et al. 2022) which proposing a state-of-the-art result in synthesizing novel views based on a single or Multi-view image. Extended the new possibility for machines to understand from a visual perspective. To summarize what I have done for the scene construction, the first stage is to train a generative model that's based on NeRF representation of light fields, with a rendering process mimic by an Multi Layer Perceptron (MLP), which provide us with differentiable feed-forward process that can be back-propagate with loss from image.

As NeRF use a process of reverse the common rendering process of geometry to image pipeline, which can be classified as a way of differentiable rendering (Kato et al. 2020). In a classical rendering process Rasterization (Pineda 1988) and Ray-tracing (Appel 1968) are commonly used algorithm that samples image first from meshes in a certain camera view to obtain the transformation coordinate with its experience called the vertex shader, and then go through a rasterizer and fragment

shader to obtain the final image. For the purpose of reverse the whole process, I have to make sure every step within the rendering process is differentiable therefore I am able to apply chain rule to backpropagate images back to 3D representation. In the rendering process, all the coordinate transformation can be defined by matrix multiplication and the texturing process are defined by the user thus they are differentiable, while in the rasterization process the color need to be considered as a constant, although I might know that if there is a great shift between the color value, this might indicate an edge of the 3D representation I intended to reconstruct, therefore the whole process is not reversible. One way of approximating the gradient is to consider the image as a field of color (Loper and Black 2014) otherwise there is also some methods that rebuild the whole rendering process with each stage differentiable (Liu et al. 2019) for explicit 3D representation like mesh. NeRF on the other hand, stimulate this process by using an MLP (although it's been proven that neural network is not necessary (Fridovich-Keil et al. 2022)) to mimic the implicit representation function with an input of 5D coordinate and output a voxel scene (the detail explanation of how it functions please refer 3.1).

For an imaginary scene created by artist, it's impossible to obtain the camera position which are crucial to NeRF, compare to real picture capture by the camera. Therefore, I need to estimate the camera status from the drawing beforehand. To accomplish this, visual optometry that estimate the transformation between cameras

based on the images that I captured.

In this research, I produce a pipeline of generating both realistic and imagery content that are represented in 3D by using NeRF with camera pose estimation. Producing a next generation workflow that might empowers many areas include VR/AR, conceptual art, or even advertisement.

## 2. Literature review

### 2.1. Text Guided/shape awareness Object Generation

It's essential to be aware that to generate a 3D representation from art work, there is no shapes and multi-view data to optimize the Neural radiance field. Text alignment to 3D shapes in NeRF are popular researching field (Jain et al. 2022)(Poole et al. 2022) (Lin et al. 2023) as it provides more flexibility and modifiability for human user to interact with while nature language input is considering been easier to use compare to drawing. DreamField (Jain et al. 2022) provide a simple and constrained 3D representation with neural guidance that supports diverse generation from caption zero-shot, the text-alignment process is achieve by pretrain representation using CLIP (Radford et al. 2021) which consist of transformer image encoder with a masked transformer text encoder. Further more Dreamfusion (Poole et al. 2022) empower its process with the methods of stable diffusion model (Rombach et al. 2022). While providing a shape awareness methods with another approach with stochastic process over Variational Auto Encoder.

### 2.2. Scene-scale Contraction

Enable to model the 2D representation into a large-scale 3D scene representation instead of object centric reconstruction in order to understand the special relationship of object in space. Block-Nerf (Tancik et al. 2022), decompose the scene into individual trained NeRFs to represent a large-scale environments. Still facing the difficulties in presenting transient object like cars, lack of connection between each block, and constrain by the computational power. While Neuralangelo (Li et al. 2023), achieve the fine details for a large scene using a coarse-to-fine optimization by fine to the hyperparameter within training the MLP for NeRF, preventing over-fit to local minimum. Since egocentric scene is highly required in the robotics scenario, map reconstruction is also a significant process in understanding the environment. Nice-SLAM (ibid.) a simultaneous localization and mapping (SLAM) (Fuentes-Pacheco, Ruiz-Ascencio, and Rendón-Mancha 2015) algorithm, which introduce a hierarchical reconstruction methods

involve a pretrained differentiable renderer from coarse to fine level feature grid with Mean Square Error (MSE) loss within depth and image.

### 2.3. Reconstruction without pose prior

To train a differentiable rendering frame work, a camera poses is usually a curtail component that contribute to matrix transformation and bridging image with scene. However, when dealing with conceptual artwork, I barely have the ability to obtain the camera poses for a drawing. Therefore, building the scene from a prior without a pose is the main challenge here. The most common way to optimize pose is by self-supervised learning, where researcher can append posture optimization into NeRF's volume rendering process. Where iNeRF (Yen-Chen et al. 2021) optimize pose within the gradient flow, updating the hyperparameters during backpropagation. Although the camera pose calculation in NeRF is non-linear, as it involve affine transformation matrix multiplication, iNeRF addresses this by replacing the matrix with Lie algebra, changing multiplication to variable addition which allows gradient decent applied. For complex trajectory estimation, Nope-NeRF (Bian et al. 2023) demonstrate using mono-depth maps to optimize NeRF, regularize geometry from sequence of images with massive amount of motion between cameras.

## 3. 3 Methodology

### 3.1. NeRF

#### 3.1.1. Scene in Neural radiance field

NeRF representing scene implicitly by a mapping  $F_\theta(x, d) \rightarrow (c, \sigma)$  which input special location  $x \in \mathbb{R}^3$  with special Cartesian unit vector  $d$  replacing viewing direction  $(\theta, \phi)$  in application and output RGB color  $c$  with it's volume density  $\sigma$ .

