

# Just Flip: Flipped Observation Generation and Optimization for Neural Radiance Fields to Cover Unobserved View

Minjae Lee, Kyeongsu Kang and Hyeonwoo Yu

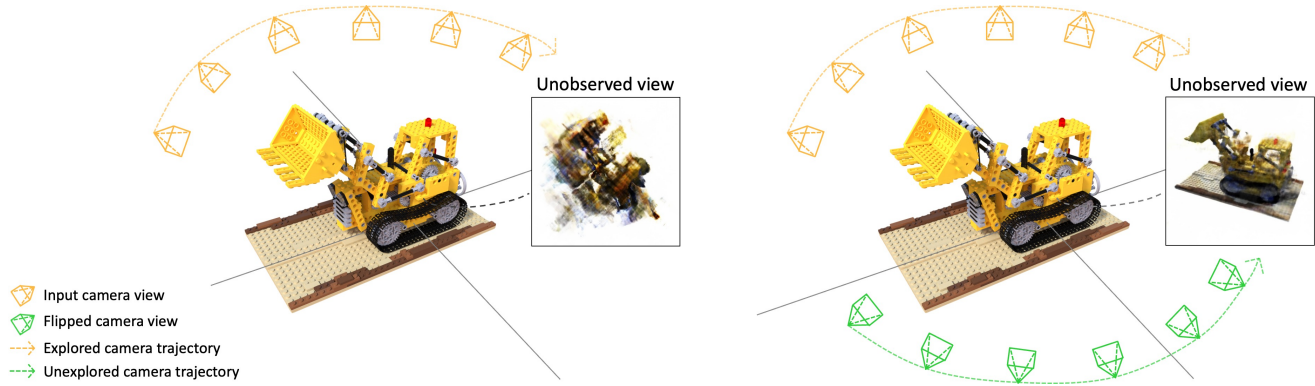


Fig. 1: Overview of our method. (Left) the baseline approach where the robot only observes one side of an object while driving. This case does not yield good rendering results in unobserved views that the robot has not explored. (Right) our method generates the flipped observations from the actual observations. The robot exploits both input images and flipped images and estimated camera poses to learn 3D space using NeRF for unexplored regions as well. Our method obtains qualified rendering results in unobserved views, even without providing images from unobserved views as a training set.

**Abstract**—With the advent of Neural Radiance Field (NeRF), representing 3D scenes through multiple observations has shown remarkable improvements in performance. Since this cutting-edge technique is able to obtain high-resolution renderings by interpolating dense 3D environments, various approaches have been proposed to apply NeRF for the spatial understanding of robot perception. However, previous works are challenging to represent unobserved scenes or views on the unexplored robot trajectory, as these works do not take into account 3D reconstruction without observation information. To overcome this problem, we propose a method to generate flipped observation in order to cover unexisting observation for unexplored robot trajectory. To achieve this, we propose a data augmentation method for 3D reconstruction using NeRF by flipping observed images, and estimating flipped camera 6DOF poses. Our technique exploits the property of objects being geometrically symmetric, making it simple but fast and powerful, thereby making it suitable for robotic applications where real-time performance is important. We demonstrate that our method significantly improves three representative perceptual quality measures on the NeRF synthetic dataset.

\*This work was supported by Institute of Information & communications Technology Planning & Evaluation(IITP) grant funded by the Korea government(MSIT) (No.2020-0-01336, Artificial Intelligence Graduate School Program(UNIST)), and by Basic Science Research Program through the National Research Foundation of Korea(NRF) funded by the Ministry of Education(2022R1F1A1066488).

Minjae Lee, Kyeongsu Kang and Hyeonwoo Yu are with the Department of Electrical Engineering & Graduate School of AI, Ulsan National Institute of Science and Technology (UNIST), Ulsan, South Korea. {lmjbsj, thithin, hyeonwoo.yu}@unist.ac.kr

The code will be available at: <https://github.com/minjae-lulu/Just-Flip>

## I. INTRODUCTION

Simultaneous Localization and Mapping (SLAM) is one of the key methods both in the fields of robotics and computer vision. SLAM is the process of generating a 3D map of an unknown environment while simultaneously estimating the position and movement of a robot or camera. This technique is being actively researched due to its potential applications in various real-world scenarios such as autonomous vehicles, unmanned robots, etc [1], [2], [3], [4]. Recently, in computer graphics, a novel technology called Implicit Neural Representation (INR) has emerged that utilizes neural networks to parameterize continuous and differentiable signals. As one of the INR, Neural Radiance Field (NeRF) [5] is a deep learning-based approach for 3D scene representation. In this approach, network learns 3D space by projecting a set of 2D views of a scene into a continuous 3D space, such that the color and density information from each view, thus new views of the scene is interpolated for 2D rendering. This technique can be applied to 3D mapping for SLAM as it can be exploited as a 3D scene representation. It also has several advantages over existing SLAM systems in that it can represent space continuously and is more memory-efficient as it utilizes neural networks to represent space. Due to these features, many recent works [6], [7], [8], [9], [10], [11] have applied NeRF to SLAM and 3D mapping.

