
ENHANCEMENT OF NOVEL VIEW SYNTHESIS USING OMNIDIRECTIONAL IMAGE COMPLETION

Takayuki Hara
The University of Tokyo
hara@mi.t.u-tokyo.ac.jp

Tatsuya Harada
The University of Tokyo / RIKEN
harada@mi.t.u-tokyo.ac.jp

ABSTRACT

We present a method for synthesizing novel views from a single 360-degree image based on the neural radiance field (NeRF). Prior studies rely on the neighborhood interpolation capability of multi-layer perceptrons to complete missing regions caused by occlusion and zooming, and this leads to artifacts. In the proposed method, the input image is reprojected to 360-degree images at other camera positions, the missing regions of the reprojected images are completed by a self-supervised trained generative model, and the completed images are utilized to train the NeRF. Because multiple completed images contain inconsistencies in 3D, we introduce a method to train NeRF while dynamically selecting a sparse set of completed images, to reduce the discrimination error of the synthesized views with real images. Experiments indicate that the proposed method can synthesize plausible novel views while preserving the features of the scene for both artificial and real-world data.

Keywords scene representation, view synthesis, neural radiance field, image completion, 360-degree image, 3D deep learning

1 Introduction

Novel view synthesis from a set of captured images has a wide range of application which include AR/VR and immersive 3D photography. Conventionally, structure-from-motion [19] and image based rendering [45] have been employed for this task. In recent years, neural networks-based rendering methods have been rapidly developed, and the neural radiance field (NeRF) [37] is a promising method for synthesizing photorealistic views. However, the NeRF requires tens to hundreds of images with known relative positions and the same shooting conditions to be given as input, and such imaging is a large and time-consuming process. Accordingly, various efforts have been made to reduce the number of input images [62, 58, 57, 23] or ease the shooting conditions [36, 59, 32, 24].

With this background, we attempt to learn a 3D scene model from a single 360-degree image. Recently, 360-degree cameras, which can capture the entire field of view in a single shot, have become more easily accessible, and significant amounts of 360-degree image data are also available through social media and public datasets. Learning NeRF from a single 360-degree image is advantageous, in that we do not need to align the shooting conditions between images, nor do we need to know the relative positions between images, because we use only one image that contains a wealth of omnidirectional information. OmniNeRF [21] is a prior study of this approach, however, it relies only on the neighborhood interpolation capability of the multi-layer perceptron to complete the missing regions caused by occlusion and zooming. This leads to artifacts, and the image quality is significantly reduced when moving away from the camera position of the input image.

In contrast, the technology of completing the missing regions of 2D images has been studied for a long time, such as inpainting or image completion. Recent learning-based such as generative adversary networks (GAN) [16], variational autoencoders (VAE) [27, 43], and diffusion models [46, 20] have made it possible to generate semantically high-quality images. In addition, 360-degree image completion has been well-researched [51, 1, 17, 18], and high-quality image completion is possible over the entire field of view. In this study, we attempt to synthesize more plausible novel views from a single 360-degree image by combining NeRF and image completion using 360-degree images.

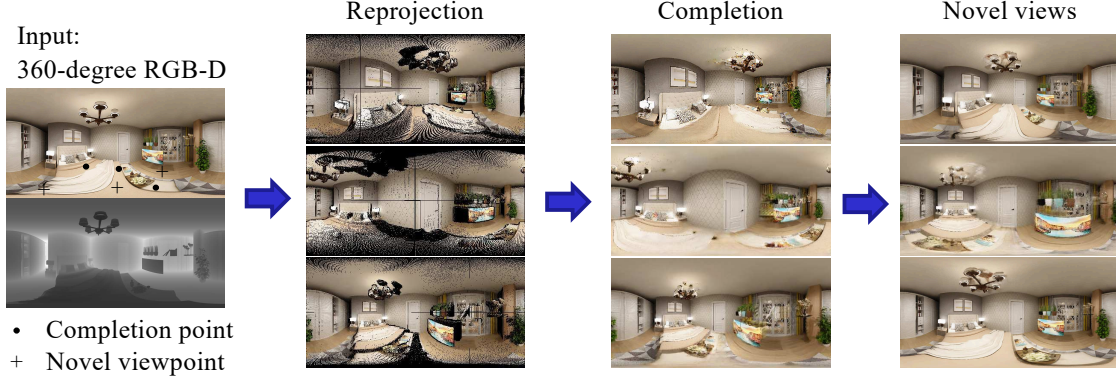


Figure 1: Given only a single 360 RGB-D image, our method is able to render novel views. The input image is reprojected to 360-degree images at another camera positions, and the missing regions of the reprojected images are completed. A sparse set of the dynamically selected completed images is used to train the NeRF, and novel views are synthesized. Note that the 360-degree image is represented by the equirectangular projection that maps the longitude of the viewing direction to the horizontal coordinate and the latitude to the vertical coordinate.

However, 2D image completions generally do not consider the 3D structure, and thus do not have 3D consistency when simply combined with NeRF. To maintain 3D consistency, we introduce a method to learn NeRF while dynamically selecting a sparse set of the completed images to reduce the discrimination error with real images. Figure 1 illustrates the overview of our method.

Our contributions in this paper are as follows.

- We propose a method for synthesizing novel views by learning NeRF from a single 360-degree RGB-D image, which improves image quality by employing the 360-degree image completion.
- We design a novel architecture that can select a sparse set of completed images at different camera positions and train NeRF simultaneously.
- We demonstrate that our proposed method can synthesize more plausible views for artificial and real-world data, using new metrics for evaluating the quality of novel views for which no ground truth image exists.