#### 3.1.2. Volume rendering

Volume rendering are a ray tracing based algorithm that core in NeRF, rendering image both used in the calculating the optimisation loss and synthesis final result. (Kajiya and Von Herzen 1984) The intensity (color intensity) of light along a ray can be given by a function  $L(r(t), d)$  for  $r(t) = o + t * d$  set of particle on the ray Which change in  $L$  after the light transmit through a unit volume along direction  $d$  is given by *emission + inscattering-absorption* In NeRF the scene is been modeled as a cloud of space with particle emitting light whose intensity change though out the motion of photon. In practical I only use the absorption

plus emission to done the volume rendering process to reduce massive calculation done on scattering light between particles. (However, this leads to a constrain that my scene can only stay on the original lighting condition, without the ability to modify as well as combine multiple scene. (Zheng, Singh, and Seidel 2021))

For a space given front position  $t_n$  to back position  $t_f$ , lets assume  $s$  to be the distance from the front not overpassing  $t_f$  along the ray. With  $\sigma(s)$  as the volume density, defined as probability density of photon been absorb by the space, with reciprocal unit of length. Therefore it's easy to derive that for a unit length of  $\delta s$  at a given position of  $s$ , the probability of the light been absorb is  $\delta s * \sigma(s)$ . So the intensity of the light decrease along a minimal change with function of

$$L(s + \delta s) = L(s) * (1 - \delta s * \sigma(s)) \quad (1)$$

Reorganize the formula to rate of change in intensity gives

$$\begin{aligned} \frac{L(s + \delta s) - L(s)}{\delta s} &= -L(s) * \sigma(s) \\ L(s)' &= -L(s) * \sigma(s) \end{aligned}$$

Solving this Ordinary Differential Equation ODE gives an expression of  $L(s)$

$$L(s) = e^{-\int_{t_n=0}^s \sigma(x) \delta x} \quad (2)$$

Defining the transparency along a give distance  $s$  to  $s'$ , as  $T = \frac{L_o}{L_i}$  which  $L_o$  indicate the input intensity and  $L_o$  as output. Having this, the RGB value gained after rendering can be calculate by adding up color contribute by each particle along the ray.

$$C(r) = \int_{t_n}^{t_f} T(s) * \sigma(r(s)) * c(r(s), d) \delta s \quad (3)$$

To calculate this in practice, the sample position  $s$  should be discrete to

$$s_i = U[t_n + \frac{i-1}{N} * (t_f - t_n), t_n + \frac{i}{N} * (t_f - t_n)] \quad (4)$$

thus estimate the color with quadrature rule,  $\delta_i = s_i + 1 - s_i$ , with  $\alpha_i = 1 - e^{\sigma_i * \delta_i}$  be the density of color in the given length (Max 1995)

$$T_i = e^{-\sum_{j=1}^{i-1} \sigma_j * \delta_j} \quad (5)$$

$$C'(r) = \sum_{i=1}^N T_i * \alpha_i * c_i \quad (6)$$

If you interest in detail explanation, please have a look at further conclude derivation. (Tagliasacchi and Mildenhall 2022)

### 3.2. Optimize MLP

With the design of the general structure, however, there is still problems of not been able for the MLP to fit a fine-grained detail as well as the computational inefficiency. NeRF introduce two methods following to address the two problems.

#### 3.2.1. Positional encoding

Recent studies had shown that Polynomial Neural Networks (PNN) has a specific bias over the low frequency functions, which shows a slower learning rate in high frequency domain. The 5-dimensional input for NeRF is not able to train the detail information efficiently from the image. However, mapping the input into high dimensional space using high frequency function which allows it to fit the fine-grained information within the data. Intuitively NeRF modify the original MLP  $F_\theta$  with a encoding function  $\gamma$  represent  $R$  to  $R^{2L}$

$$\gamma(p) = (\sin(2^{0*\pi*p}), \dots, \cos(2^{(L-1)*\pi*p})) \quad (7)$$

$$F_\theta = F'_\theta \circ \gamma \quad (8)$$

Which gives the input Fourier feature with location normalize to  $[-1, 1]$ , thus in an iteration with same loss propagate back to hyper-parameters, the high frequency part contribute a much more massive weight in the final outcome, which allows learning the much finer details.

#### 3.2.2. Hierarchical volume sampling

Since it's time-consuming to direct sample along the ray using N step. It's better to strategies the sampling process from coarse to fine. Which propose a hierarchical way sampling along the ray. In coarse sampling stage, I nomalize the alpha value to obtain a piecewise-constant Probability Density Function (PDF) of which part along the ray contribute more on the finial rendered image.

$$\alpha'_i = \frac{\alpha_i}{\sum_{j=1}^{N_c} \alpha_j} \quad (9)$$

applying inverse transform sampling, by sample's a piecewise-constant set along the range to get a new set

of sample point that obey the distribution of the PDF, which than added with the first set of sample point to perform a second round of rendering.

### 3.3. Camera Calibration

To work out the rotation and translation of a camera from the 2D images is usually defined as a pose-estimation problem which is fairly important and well-studied in the world of graphics, computer vision and robotics. To simplified the calculation and considering the most common cases, my implementation only stick on the pinhole camera model for simplification, if you are interested in different camera model, please refer to the following passage (Sturm et al. 2011).