However, despite many advantages of NeRF, similar to the previous 3D perception methods it also has limitations

in synthesizing unobserved views; from here, we define an unobserved view as a view that the robot has not explored or cannot explore. We also define the unobserved view problem as the task of estimating the synthesized view of an unexplored area. For example, in Fig. 1, the robot explores the orange camera trajectory while observing the lego truck, and the green cameras and trajectory represent unobserved views and the unexplored trajectory respectively. The view from the green cameras is unobserved views because they have not been explored yet. These unobserved views cannot be obtained from the observation trajectory due to the 2.5D observation of the robot and the self-occlusions of the objects. Since it is challenging to obtain observations of the unexplored region, various approaches based on NeRF show poor rendering performance on the unexplored scene. Previous researches such as [12] tried to overcome the occlusion problem by capturing visibility patterns of 3D voxel exemplars, but they are learning-based methods which require massive training dataset. The robot needs to be adaptive in environments where it encounters various objects, so NeRF achieves some advantages compared to the method based on the object dataset [12] while it does not require pretrainings in such challenging environments. However, as mentioned above, the robot is hard to achieve unobserved views while navigating and thus obtains limited input image data. Therefore, training NeRF with a few images become important and various data augmentation methods such as [13], [14] have been studied. Although these methods use geometric approaches to acquire views between or near explored views, they still have limitations in that they do not consider any information about unexplored areas, resulting in the unobserved view problem. Moreover, since they require model training for data augmentation, they may not be suitable for robotics applications that require real-time processing.

To address the unobserved view problem in NeRF, we propose a new unobserved view generation and optimization method. We assume that artificially created objects in indoor and outdoor environments mostly exhibit symmetrical 3D shapes [15], [16], [17]. Typical representative datasets also exhibit symmetric characteristics of 3D objects. Under this assumption, by flipping object observation we can generate additional significant information of unobserved views. Since training of NeRF requires both the images and camera poses, the challenge with just using the flipped image is that estimating camera poses of the flipped observation image are also required. Fortunately, some works [18], [19], [20], [21] that can train NeRF without camera poses have been introduced. However, these methods cannot guarantee finding the global minimum of the camera poses where the camera pose is completely randomly initialized. To relax this problem and improve the robustness of our method, we also propose a method to estimate the 6DOF camera pose of the flipped image using the camera pose of the input image. Our experiments show that flip and pose estimation method improves the performance of neural radiance fields on the NeRF synthetic dataset.

In summary, we propose a flipped observation generation method by considering robot trajectory to improve the understanding of unobserved views. Since many objects are geometrically symmetric, we use image flipping as a way to predict unobserved views. To train the NeRF model, both the image and its corresponding camera pose are required. However, the flipped image alone does not provide any information about the camera pose, which poses a challenge for using image flipping to predict unobserved views. To address this issue, we also propose a method to estimate the 6DOF camera pose of the flipped image to complete the flipped observation generation. Using this flipped observation, we perform NeRF training and camera pose fine-tuning simultaneously. Since we exploit the geometric characteristics of both 3D objects and robot trajectories, we can achieve the unobserved view training of NeRF, and prevent camera poses from falling into the local minimum. We demonstrated that our method shows significant improvements in performance of neural radiance fields on the NeRF synthetic dataset, especially for unobserved views. Our method is simple but strong and fast, making it suitable for robotics applications that require real-time performance.

## II. RELATED WORK

In the traditional SLAM techniques, 3D representation was achieved using methods such as mesh [22], [23], point clouds [24], [25], and depth maps [26], [27]. However, these approaches have the drawback of having memory limitations since they require storing discrete information, resulting in the map being represented in a sparse manner. To address these limitations, recent research [6], [7], [8], [9], [10], [11] has applied NeRF to SLAM for 3D mapping. NeRF defined a continuous function by a neural network that takes in a spatial location and viewing direction as inputs and outputs the RGB values and volume density. By continuous neural network function, their approach has the advantages of high resolution and low memory. NeRF has shown remarkable performance in synthesizing photorealistic novel views of real-world scenes or objects. Afterward, NeRF have been extended to achieve fast training [28], [29], [30], few input data [31], [32], large scale scene [33], [34], [35], dynamic scene capture [36], [37], scene editing [38], [39].

In addition, some works [18], [19], [20], [21] have been introduced that can train NeRF model in the unknown camera parameter condition. Since the flipped image does not have a camera pose, NeRF needs to be trained without camera pose or simultaneously estimate the camera pose and neural radiance fields. [18] performs pose estimation by inverting a trained neural radiance field. [19] optimizes both camera parameters and radiance field by two-stage approach. But they can only optimize camera pose for relatively short camera trajectories. [21] uses the bundle adjustment method and requires imperfect camera poses. GAN-based [20] can reconstruct neural radiance fields and estimate camera poses when the camera poses are randomly initialized.

To mitigate the over-fitting problem under few-shot settings, data augmentation works [13], [14] proposed to increase the number of input data. [13] uses geometric method view morphing [40] to generate images between views. And [14] provides artificial pose and image pairs for training. The rendering of pose is warped to the nearby training view with depth map and relative pose to match the image supervision. However, these methods are dependent on the input image pose, resulting in only generating nearby or between training views. In contrast to these methods, our flip method is capable of predicting unobserved views in regions that the robot has yet to explore.

### III. PRELIMINARY

NeRF is a deep learning-based approach for 3D scene reconstruction from 2D images. It aims to model a 3D scene as a continuous function that maps a 3D point in space to its color and intensity in the image plane. This function is represented as a neural network, which is trained on a set of input images  $\mathcal{I} = \{I_1, I_2, \dots, I_N\}$  of a scene, with their associated camera parameters  $\Pi = \{\pi_1, \pi_2, \dots, \pi_N\}$ . NeRF model is to learn a mapping from the 3D world coordinates  $\mathbf{x} = (x, y, z)$  and 2D viewing direction  $\mathbf{d} = (\theta, \phi)$  to the image radiance value  $\mathbf{c} = (r, g, b)$  and intensity  $\sigma$ . This mapping is represented by a neural network and the mapping function can be represented mathematically as  $F_\Theta : (\mathbf{x}, \mathbf{d}) \rightarrow (\mathbf{c}, \sigma)$ , where  $\Theta$  are network parameters.