#### 3.3.1. Camera and Image

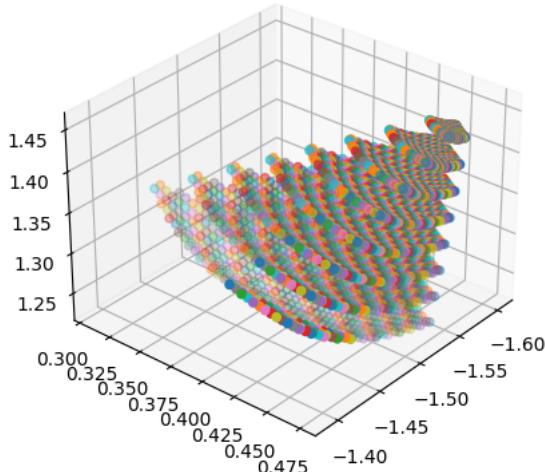


Figure 1: Light ray sample from a pinhole camera model

Pinhole camera model is the most simple and straight forward camera model used for a mono camera which is developed based on the physics phenomenon of candle's image appear on a plane at the other side of a pin-hole. The model consider light as a straight ray, allowing us to calculate the translation relation between the pixel  $(u, v)$  from image coordinate of the image, from the camera intrinsic matrix  $K$  which define based on its internal parameter include focal lens  $f$ , and the extrinsic matrix  $E \in \mathbb{R}^{4 \times 4}$  which describe the rotation  $R$  and translation  $t$  of the camera in the world coordinate. Given by following equation:

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \frac{1}{z} \begin{bmatrix} f_x & 0 & c_x & 0 \\ 0 & f_y & c_y & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} R & t \\ 0^T & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} \quad (10)$$

#### 3.3.2. Visual Odometry (VO)

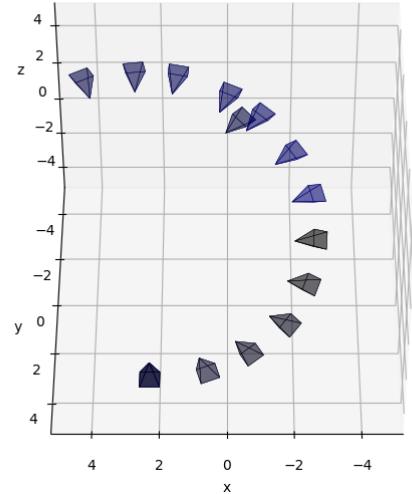


Figure 2: Visualisation of trajectory of multiple camera from circular path

Visual Odometry general term referring to the methods used in motion estimation based on the image taken by the sensor though out the time domain. In order to achieve a concordant result in poses estimation, I must produce a solid way in evaluating it. When a organism is locating it's direction in a map, it usually considering the points that features the location and recognizable after movement. Similarly, VO done this by comparing the chosen feature point, consist of Key-point and Descriptor (a vector designed to describe the information around the chosen pixel). By cleverly choosing and pairing those point in continues frame, it's then capable to find the motion of the camera. A common implemented algorithm is Oriented FAST and Rotated BRIEF (Aglave and Kolkure 2015)

To minimize the lose design for the pairing point, 2D-2D: Epipolar Geometry (Larsson, Pollefeys, and Os-karsson 2021) is used to give out a linear equation of translation follow by epipolar constraint, essential Matrix  $E' = t^\wedge * R$  decomposed by to it's singular value (SVD) (Kalman 1996) which leave a terms containing  $E'$  and produce least squared methods or Random Sample Consensus (RANSAC) (Bahraini, Bozorg, and Rad

2018) to optimize it. With all this been done, I am able to apply the classical methods of Triangulation (Hartley and Sturm 1997) to calibrate my camera with depth from the images.

## 4. Experiment

### 4.1. Scope estimation and difficulty indication

The central target of this research is the development and utilization of the Neural radiance field, a state-of-the-art technique in computer vision and world of machine learning used in novel scene synthesis and 3D object generation. Thus, I need to investigating methods to enhance the performance, robustness and scalability of NeRF. One primarily difficulty lies in the optimizing the NeRF algorithm, in other word a clever design of the training pipeline and the generation process. Base on my experiment on other's methods in generating high resolution object (Lin et al. 2023) (referring to 4.2) can be time-consuming as it involves numerous iterations to reach desired outcomes. training a NeRF model in a 3D scene can take up super long time (10 days in one RTX-3080) and gives a poor result. So, I should carefully choose the back bone I am using and aim in a minimal scale. Collecting and curating 2D conceptual art for conversion into 3D scenes can be a labor-intensive process. It may involve digitizing physical artwork, finding high-quality digital images, or creating artwork from scratch, which can consume a significant amount of time and effort. While there might take time to apply a robustness backbone and transfer it onto the scene of Art-2-3D task

### 4.2. Implementation

```
/* file structure */
|- config file #setup parameter
  |-- get config.py
  |-- config.txt
|- data
  |-- input #data used for training and testing
    |-- xxx.jpg
    |-- xxx.ply
  |-- output #output images
    |-- xxx.jpg
|- src #the helper function which defined as the back end of the program
  |-- Utils # helper libraries like colmap
  |-- visualization # visualization of camera
  |-- Nerf.py
    |-- rendering.py #function used in volume rendering
|- train.py #running file to run the file
```

Figure 3: General file structure of the program

my frame work was mostly developed using python programming language and pytorch frame work in building neural network and implementing tensor operation based on linear algebra. For the purpose of convenient

dependencies management and storage. This early-stage experiment are all test on Google-Colab and Linux Ubuntu 22.04 release, which provide CUDA environment. The training stage involve the use of T4 GPU, V100 GPU Nvidia provided by colab and RTX3080 Nvidia on Ubuntu. The core NeRF structure is based on the original design with 11 hidden layers with dimension 128 and reduce by halve for the last two. Activation function (Sagar Sharma, Simone Sharma, and Athaiya 2017) of ReLU is used between each layer, and Sigmoid for the last output RGB value. An extra dim is given at 9th layer to give out the density (sigma) before head which allow the front layer focus on finding the geometrical distribution and last 2 on colors and texture. Besides the uses of the human synthetic data "Tractorn" provided by the original NeRF which data is well distributed around the hemisphere of the parametrized 3D model. I further trained the model using real-world images collect by easy access mobile device. The images of individual object are taken using a circular trajectory around it



Figure 4: User input image

which helps to eliminate the scale ambiguity in the mono camera by providing translation, In addition capture more information from the object. As long as I gather the images, Colmap, a general implementation to solve the problem of Structure From Motion (SFM) (Özyeşil et al. 2017) as well as the Multi-View Stereo (MVS) (Furukawa, Hernández, et al. 2015), is used as a pre-built wheel to estimate the camera poses and in-

ternal parameter. The data are sampled as a field of light rays radiate out from individual cameras by a transformation matrix from pixel coordinate to  $\mathbb{R}^3$  structure space. Each ray are presented as a collaboration of ray origins ( $\mathbb{R}^3$ ), ray directions ( $\mathbb{R}^3$ , vector form), and the pixel value as ground truth ( $\mathbb{R}^3$ ). In addition, to maintain the consistency of gradient in back propagation stage, the input value should be normalized.

## 5. Discussion

### 5.1. Result and Evaluation

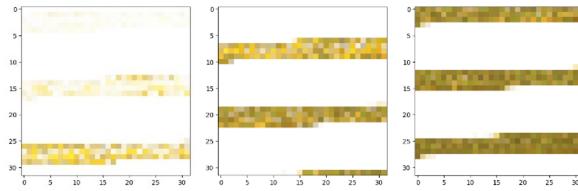


Figure 5: Failure case with lack in training data

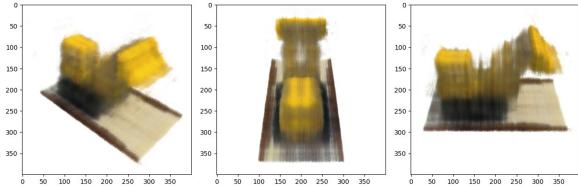


Figure 6: Course sampling example

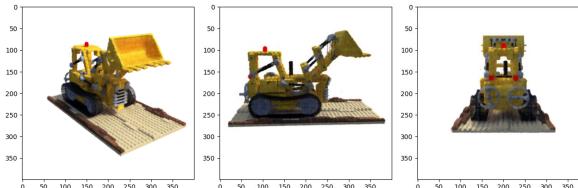


Figure 7: Hierarchical sampling example

The first experience are tested on human synthetic dataset used in the original implication of NeRF model. With the upper three picture showing the result with a failure case with a minimum number of light rays, while the middle one shows the course rendering stage without empowering by the hierarchical sampling strategy. The bottom third shows the rendering result with one more epoch of hierarchical sampling which shows a detailed generated novel scene.

Next I test my model on realistic photo taken though normal live mobile device, the first row shows failure

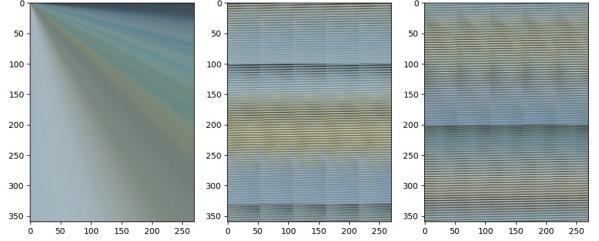


Figure 8: Failure case with wrong coordinate and lack in training data

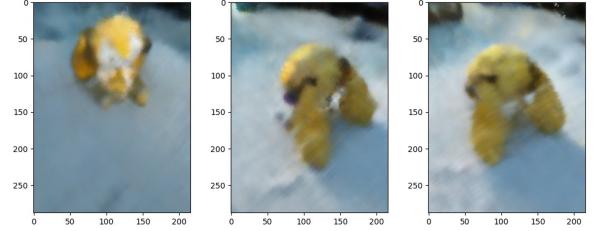


Figure 9: Finial results

cases with wrong calculation pixel to world sampling and not sufficient amount of data. While the last row demonstrate the capability of NeRF in rendering novel view as well as edit color though it's implicit representation.

### 5.2. Limitation

The sample algorithm calculate 3 for-loop iteratively. The number of process is calculated by num\_Of\_Camera \* H \* W which produce low efficiency in data sampling when the inputs is getting larger (current cases:  $15 * 1420 * 1080 > 1e07$ ). To imply the methods onto real world cases, and process in real time, there is clearly a needs for further algorithm design, more Computation power or even a customized hardware, and microprocessor.

### 5.3. Conclusion

In this work, I build an architecture that allow it's user to upgrade a picture taken from the real world into synthetic 3d spaces without geometrical prior given. despite that it's now stays in a prototype state without a user friendly interface and highly constrain in computational power and time usage. I believe that in the future, with time consistence diffusion workflow been implemented in, we will be able to construct 3d structure from simply your imagination, which shows a great potential in the fields of VR, AR and meta verse. Application

involves producing interactive collections in cloud museum; Creating digital class room, favor the student in young age or learning design and architecture; Providing a new methods and tool set for artist to explore the world.