In order to render an image, the color at each pixel  $\mathbf{p} = (u, v)$  on the image  $\hat{I}_i$  is obtained by a rendering function  $\mathcal{R}$ . The expected image color  $\hat{I}_i(p)$  of camera ray  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$  that starts from camera origin  $\mathbf{o}$  with near and far bound  $t_n$  and  $t_f$  is:

$$\hat{I}_i(p) = \mathcal{R}(\mathbf{p}, \pi_i | \Theta) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt,$$

where  $T(t) = \exp(-\int_{t_n}^t \sigma(\mathbf{r}(s)) ds)$ .

The function  $T(t)$  denotes the accumulated transmittance along the ray from  $t_n$  to  $t$ , *i.e.*, the probability that the ray travels from  $t_n$  to  $t$  without hitting any other particle. The network is trained to minimize the difference between the predicted image radiance values from the rendering function and the ground truth values obtained from the input images. The photometric loss  $L_p$  defined as follows:

$$L_p(\Theta, \Pi) = \sum_i^N \|I_i - \hat{I}_i\|_2^2. \quad (1)$$

### IV. METHOD

This section introduces a flip data generation technique for NeRF to accurately predict unobserved views. Symmetry is prevalent in real-world objects due to its functional and aesthetic benefits, and can be manufactured easily in some cases [15], [16], [17]. Based on this premise, our idea starts

flipping the input image to predict unobserved views. Our flip approach is simple, efficient, and doesn't require complex calculations or additional training, allowing it to operate quickly. Furthermore, our method does not depend on the relationship of the input images, it can be applied regardless of the number of input images. However, the flipped image does not have a camera pose, training NeRF in this scenario can be challenging.

To resolve this issue, we used the [20] algorithm, which simultaneously optimizes both the camera pose and NeRF network parameters. This algorithm starts the training by initializing the camera pose near the zero vector. Nonetheless, this initialization may lead to unstable pose estimation. To deal with this problem, we also propose a method for estimating the 6DOF camera pose of the flipped image using the existing camera pose. In the following, we present (i) NeRF without camera pose and (ii) flipped camera pose estimation.

#### A. NeRF without camera pose

To overcome challenging situations where camera poses are not available, we used an algorithm that can jointly estimate both camera poses  $\Pi$  and neural radiance fields  $F_\Theta$  when the cameras are initialized at random poses.

To generate a coarse estimate of the implicit scene representation and the camera poses, used a generator  $G$ , a discriminator  $D$ , and an inversion network  $E$ . The generator takes a random camera pose  $\Pi$  and synthesizes an image of the scene, which is then decomposed into patches and compared to real patches by the discriminator  $D$ . The generator and discriminator are then trained adversarial manner. An inversion network takes as input a patch from a generated image and outputs the corresponding camera pose. To train the inversion network, we used a self-supervised approach, which exploits the synthetic image patches and their corresponding camera poses as the training data. We denote the parameters of the inversion network  $E$  as  $\theta_E$ , and MSE loss function  $\mathcal{L}_E$  defined as follows:

$$\mathcal{L}_E(\theta_E) = \mathbb{E}_{\pi \sim P(\pi)} [\|E(G(\pi; F_\Theta); \theta_E) - \pi\|_2^2]. \quad (2)$$

Once we have obtained the coarse estimates of the implicit scene representation and camera poses, we further optimize the pose embedding and NeRF model by minimizing the photometric loss  $\mathcal{L}_p$ , as defined in Equation 1. To achieve more accurate results, we used a loss function that interleaves the refinement process and the estimation process. When we denote  $\lambda$  is the weighting coefficient, the loss is as follows.

$$\mathcal{L}_R(\Theta, \Pi) = \mathcal{L}_p(\Theta, \Pi) + \frac{\lambda}{n} \sum_{i=1}^n \|E(I_i; \theta_E) - \pi_i\|_2^2. \quad (3)$$

#### B. flipped camera pose estimation

While the above approach is capable of learning a model without requiring the camera pose, it still has certain limitations. One of the main issues is it may produce inaccurate pose estimates, especially when dealing with symmetrical objects. In such cases, the camera pose  $\pi^{init}$  initialized near

the zero vector tends to converge to the local minima, which may not be the global minima. This occurs because the learning process may not take place equally on both sides, resulting in a biased learning outcome. Another limitation is it may take a substantial amount of time to obtain a pose that is close to the true value. To address these issues, we propose a method for determining the 6DOF camera pose of the flipped image using the input information.

1) *Find the optimal sphere using least squares to pass through the input camera pose:* The optimal radius and center of the sphere can be found using the least square method when the input camera coordinates are given. The general equation of a sphere with its center at  $x_0, y_0, z_0$  and radius  $r$  is represented by  $(x - x_0)^2 + (y - y_0)^2 + (z - z_0)^2 = r^2$ . To use the least square method, the expanded and rearranged equation written as,  $x^2 + y^2 + z^2 = 2xx_0 + 2yy_0 + 2zz_0 + r^2 - x_0^2 - y_0^2 - z_0^2$ . The terms  $x_i, y_i, z_i$  refer to the coordinates of the first input camera pose, and the terms  $x_n, y_n, z_n$  refer to the coordinates of the final input camera pose in the training dataset. Using the  $\vec{f}$  vector, the  $A$  matrix, and the  $\vec{c}$  vector, we can consolidate the terms of the expanded sphere equation, followed by expressing the equation of a sphere as a matrix form  $\vec{f} = A\vec{c}$ , where:

$$\vec{f} = \begin{bmatrix} x_i^2 + y_i^2 + z_i^2 \\ x_{i+1}^2 + y_{i+1}^2 + z_{i+1}^2 \\ \vdots \\ x_n^2 + y_n^2 + z_n^2 \end{bmatrix},$$

$$A = \begin{bmatrix} 2x_i & 2y_i & 2z_i & 1 \\ 2x_{i+1} & 2y_{i+1} & 2z_{i+1} & 1 \\ \vdots & \vdots & \vdots & \vdots \\ 2x_n & 2y_n & 2z_n & 1 \end{bmatrix},$$

$$\vec{c} = \begin{bmatrix} x_0 \\ y_0 \\ z_0 \\ r^2 - x_0^2 - y_0^2 - z_0^2 \end{bmatrix}.$$