## References

- Aglave, Prashant and Vijaykumar S Kolkure (2015). “Implementation Of High Performance Feature Extraction Method Using Oriented Fast And Rotated Brief Algorithm”. In: *Int. J. Res. Eng. Technol* 4, pp. 394–397.
- Appel, Arthur (1968). “Some techniques for shading machine renderings of solids”. In: *Proceedings of the April 30–May 2, 1968, spring joint computer conference*, pp. 37–45.
- Bahraini, Masoud S, Mohammad Bozorg, and Ahmad B Rad (2018). “SLAM in dynamic environments via ML-RANSAC”. In: *Mechatronics* 49, pp. 105–118.
- Bian, Wenjing et al. (2023). “Nope-nerf: Optimising neural radiance field with no pose prior”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4160–4169.
- Fridovich-Keil, Sara et al. (2022). “Plenoxels: Radiance fields without neural networks”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5501–5510.
- Fuentes-Pacheco, Jorge, José Ruiz-Ascencio, and Juan Manuel Rendón-Mancha (2015). “Visual simultaneous localization and mapping: a survey”. In: *Artificial intelligence review* 43, pp. 55–81.
- Furukawa, Yasutaka, Carlos Hernández, et al. (2015). “Multi-view stereo: A tutorial”. In: *Foundations and Trends® in Computer Graphics and Vision* 9.1–2, pp. 1–148.
- Gao, Kyle et al. (2022). “Nerf: Neural radiance field in 3d vision, a comprehensive review”. In: *arXiv preprint arXiv:2210.00379*.
- Hartley, Richard I and Peter Sturm (1997). “Triangulation”. In: *Computer vision and image understanding* 68.2, pp. 146–157.
- Jain, Ajay et al. (2022). “Zero-shot text-guided object generation with dream fields”. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 867–876.
- Kajiya, James T and Brian P Von Herzen (1984). “Ray tracing volume densities”. In: *ACM SIGGRAPH computer graphics* 18.3, pp. 165–174.
- Kalman, Dan (1996). “A singularly valuable decomposition: the SVD of a matrix”. In: *The college mathematics journal* 27.1, pp. 2–23.
- Kato, Hiroharu et al. (2020). “Differentiable rendering: A survey”. In: *arXiv preprint arXiv:2006.12057*.
- Larsson, Viktor, Marc Pollefeys, and Magnus Oskarsson (2021). “Orthographic-perspective epipolar geometry”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 5570–5578.
- Li, Zhaoshuo et al. (2023). “Neuralangelo: High-fidelity neural surface reconstruction”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8456–8465.
- Lin, Chen-Hsuan et al. (2023). “Magic3d: High-resolution text-to-3d content creation”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 300–309.
- Liu, Shichen et al. (2019). “Soft rasterizer: A differentiable renderer for image-based 3d reasoning”. In: *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 7708–7717.
- Loper, Matthew M and Michael J Black (2014). “OpenDR: An approximate differentiable renderer”. In: *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part VII 13*. Springer, pp. 154–169.
- Max, Nelson (1995). “Optical models for direct volume rendering”. In: *IEEE Transactions on Visualization and Computer Graphics* 1.2, pp. 99–108.
- Mildenhall, Ben et al. (2021). “Nerf: Representing scenes as neural radiance fields for view synthesis”. In: *Communications of the ACM* 65.1, pp. 99–106.
- Özyeşil, Onur et al. (2017). “A survey of structure from motion\*”. In: *Acta Numerica* 26, pp. 305–364.
- Pineda, Juan (1988). “A parallel algorithm for polygon rasterization”. In: *Proceedings of the 15th annual conference on Computer graphics and interactive techniques*, pp. 17–20.
- Poole, Ben et al. (2022). “Dreamfusion: Text-to-3d using 2d diffusion”. In: *arXiv preprint arXiv:2209.14988*.
- Radford, Alec et al. (2021). “Learning transferable visual models from natural language supervision”. In: *International conference on machine learning*. PMLR, pp. 8748–8763.
- Rombach, Robin et al. (2022). “High-resolution image synthesis with latent diffusion models”. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695.
- Sharma, Sagar, Simone Sharma, and Anidhya Athaiya (2017). “Activation functions in neural networks”. In: *Towards Data Sci* 6.12, pp. 310–316.
- Sturm, Peter et al. (2011). “Camera models and fundamental concepts used in geometric computer vision”. In: *Foundations and Trends® in Computer Graphics and Vision* 6.1–2, pp. 1–183.
- Tagliasacchi, Andrea and Ben Mildenhall (2022). “Volume rendering digest (for nerf)”. In: *arXiv preprint arXiv:2209.02417*.
- Tancik, Matthew et al. (2022). “Block-nerf: Scalable large scene neural view synthesis”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8248–8258.
- Yen-Chen, Lin et al. (2021). “inerf: Inverting neural radiance fields for pose estimation”. In: *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, pp. 1323–1330.
- Zheng, Quan, Gurprit Singh, and Hans-Peter Seidel (2021). “Neural relightable participating media rendering”. In: *Advances in Neural Information Processing Systems* 34, pp. 15203–15215.