Now, through the above matrix equation, we can find the value of  $\vec{c}$  that minimizes the norm of the residual as the following:

$$\vec{c} = \underset{\vec{c}}{\operatorname{argmin}} (A\vec{c} - \vec{f})^T (A\vec{c} - \vec{f}).$$

With the  $\vec{c}$  value obtained by the least square method, we can obtain the optimized center point and radius of the sphere.

2) *Symmetrically transforming the input camera pose with respect to the symmetrical plane:* Once the optimized sphere has been obtained, the input camera pose is symmetrically transformed with respect to the symmetrical plane. Symmetrical objects have the property of having the same shape with respect to the symmetrical plane. As a result, the camera pose of the flipped image also becomes symmetrical with respect to the symmetrical plane. When given the input camera pose point  $(x_i, y_i, z_i)$  and the equation of the symmetrical plane  $Ax + By + Cz + D = 0$ , the coordinate of the reflected point  $(x', y', z')$  on the plane can be expressed as follows:

$$x'_i = x_i + 2A\left(-\frac{(Ax_i + By_i + Cz_i + D)}{A^2 + B^2 + C^2}\right),$$

$$y'_i = y_i + 2B\left(-\frac{(Ax_i + By_i + Cz_i + D)}{A^2 + B^2 + C^2}\right),$$

$$z'_i = z_i + 2C\left(-\frac{(Ax_i + By_i + Cz_i + D)}{A^2 + B^2 + C^2}\right).$$

3) *Project the symmetrically transformed points onto the sphere:* The sphere determined by the least square is only optimal, but it does not pass through all input camera pose points. Same way, the flipped camera points do not pass through the optimized sphere, so we need to project the symmetrical point  $(x', y', z')$  onto the sphere  $(x_0, y_0, z_0, r)$ . The coordinates of the projected point  $(x'', y'', z'')$  are as follows:

$$x'' = \alpha x', y'' = \alpha y', z'' = \alpha z',$$

where  $\alpha = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$ ,  $a = x'^2 + y'^2 + z'^2$ ,  $b = x'x_0 + y'y_0 + z'z_0$  and  $c = x_0^2 + y_0^2 + z_0^2 - r^2$ . Here, the obtained  $(x'', y'', z'')$  points represent a  $3 \times 1$  translation transformation. As a result, the translation  $T$  of the flipped camera pose can be written as follows:

$$T = [x'', y'', z'']^T.$$

4) *Estimate the rotation matrix of a flipped camera:* The camera pose should include not only translation information but also rotation information about the direction in which it is looking. If we have the camera coordinates  $\mathbf{c}$ , coordinates of the point being looked at, and the look-up vector  $\mathbf{up}$ , then we can calculate the rotation matrix of the camera. Our method set the estimated points as camera coordinates  $\mathbf{c} = (x'', y'', z'')$ , set the  $\mathbf{at}$  coordinates as the center point of the sphere  $\mathbf{at} = (x_0, y_0, z_0)$ , and set the look-up vector as  $\mathbf{up} = (0, 0, 1)$ . The  $3 \times 3$  rotation matrix  $R$  can be calculated using  $\mathbf{c}$ ,  $\mathbf{at}$ ,  $\mathbf{up}$  as follows:

- 1)  $\mathbf{z} = (\mathbf{c} - \mathbf{at}) / \|\mathbf{c} - \mathbf{at}\|$
- 2)  $\mathbf{x} = (\mathbf{up} \times \mathbf{z}) / \|\mathbf{up} \times \mathbf{z}\|$
- 3)  $\mathbf{y} = \mathbf{z} \times \mathbf{x}$
- 4)  $R = [\mathbf{x} \ \mathbf{y} \ \mathbf{z}]$

We obtained the camera's translation vector  $T$  and rotation matrix  $R$ . Through this, we can construct the estimated camera pose  $\pi^*$  as a 4x4 matrix, which then allows us to redefine the Equation 2 and 3 accordingly. The estimated camera pose  $\pi^*$  and the redefined equations are as follows:

$$\pi^* = \begin{bmatrix} R & T \\ 0 & 1 \end{bmatrix},$$

$$\mathcal{L}_E(\theta_E) = \mathbb{E}_{\pi' \sim P(\pi')} [\|E(G(\pi'; F_\Theta); \theta_E) - \pi'\|_2^2],$$

$$\mathcal{L}_R(\Theta, \Pi') = \mathcal{L}_p(\Theta, \Pi') + \frac{\lambda}{n} \sum_{i=1}^n \|E(I'_i; \theta_E) - \pi'_i\|_2^2,$$

TABLE I: Quantitative comparisons of our flip methods with baseline method and baseline with ground-truth method on the NeRF synthetic dataset. Ground-truth image is an image observed from a viewpoint that the robot has not explored. Our method shows improved performance, even without estimating the flipped camera poses. In addition, we demonstrate that applying camera pose estimation to flipped images yields a more substantial improvement in performance compared to using the simple flip method.