## 6. Appendix

### Code

The main training loop of the program (train.py):

```
1 import torch
2 import numpy as np
3 from tqdm import tqdm
4 import matplotlib.pyplot as plt
5 from torch.utils.data import DataLoader
6
7 from configs.get_configs import get_args
8 #from datasets import dataset_dict
9
10 from srcs.Nerf import NerfModel
11 from srcs.render import render_rays, compute_accumulated_transmittance
12
13 def train(nerf_model, optimizer, scheduler, data_loader, device='cpu', hn=0, hf=1,
14           nb_epochs=int(1e5),
15           nb_bins=192, H=400, W=400, fine_stage=2):
16
17     training_loss = []
18
19     for _ in tqdm(range(nb_epochs)):
20         for batch in data_loader:
21             #for each training point in the batch it exist in form [origin(x,y,z), direction(a1,
22             a2,a3), gt_pixel(r,g,b)]
23             ray_origins = batch[:, :3].to(device)
24             ray_directions = batch[:, 3:6].to(device)
25             ground_truth_px_values = batch[:, 6:].to(device)
26             if _ < fine_stage:
27                 regenerated_px_values = render_rays(nerf_model, ray_origins, ray_directions, hn=hn,
28                 hf=hf, nb_bins=nb_bins, stage='coarse')
29             else:
30                 regenerated_px_values = render_rays(nerf_model, ray_origins, ray_directions, hn=hn,
31                 hf=hf, nb_bins=nb_bins, stage='fine')
32             loss = ((ground_truth_px_values - regenerated_px_values) ** 2).sum()
33
34             optimizer.zero_grad()
35             loss.backward()
36             optimizer.step()
37             training_loss.append(loss.item())
38             scheduler.step()
39             for img_index in range(3):
40                 test(hn, hf, testing_dataset, img_index=img_index, nb_bins=nb_bins, H=H, W=W)
41
42     return training_loss
43
44 @torch.no_grad()
45 def test(hn, hf, dataset, chunk_size=5, img_index=0, nb_bins=192, H=400, W=400):
46     ray_origins = dataset[img_index * H * W: (img_index+1)* H * W, :3]
47     ray_directions = dataset[img_index * H * W: (img_index+1)* H * W, 3:6]
48
49     data = []
50
51     for i in range(int(np.ceil(H / chunk_size))):
52
53         ray_origins_ = ray_origins[i * W * chunk_size : (i+1) * W * chunk_size].to(device)
54         ray_directions_ = ray_directions[i * W * chunk_size : (i+1) * W * chunk_size].to(device)
55
56         regenerated_px_values = render_rays(model, ray_origins_, ray_directions_, hn=hn, hf=hf,
57             nb_bins = nb_bins)
58         #print("rpv:",regenerated_px_values)
59         data.append(regenerated_px_values)
60
61     img = torch.cat(data).data.cpu().numpy().reshape(H,W,3)
```

```

55 plt.figure()
56 plt.imshow(img)
57 plt.savefig(f'output/novel_views_2/img_{img_index}.png', bbox_inches='tight')
58 plt.close()
59
60 if __name__ == '__main__':
61     # get training argument
62     args = get_args()
63     device = args.device
64
65     # get dataset
66     # todo - load data to correct format [img.jpg + poses + ] -> [origin + direction +
67     # ground_truth] for training
68     training_dataset = torch.from_numpy(np.load('datasets/fufu/short_cut/training_data.pkl',
69     allow_pickle=True))
70     testing_dataset = torch.from_numpy(np.load('datasets/fufu/short_cut/training_data.pkl',
71     allow_pickle=True))
72
73     idx = torch.randperm(training_dataset.size(0))
74     training_dataset = training_dataset[torch.sort(idx).indices,:]
75     num_data_use = int(int(idx.size(0))/10)
76
77     #training_dataset = training_dataset[:num_data_use,:]
78     #print(testing_dataset)
79
80     # initialize NeRF
81     model = NerfModel(hidden_dim=256).to(device)
82     model_optimizer = torch.optim.Adam(model.parameters(), lr=5e-4)
83     scheduler = torch.optim.lr_scheduler.MultiStepLR(model_optimizer, milestones=[2, 4, 8],
84                                                     gamma=0.5)
85
86     data_loader = DataLoader(training_dataset, batch_size=16, shuffle=True) # fit data to
87     # batches
88
89     train(model, model_optimizer, scheduler, data_loader, nb_epochs=args.num_epochs, device=
90             device, hn=2, hf=6, nb_bins=100, H=288, W=216, fine_stage = args.fine_stage)

```

NeRF architecture (Nerf.py):

```

1 import torch
2 import torch.nn as nn
3
4 class NerfModel(nn.Module):
5     def __init__(self, embedding_dim_pos=10, embedding_dim_direction=4, hidden_dim=128):
6         super(NerfModel, self).__init__()
7
8         self.block1 = nn.Sequential(nn.Linear(6*embedding_dim_pos+3,hidden_dim), nn.ReLU(), #
9                                     including original pos
10                                     nn.Linear(hidden_dim,hidden_dim), nn.ReLU(),
11                                     nn.Linear(hidden_dim,hidden_dim), nn.ReLU(),
12                                     nn.Linear(hidden_dim,hidden_dim), nn.ReLU(),)
13
14         self.block2 = nn.Sequential(nn.Linear(6*embedding_dim_pos+hidden_dim+3,hidden_dim), nn.
15             ReLU(),
16                                     nn.Linear(hidden_dim,hidden_dim), nn.ReLU(),
17                                     nn.Linear(hidden_dim,hidden_dim), nn.ReLU(),
18                                     nn.Linear(hidden_dim,hidden_dim+1), )
19
20         self.block3 = nn.Sequential(nn.Linear(6*embedding_dim_direction+hidden_dim+3,hidden_dim
21 //2), nn.ReLU(),)
22         self.block4 = nn.Sequential(nn.Linear(hidden_dim // 2,3), nn.Sigmoid(),)
23
24         self.embedding_dim_pos = embedding_dim_pos
25         self.embedding_dim_direction = embedding_dim_direction
26         self.relu = nn.ReLU()

```

```

22
23     @staticmethod
24     def positional_encoding(x, L):
25         out = [x]
26         for j in range(L):
27             out.append(torch.sin(2**j * x))
28             out.append(torch.cos(2**j * x))
29         return torch.cat(out, dim=1)
30
31     def forward(self, o, d):
32         emb_x = self.positional_encoding(o, self.embedding_dim_pos)
33         emb_d = self.positional_encoding(d, self.embedding_dim_direction)
34         h = self.block1(emb_x)
35         tmp = self.block2(torch.cat((h, emb_x), dim=1))
36         h, sigma = tmp[:, :-1], self.relu(tmp[:, -1])
37         h = self.block3(torch.cat((h, emb_d), dim=1))
38         c = self.block4(h)
39         return c, sigma

```

rendering stage (render.py):

```

1   import torch
2   # utils
3   ###
4   def compute_accumulated_transmittance(alphas):
5       accumulated_transmittance = torch.cumprod(alphas, 1) #cumprod get the cumulative product
6       of alphas
7       return torch.cat((torch.ones((accumulated_transmittance.shape[0], 1), device=alphas.device
8           ),
9           accumulated_transmittance[:, :, -1]), dim = -1)
10
11
12 def PDF_reverse(sigma, hn, hf):
13     # x_axis: sigma [nb_samples, nb_bins]
14     device = sigma.device
15     nb_samples = sigma.shape[0]
16     nb_bins = sigma.shape[1]
17
18     # normalization
19     sigma = torch.cumsum(sigma, -1)
20     n_sigma = sigma.div((1e-10)+sigma.sum(-1).unsqueeze(-1))
21
22
23     t = torch.linspace(0, hf-hn, nb_bins, device=device).expand(nb_samples, nb_bins)
24     off = torch.linspace(hn, hn, nb_bins, device=device).expand(nb_samples, nb_bins)
25
26     t = t.mul(n_sigma)+off
27
28     return t
29
30
31 def coarse_sample(ray_origins, ray_directions, hn=0, hf=0.5, nb_bins=192, nb_samples=100):
32     device = ray_origins.device
33
34     t = torch.linspace(hn, hf, nb_bins, device=device).expand(nb_samples, nb_bins)
35     mid = (t[:, :-1]+t[:, 1:]) / 2
36     lower = torch.cat((t[:, :1], mid), -1)
37     upper = torch.cat((mid, t[:, -1:]), -1)
38     u = torch.rand(t.shape, device=device)
39     t = lower + (upper - lower) * u
40
41     delta = torch.cat((t[:, 1:] - t[:, :-1], torch.tensor([1e10], device=device).expand(
42         nb_samples, 1)), -1)
43
44     x = ray_origins.unsqueeze(1) + t.unsqueeze(2) * ray_directions.unsqueeze(1)