| Scene      | PSNR $\uparrow$ |       |              |         | SSIM $\uparrow$ |             |             |         | LPIPS $\downarrow$ |      |             |         |
|------------|-----------------|-------|--------------|---------|-----------------|-------------|-------------|---------|--------------------|------|-------------|---------|
|            | Base            | Ours  | Ours+Pose    | Base+GT | Base            | Ours        | Ours+Pose   | Base+GT | Base               | Ours | Ours+Pose   | Base+GT |
| Chairs     | 11.10           | 12.80 | <b>14.87</b> | 16.17   | 0.64            | 0.74        | <b>0.81</b> | 0.86    | 0.48               | 0.34 | <b>0.24</b> | 0.24    |
| Ficus      | 17.05           | 16.17 | <b>17.73</b> | 18.43   | 0.82            | <b>0.82</b> | 0.75        | 0.84    | <b>0.21</b>        | 0.26 | 0.29        | 0.16    |
| Hotdog     | 13.37           | 13.54 | <b>16.27</b> | 17.20   | 0.67            | 0.74        | <b>0.76</b> | 0.80    | 0.45               | 0.36 | <b>0.35</b> | 0.30    |
| Lego truck | 11.25           | 13.53 | <b>14.17</b> | 23.13   | 0.68            | 0.74        | <b>0.77</b> | 0.88    | 0.48               | 0.39 | <b>0.30</b> | 0.11    |
| Materials  | 15.24           | 15.87 | <b>16.88</b> | 17.42   | 0.66            | 0.68        | <b>0.80</b> | 0.81    | 0.42               | 0.42 | <b>0.29</b> | 0.30    |
| Ship       | 10.67           | 10.93 | <b>17.98</b> | 23.59   | 0.63            | 0.64        | <b>0.72</b> | 0.85    | 0.54               | 0.51 | <b>0.31</b> | 0.16    |
| Mean       | 13.11           | 13.81 | <b>16.32</b> | 19.32   | 0.68            | 0.73        | <b>0.77</b> | 0.84    | 0.43               | 0.38 | <b>0.30</b> | 0.21    |

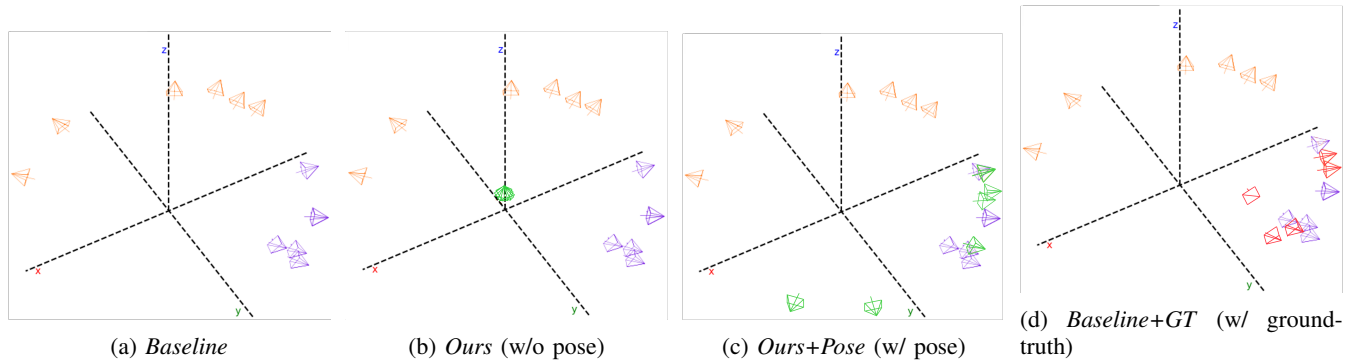


Fig. 2: Visualization of initial camera poses setting on each method in the Ship scene. In our experiment, the NeRF model is trained using the images from robot camera views (orange) and the flipped camera views (green) as inputs to predict the test image set of the test camera views (purple). (a) *Baseline* method only uses the images that the robot has acquired while navigating. (b) *Ours* method uses both the images obtained by the robot and their horizontally flipped versions for training. The flipped images are initialized with a zero vector as flipped camera pose is unknown. (c) *Ours+Pose* method estimates the flipped camera pose and adds it for training. (d) *Baseline+GT* method uses the images obtained by both robot and camera view (red) of the ground-truth image. We set ground-truth images as close as possible to the purple camera view because there is no one-to-one mapping between the train and test dataset.

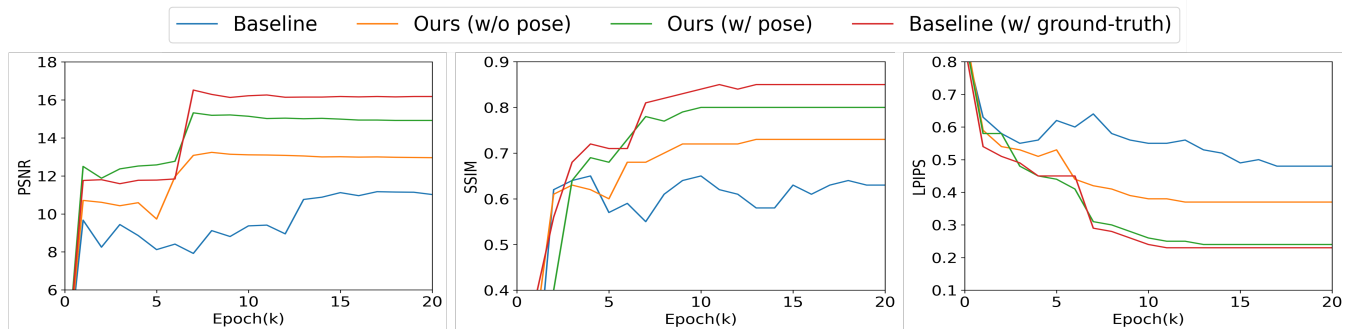


Fig. 3: Comparisons converge time of our flip methods with baseline method and baseline with ground-truth method. The baseline method takes around 14k epochs to converge, while our methods that flip the images take half of that, about 7k epochs.

where  $\mathcal{I}'$  is a new image set consisting of the original images  $\{I_i, \dots, I_n\}$  and flipped images  $\{I'_i, \dots, I'_n\}$ ,  $\Pi'$  is a new set including the original camera poses  $\{\pi_i, \dots, \pi_n\}$  and estimated camera poses  $\{\pi'_i, \dots, \pi'_n\}$ . In the next section, we compare our method that uses the redefined sets  $\mathcal{I}'$  and  $\Pi'$ , with the baseline methods.

## V. EXPERIMENTS

### A. Implementation and Setup

We used the NeRF synthetic dataset [5] as our data set and conducted experiments on Chairs, Ficus, Hotdog, Lego truck, Materials, and Ship. We compared four methods; baseline, ours(flip w/o pose estimation), ours(flip w/ pose estimation), baseline with ground-truth (using ground-truth images of unexplored areas), which will be discussed in

detail in the analysis. We use [20] as the baseline, which can learn both the neural radiance field and optimize the camera pose simultaneously. Considering the input data scenario, we make the assumption that the robot moves at a steady speed when collecting image data and is restricted to observing only one side of the object. Consequently, the  $z$  value of the camera pose in the input image remains mostly constant and aligned with the robot trajectory as shown in Fig. 1 and Fig. 2. As our approach strives to predict the view that cannot be observed, we utilize the unobserved view from the opposite side as the validation and ground-truth image. For quantitative comparison, we used three representative metrics: Peak Signal-to-Noise Ratio (PSNR) [41], Structural Similarity Index Measure (SSIM) [42], and Learned Perceptual Image Patch Similarity (LPIPS) [43]. We also compared the convergence speed of these values for each method. The parameters used for training were set the same way as our baseline model, except the learning rate for the camera pose which was adaptively applied to allow for more flexible adjustments.

### B. Analysis of the results

Before conducting the quantitative analysis, we provide a description of the four methods in detail. Our dataset consists of four sets: the input image set obtained from the robot exploration, the flipped image set, the test image set, and the ground-truth image set for the unexplored scene. Each set contains six images, and we assume that the corresponding camera poses are already known except for the camera pose of the flipped image. The input image set consists of data that the robot has acquired while navigating, whereas the flipped image set is created by flipping the input images. The test set comprises unobserved views that the robot cannot explore. Since there is no one-to-one mapping between the train and test datasets, we construct the ground-truth data by choosing the closest to the test data views as shown in Fig. 2d. The baseline method in Fig. 2a is trained using the image set obtained while the robot is navigating through the  $y < 0$  region. The proposed method without camera pose estimation in Fig. 2b, on the other hand, leverages both the input image and its flipped version. However, since there is no information available about the flipped camera pose, we initialize the flipped camera pose with a zero vector for training. In contrast, the proposed method with camera pose flipping in Fig. 2c uses both datasets and utilizes the least squares method to estimate the flipped camera pose for training. Meanwhile, the other method in Fig. 2d trains the model using both the input image, ground-truth image, and corresponding ground-truth camera pose. This approach is an upper bound since it directly uses the ground-truth image of the unobserved view in training. Based on this, the four methods can be summarized in the table below.

In TABLE I, we presented the metric results for the four methods on the NeRF synthetic dataset. The result demonstrates that (c) *Ours+Pose* method has achieved significant performance improvements compared to the other approaches. Nevertheless, we have observed that performance

TABLE II: Summary of the four experimental methods.

|                      |   |
|----------------------|---|
| (a) <i>Base</i>      | input image(N=6)                                |
| (b) <i>Ours</i>      | input image(N=6) + flipped image w/o pose (N=6) |
| (c) <i>Ours+Pose</i> | input image(N=6) + flipped image w pose (N=6)   |
| (d) <i>Base+GT</i>   | input image(N=6) + ground-truth image (N=6)     |

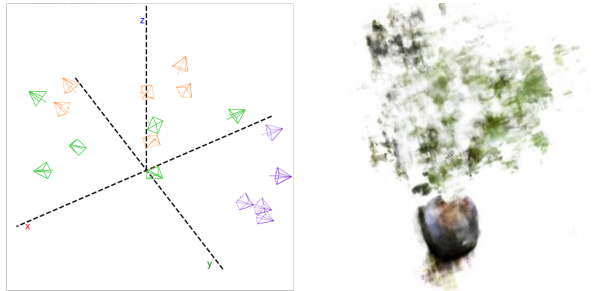


Fig. 4: Example of fail cases. (Left) The camera pose results when training with the initial camera pose set as shown in Fig. 2b in the Ship scene. (Right) The result of learning in the wrong direction due to the non-symmetric property in the Ficus scene.

improvement was not achieved in some cases. The left image of Fig. 4 shows the camera pose estimation result by zero-initialization as shown in Fig. 2b. Intuitively, the flipped images should be moved to the unexplored area,  $y > 0$ . However, it is failed to achieve the global minimum of the flipped camera poses due to the zero-vector initialization. In addition, the ship dataset was only illuminated on one side, leading to the opposite side being shadowed. In this situation, the camera pose has difficulty moving towards  $y > 0$  during training, which resulted in the model being trained to the local minimum point. Therefore, as shown in TABLE I, the ship scene did not achieve performance improvements in (b) *Ours* method compared to (a) *Baseline* method. To relax this problem, we use the (c) *Ours+Pose* method which involves performing pose estimation. As a result, we finally obtain an increase in the PSNR value from 10 to 17. From the results, we found that accurate camera pose initialization is crucial for achieving better performance. In the case of Ficus, it was observed that the results of (b) *Ours*, which did not involve pose estimation, had lower PSNR and LPIPS performance compared to (a) *Baseline*. This can be attributed to the fact that Ficus has complex elements such as branches and leaves that are not symmetric. These non-symmetric elements make it challenging for the model to be trained correctly, resulting in lower performance. Regarding these issues, we propose future work in the discussion.