```

```

42 ray_directions = ray_directions.expand(nb_bins, ray_directions.shape[0], 3).transpose
43 (0,1)
44
45 return x, ray_directions, delta
46
47 def fine_sample(sigma, ray_origins, ray_directions, hn=0, hf=0.5, nb_bins=192, nb_samples
48 =100):
49 device = ray_origins.device
50
51 t = PDF_reverse(sigma, hn, hf)
52
53 mid = (t[:, :-1] + t[:, 1:]) / 2
54 lower = torch.cat((t[:, :1], mid), -1)
55 upper = torch.cat((mid, t[:, :-1]), -1)
56 u = torch.rand(t.shape, device=device)
57 t = lower + (upper - lower) * u
58
59 delta = torch.cat((t[:, 1:] - t[:, :-1], torch.tensor([1e10], device=device).expand(
60 nb_samples, 1)), -1)
61
62 x = ray_origins.unsqueeze(1) + t.unsqueeze(2) * ray_directions.unsqueeze(1)
63 ray_directions = ray_directions.expand(nb_bins, ray_directions.shape[0], 3).transpose
64 (0,1)
65
66 return x, ray_directions, delta
67
68 #####
69
70 def render_rays(nerf_model, ray_origins, ray_directions, hn=0, hf=0.5, nb_bins=192, stage =
71 'coarse'):
72 device = ray_origins.device
73 nb_samples = ray_origins.shape[0]
74
75 # coarse sampling
76 x_coarse, r_d_coarse, delta_coarse = coarse_sample(ray_origins, ray_directions, hn, hf,
77 nb_bins, nb_samples)
78 x, r_d, delta = x_coarse, r_d_coarse, delta_coarse
79
80 if stage == 'fine':
81     # fine sampling
82     c, sigma = nerf_model(x.reshape(-1,3), r_d.reshape(-1,3))
83     sigma = sigma.reshape(nb_samples, nb_bins)
84
85     x_fine, r_d_fine, delta_fine = fine_sample(sigma, ray_origins, ray_directions, hn, hf,
86 nb_bins, nb_samples)
87
88     # concatenating sampling
89
90     x = torch.cat((x_fine, x_coarse), -1)
91     r_d = torch.cat((r_d_fine, r_d_coarse), -1)
92     delta = torch.cat((delta_fine, delta_coarse), -1)
93
94     x, r_d, delta = x_fine, r_d_fine, delta_fine
95
96     #print("delta",delta)
97     # final rendering
98     colors, sigma = nerf_model(x.reshape(-1,3), r_d.reshape(-1,3))
99
100    colors = colors.reshape(x.shape)
101    sigma = sigma.reshape(x.shape[:-1])
102    #print("sigma",sigma)
103
104    alpha = 1 - torch.exp(-sigma * delta)
105
106    #print("alpha:",alpha)

```

```

100 weights = compute_accumulated_transmittance(1-alpha).unsqueeze(2) * alpha.unsqueeze(2)
101 c = (weights * colors).sum(dim=1)
102 weight_sum = weights.sum(-1).sum(-1)
103 return c + 1 - weight_sum.unsqueeze(-1)

```

Sampling from image (samples.py):

```

1
2 import torch
3 import numpy as np
4 import os
5 import cv2 as cv
6
7 def resize_img(img,scale):
8     scale_percent = scale # percent of original size
9     width = int(img.shape[1] * scale_percent / 100)
10    height = int(img.shape[0] * scale_percent / 100)
11    dim = (width, height)
12
13    # resize image
14    resized = cv.resize(img, dim, interpolation = cv.INTER_AREA)
15    return resized
16
17 def get_odg(img,u,v,E,hwf,z):
18     #o = -E[:3,:3].matmul(E[:3,3])
19     o = E[:3,3]
20     #px_2_world(hwf[0]/2,hwf[1]/2, E, hwf, 0)
21     #d = fp-o
22     d = torch.zeros(3,1)
23     d[0],d[1],d[2] = u-hwf[0]/2,v-hwf[1]/2,hwf[2]
24     d = d.to(torch.float32)
25     #d = E[:3,:3].transpose(0,1).matmul(d).reshape(-1)
26     d = E[:3,:3].matmul(d).reshape(-1)
27
28     #fp = px_2_world(u,v, E, hwf, z)
29     d = -d/(torch.sqrt((d*d).sum(-1))+1e-10) + 1e-10
30
31     g = img[u,v]/256
32     return o,d,g
33
34 def sample_tt(basedir, nm_cams, to_test=5):
35     #basedir = "/content/drive/MyDrive/github-project/My_Nerf/datasets/fufu"
36     nm_testing = nm_cams//to_test
37     nm_training = nm_cams - nm_testing
38
39     poses_arr = np.load(os.path.join(basedir, 'poses_bounds.npy'))
40
41     info = poses_arr[:, :-2].reshape([-1, 3, 5])#.transpose([1,2,0])
42     #bds = poses_arr[:, -2:].transpose([1,0])
43     z = 10 #int(torch.mean(bds))
44
45     imgs = [os.path.join(basedir, 'images', f) for f in sorted(os.listdir(os.path.join(
46         basedir, 'images')))\ \
47         if f.endswith('JPG') or f.endswith('jpg') or f.endswith('png')]
48
49     scale = 20
50     hwf = torch.from_numpy(info[0,:3,-1])
51     hwf[0], hwf[1] = hwf[0]*(scale/100), hwf[1]*(scale/100)
52     print(int(hwf[0]))
53     print(int(hwf[1]))
54
55     poses = torch.from_numpy(info[:, :3, :4])
56
57     Es = torch.zeros((poses.shape[0],1,4))