Fig. 3 provides information on the convergence speed times for each of the methods. (a) *Baseline* converged at about 14k epochs, while (b) *Ours* achieved converged at around 7k epochs. It can be observed that (b) *Ours* and (c) *Ours+Pose* converge at a similar speed to the (d) *Baseline+GT*. In short, the experimental results indicate that using our flipped generation method increases the convergence speed about two times, which is comparable to the method that includes ground-truth images.



Fig. 5: Input image efficiency example. (Left) Flipping the image of the front view only provides information about the already-known view. (Right) In contrast, flipping the image of the side view provides additional information about the unobserved view.

Our flip method has demonstrated promising results on the novel view synthesis task, assuming that the objects are symmetrical and the symmetrical plane is known. However, for scenes with non-symmetrical objects like ficus scenes, we cannot expect improvement in performance. As a result, we propose two future directions for improvement. Firstly, we need to determine whether an object is symmetrical and obtain information about the symmetrical plane through input image data. This approach provides us with an assessment of whether the flip method is suitable and prevents the model from learning in the wrong direction. Secondly, we need to obtain information about which side of the object is being viewed from the input image. The effectiveness of data generation may depend on the side of the object being viewed. Consider an example where we are given a lego truck image of the front side, as shown in Fig. 5. In such a scenario, the flipped image will not provide any additional information about the unobserved view and such image is not useful for training. To provide efficient images for training, we need a method that can determine which side of the object is being viewed.

## VI. CONCLUSION

In this paper, we presented a flipped observation generation method of NeRF to predict unobserved views. By flipping the observed images, we demonstrate that our method can efficiently generate unobserved views. Since the flipped images do not have camera poses, both camera poses and neural radiance fields should be optimized simultaneously. However, when the camera poses are randomly initialized, it is not guaranteed to find the global minimum of the camera poses which leads to the failure of training NeRF. To overcome these problems and achieve a robust approach, we also proposed a flipped camera 6DOF pose estimation method. Experimental results show that our method enables significant performance improvements in predicting scenes of unobserved view. Our flipped observation method is suitable for robotic applications such as SLAM where real-time is important since our approach is simple but robust, and also fast.

- [1] M. Filipenko and I. Afanasyev, “Comparison of various slam systems for mobile robot in an indoor environment,” in *2018 International Conference on Intelligent Systems (IS)*. IEEE, 2018, pp. 400–407.
- [2] C. Yu, Z. Liu, X.-J. Liu, F. Xie, Y. Yang, Q. Wei, and Q. Fei, “Ds-slam: A semantic visual slam towards dynamic environments,” in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 1168–1174.
- [3] H. Yu, J. Moon, and B. Lee, “A variational observation model of 3d object for probabilistic semantic slam,” in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 5866–5872.
- [4] M. Mittal, R. Mohan, W. Burgard, and A. Valada, “Vision-based autonomous uav navigation and landing for urban search and rescue,” in *Robotics Research: The 19th International Symposium ISRR*. Springer, 2022, pp. 575–592.
- [5] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, “Nerf: Representing scenes as neural radiance fields for view synthesis,” *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.
- [6] E. Sucar, S. Liu, J. Ortiz, and A. J. Davison, “imap: Implicit mapping and positioning in real-time,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 6229–6238.
- [7] Z. Zhu, S. Peng, V. Larsson, W. Xu, H. Bao, Z. Cui, M. R. Oswald, and M. Pollefeys, “Nice-slam: Neural implicit scalable encoding for slam,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 786–12 796.
- [8] A. Rosinol, J. J. Leonard, and L. Carlone, “Nerf-slam: Real-time dense monocular slam with neural radiance fields,” *arXiv preprint arXiv:2210.13641*, 2022.
- [9] X. Yang, H. Li, H. Zhai, Y. Ming, Y. Liu, and G. Zhang, “Vox-fusion: Dense tracking and mapping with voxel-based neural implicit representation,” in *2022 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. IEEE, 2022, pp. 499–507.
- [10] X. Zhong, Y. Pan, J. Behley, and C. Stachniss, “Shine-mapping: Large-scale 3d mapping using sparse hierarchical implicit neural representations,” *arXiv preprint arXiv:2210.02299*, 2022.
- [11] X. Kong, S. Liu, M. Taher, and A. J. Davison, “vmap: Vectorised object mapping for neural field slam,” *arXiv preprint arXiv:2302.01838*, 2023.
- [12] Y. Xiang, W. Choi, Y. Lin, and S. Savarese, “Data-driven 3d voxel patterns for object category recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1903–1911.
- [13] D. Chen, Y. Liu, L. Huang, B. Wang, and P. Pan, “Geoaug: Data augmentation for few-shot nerf with geometry constraints,” in *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XVII*. Springer, 2022, pp. 322–337.
- [14] M. Bortolon, A. Del Bue, and F. Poiesi, “Data augmentation for nerf: a geometric consistent solution based on view morphing,” *arXiv preprint arXiv:2210.04214*, 2022.
- [15] S. Tulsiani, A. Kar, Q. Huang, J. Carreira, and J. Malik, “Shape and symmetry induction for 3d objects,” *arXiv preprint arXiv:1511.07845*, 2015.
- [16] N. J. Mitra, M. Pauly, M. Wand, and D. Ceylan, “Symmetry in 3d geometry: Extraction and applications,” in *Computer Graphics Forum*, vol. 32, no. 6. Wiley Online Library, 2013, pp. 1–23.
- [17] S. Wu, C. Rupprecht, and A. Vedaldi, “Unsupervised learning of probably symmetric deformable 3d objects from images in the wild,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1–10.
- [18] L. Yen-Chen, P. Florence, J. T. Barron, A. Rodriguez, P. Isola, and T.-Y. Lin, “inertf: Inverting neural radiance fields for pose estimation,” in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 1323–1330.
- [19] Z. Wang, S. Wu, W. Xie, M. Chen, and V. A. Prisacariu, “Nerf: Neural radiance fields without known camera parameters,” *arXiv preprint arXiv:2102.07064*, 2021.
- [20] Q. Meng, A. Chen, H. Luo, M. Wu, H. Su, L. Xu, X. He, and J. Yu, “Gnerf: Gan-based neural radiance field without posed camera,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 6351–6361.