```

```

58 Es[:, :, -1] = 1
59 Es = torch.concat((poses, Es), axis = 1)
60
61 ray_one_img = int(hwf[0]*hwf[1])
62
63 nm_training *= ray_one_img
64 nm_testing *= ray_one_img
65
66 training = torch.ones([nm_training, 9])
67 testing = torch.ones([nm_testing, 9])
68
69 print(training.shape)
70 print(testing.shape)
71
72 off_training = 0
73 off_testing = 0
74
75 for i in range(nm_cams):
76     print(i)
77
78 E = Es[i].to(torch.float32)
79 img = cv.imread(imgs[i])
80 img = resize_img(img, scale)
81 img = torch.from_numpy(cv.resize(img, (int(hwf[1]), int(hwf[0]))))
82
83 off = 0
84 for j in range(int(hwf[0])):
85     for k in range(int(hwf[1])):
86         u, v = j, k
87         o, d, g = get_odg(img, u, v, E, hwf, z)
88
89     if (i+1)%to_test != 0:
90         # traning
91         training[off_training*ray_one_img + off, :3] = o
92         training[off_training*ray_one_img + off, 3:6] = d
93         training[off_training*ray_one_img + off, 6:] = g
94
95     else:
96         # testing
97         testing[off_testing*ray_one_img + off, :3] = o
98         testing[off_testing*ray_one_img + off, 3:6] = d
99         testing[off_testing*ray_one_img + off, 6:] = g
100
101     off += 1
102 if (i+1)%to_test != 0:
103     off_training += 1
104 else:
105     off_testing += 1
106
107 return training, testing

```

Visualise component/part (visualisation.py):

```

1 import numpy as np
2 import matplotlib as mpl
3 import matplotlib.pyplot as plt
4 from matplotlib.patches import Patch
5 from mpl_toolkits.mplot3d.art3d import Poly3DCollection
6 from pytorch3d.transforms import euler_angles_to_matrix
7
8 class Visualizer:
9     def __init__(self, xlim, ylim, zlim):
10         self.fig = plt.figure(figsize=(18, 7))
11         #self.ax = self.fig.gca(projection='3d')
12         self.ax = self.fig.add_subplot(projection='3d')

```

```

13     self.ax.set_aspect('auto')
14     self.ax.set_xlim(xlim)
15     self.ax.set_ylim(ylim)
16     self.ax.set_zlim(zlim)
17     self.ax.set_xlabel('x')
18     self.ax.set_ylabel('y')
19     self.ax.set_zlabel('z')
20     self.ax.view_init(45, 90)
21     print('initialize camera pose visualizer')
22
23 ## visulising camera
24
25 def extrinsic2pyramid(self, extrinsic, color='r', focal_len_scaled=5, aspect_ratio=-0.3):
26     vertex_std = np.array([[0, 0, 0, 1],
27                           [focal_len_scaled * aspect_ratio, -focal_len_scaled * aspect_ratio, focal_len_scaled, 1],
28                           [focal_len_scaled * aspect_ratio, focal_len_scaled * aspect_ratio, focal_len_scaled, 1],
29                           [-focal_len_scaled * aspect_ratio, focal_len_scaled * aspect_ratio, focal_len_scaled, 1],
30                           [-focal_len_scaled * aspect_ratio, -focal_len_scaled * aspect_ratio, focal_len_scaled, 1]])
31     vertex_transformed = vertex_std @ extrinsic.T
32     meshes = [[vertex_transformed[0, :-1], vertex_transformed[1, :-1],
33               vertex_transformed[2, :-1],
34               [vertex_transformed[0, :-1], vertex_transformed[2, :-1],
35               vertex_transformed[3, :-1]],
36               [vertex_transformed[0, :-1], vertex_transformed[3, :-1],
37               vertex_transformed[4, :-1]],
38               [vertex_transformed[0, :-1], vertex_transformed[4, :-1],
39               vertex_transformed[1, :-1],
40               [vertex_transformed[1, :-1], vertex_transformed[2, :-1],
41               vertex_transformed[3, :-1], vertex_transformed[4, :-1]]]
42     self.ax.add_collection3d(
43         Poly3DCollection(meshes, facecolors=color, linewidths=0.3, edgecolors=color,
44         alpha=0.35))
45
46 def customize_legend(self, list_label):
47     list_handle = []
48     for idx, label in enumerate(list_label):
49         color = plt.cm.rainbow(idx / len(list_label))
50         patch = Patch(color=color, label=label)
51         list_handle.append(patch)
52     plt.legend(loc='right', bbox_to_anchor=(1.8, 0.5), handles=list_handle)
53
54 def colorbar(self, max_frame_length):
55     cmap = mpl.cm.rainbow
56     norm = mpl.colors.Normalize(vmin=0, vmax=max_frame_length)
57     self.fig.colorbar(mpl.cm.ScalarMappable(norm=norm, cmap=cmap), orientation='vertical', label='Frame Number')
58
59 def show(self):
60     plt.title('Extrinsic Parameters')
61     plt.show()
62
63 def get_color(h):
64     hex_string = h[2:]
65     hex_string = '#' + '0' * (6 - len(hex_string)) + hex_string
66     return hex_string

```