- [21] C.-H. Lin, W.-C. Ma, A. Torralba, and S. Lucey, “Barf: Bundle-adjusting neural radiance fields,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 5741–5751.
- [22] G. Riegler and V. Koltun, “Free view synthesis,” in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIX 16*. Springer, 2020, pp. 623–640.
- [23] ———, “Stable view synthesis,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 12 216–12 225.
- [24] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, “Pointnet: Deep learning on point sets for 3d classification and segmentation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 652–660.
- [25] P. Achlioptas, O. Diamanti, I. Mitliagkas, and L. Guibas, “Learning representations and generative models for 3d point clouds,” in *International conference on machine learning*. PMLR, 2018, pp. 40–49.
- [26] F. Liu, C. Shen, G. Lin, and I. Reid, “Learning depth from single monocular images using deep convolutional neural fields,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 10, pp. 2024–2039, 2015.
- [27] P.-H. Huang, K. Matzen, J. Kopf, N. Ahuja, and J.-B. Huang, “Deepmvs: Learning multi-view stereopsis,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2821–2830.
- [28] A. Chen, Z. Xu, F. Zhao, X. Zhang, F. Xiang, J. Yu, and H. Su, “Mvs-nerf: Fast generalizable radiance field reconstruction from multi-view stereo,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 14 124–14 133.
- [29] Q. Xu, Z. Xu, J. Philip, S. Bi, Z. Shu, K. Sunkavalli, and U. Neumann, “Point-nerf: Point-based neural radiance fields,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5438–5448.
- [30] T. Müller, A. Evans, C. Schied, and A. Keller, “Instant neural graphics primitives with a multiresolution hash encoding,” *ACM Transactions on Graphics (ToG)*, vol. 41, no. 4, pp. 1–15, 2022.
- [31] A. Yu, V. Ye, M. Tancik, and A. Kanazawa, “pixelnerf: Neural radiance fields from one or few images,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4578–4587.
- [32] K. Deng, A. Liu, J.-Y. Zhu, and D. Ramanan, “Depth-supervised nerf: Fewer views and faster training for free,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 882–12 891.
- [33] Y. Xiangli, L. Xu, X. Pan, N. Zhao, A. Rao, C. Theobalt, B. Dai, and D. Lin, “Citynerf: Building nerf at city scale,” *arXiv preprint arXiv:2112.05504*, 2021.
- [34] M. Tancik, V. Casser, X. Yan, S. Pradhan, B. Mildenhall, P. P. Srinivasan, J. T. Barron, and H. Kretzschmar, “Block-nerf: Scalable large scene neural view synthesis,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 8248–8258.
- [35] H. Turki, D. Ramanan, and M. Satyanarayanan, “Mega-nerf: Scalable construction of large-scale nerfs for virtual fly-throughs,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 922–12 931.
- [36] K. Park, U. Sinha, J. T. Barron, S. Bouaziz, D. B. Goldman, S. M. Seitz, and R. Martin-Brualla, “Nerfies: Deformable neural radiance fields,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 5865–5874.
- [37] G. Yang, M. Vo, N. Neverova, D. Ramanan, A. Vedaldi, and H. Joo, “Banmo: Building animatable 3d neural models from many casual videos,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 2863–2873.
- [38] C. Wang, M. Chai, M. He, D. Chen, and J. Liao, “Clip-nerf: Text-and-image driven manipulation of neural radiance fields,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 3835–3844.
- [39] Y.-J. Yuan, Y.-T. Sun, Y.-K. Lai, Y. Ma, R. Jia, and L. Gao, “Nerf-editing: geometry editing of neural radiance fields,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 18 353–18 364.
- [40] S. M. Seitz and C. R. Dyer, “View morphing,” in *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, 1996, pp. 21–30.
- [41] J. Korhonen and J. You, “Peak signal-to-noise ratio revisited: Is simple beautiful?” in *2012 Fourth International Workshop on Quality of Multimedia Experience*. IEEE, 2012, pp. 37–38.
- [42] J. C. Yue and M. K. Clayton, “A similarity measure based on species proportions,” *Communications in Statistics-theory and Methods*, vol. 34, no. 11, pp. 2123–2131, 2005.
- [43] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, “The unreasonable effectiveness of deep features as a perceptual metric,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 586–595.