2 Related Work

2.1 Novel view synthesis

Novel view synthesis from a set of captured images has been a persistent challenge in computer vision. Structure-from-motion [19] is a method that has been employed for a long time for the novel view synthesis. It can recover point clouds and camera parameters from images of multiple viewpoints. Mesh-based methods have been employed to generate object-centric views [25, 4]; however it is difficult and impractical to represent an entire scene with a detailed mesh in the real world. The representation method using multi-plane images (MPI) [71, 44] is another way to achieve novel view synthesis in a limited field of view without the need of 3D model. Attal et al. [3] have extended MPI to multi-sphere images so that the entire field of view can be rendered, but this is not suitable for free movement within a scene. Other approaches to view synthesis include image-based rendering [45] and light-field photography [30]; however these generally require dense capture of the scene.

Recently, rendering techniques using neural networks, that map 3D spatial location to an implicit representation, are applied to this task [56]. In particular, the neural radiance field (NeRF) [37] can render objects with complex shapes and textures in a high-quality and photorealistic manner. However, the NeRF requires tens to hundreds of images with known relative positions and the same shooting conditions to be given as input, and such imaging is a large and time-consuming process. Accordingly, various efforts have been made to reduce the number of input images [62, 58, 57, 23], ease the shooting conditions [36, 59, 32, 24], speed up processing [49, 14, 41, 61, 38, 9], and express the entire scene [66, 12, 7, 42, 55].

Among them, OmniNeRF [21], which learns an entire scene from a single 360-degree RGBD image, without the need to set relative positions or identify shooting conditions. However, it only relies on the neighborhood interpolation capability of the multi-layer perceptron to complete the missing regions caused by occlusion and zooming, which

leads to artifacts, and the image quality is greatly reduced when moving away from the camera position of the input image. An alternative method to NeRF, pathdreamer [28], also synthesizes novel views from a single 360-degree RGB-D image. However, because this method is based on 2D image-to-image translation [39], it has the issue of low 3D consistency in the synthesized views.

2.2 Image Completion

Several image completion technologies have been proposed thus far for predicting the missing regions of an image. Conventionally, numerous diffusion-based methods [5, 8] diffuse the information of the visible regions into the missing regions, and multiple patch-based methods [11, 6] complete the missing regions by matching, copying, and realignment using visible regions. Both methods assume that the missing regions contain information correlated with the visible regions.

Generative models, which are trained using large-scale datasets, such as a variational autoencoder (VAE) [27, 43] and generative adversarial networks (GANs) [16], have experienced a significant boost, and both types of models have been adopted for image completion. Li et al. [31] directly generated the content of missing regions using convolutional neural networks (CNNs) [13, 29] with a combination of reconstruction loss, semantic parsing loss, and two adversarial losses. Iizuka et al. [22] employed global and local context discriminators within the framework of adversarial learning to improve the naturalness and consistency of the completed regions. Furthermore, improved methods for convolutional layers for image completion were also proposed in [33, 64]. Applying the methodology of the traditional image completion methods, Shift-Net [63] and contextual attention [60] were proposed to concatenate highly correlated features within the visible region to features within the missing region. Recently, LaMa [53], which expands the receptive field using fast Fourier convolutions, and RePaint [35], which utilizes a diffusion model [46, 20], have been proposed.

Although most image-completion methods produce only one result for each input, in [69, 67, 68, 34, 40], methods for generating multiple and diverse plausible solutions for image completion utilize conditional VAEs (CVAEs) or conditional GANs (CGANs).

360-degree image completion has also been researched. Han et al. [17] proposed an image-inpainting method for spherical structures using a cube map. In [52, 51], a 360-degree image is generated from a set of images captured in multiple directions as an input. Panoramic three-dimensional structure prediction methods have also been proposed [48, 50]. In [15, 47, 2, 1, 18], a 360-degree image is generated from a single normal field of view (NFOV) image. Thus, technology has been developed to complete the missing regions of 360-degree images with high image quality, and we apply these achievements to the novel view synthesis.

3 Preliminaries

This section provides an overview of NeRF and OmniNeRF, which forms the basis of this study. NeRF employs a multi-layer perceptron to construct a function that takes the 3D position $x \in \mathbb{R}^3$ and a unit-norm viewing direction $d \in S^2$ as the input, and outputs density $\sigma \in \mathbb{R}$ and color $c \in \mathbb{R}^3$. Then, using the approximation method of volume rendering, the density and color on the ray that corresponds to the pixel of the image are integrated to calculate the RGB value. Using images captured from multiple viewpoints, the weights of the MLP are learned to minimize the L2 error between the observed and predicted RGB values.

OmniNeRF [21] generates multiple images at virtual camera positions from a single 360-degree RGB-D image, and utilizes these images to train NeRF. A set of 3D points is generated from the given RGB-D panorama and then these 3D points are reprojected into multiple omnidirectional images that correspond to different virtual camera locations. When reprojecting the 3D points onto the virtual camera spheres, their sparsity of the 3D points causes the back part of the object, which is not originally visible to see through. To address this challenge, a median filter is applied to the depth map to mitigate the sparsity.

4 Proposed Method

In this section, we describe our proposed method that learns the NeRF model and synthesizes novel views I from a single 360-degree image I_0 . Figure 2 illustrates the training pipeline of the proposed method. The input image is first reprojected onto a 360-degree image of the virtual camera position, and the missing regions are completed. The input data to train NeRF includes a sparse set of the dynamically selected completed images along with RGB data of the input image and its reprojected image. The NeRF is trained to minimize the L2 loss of the generated and training



Figure 2: Training pipeline of the proposed method. The input image is reprojected and completed as 360-degree images at other camera positions. A sparse set of the dynamically selected completed images is utilized to train the NeRF. The NeRF is trained to minimize the L2 loss with the ground truth image. Simultaneously, the selector reselects the completed images based on the discrimination score of synthesized images.

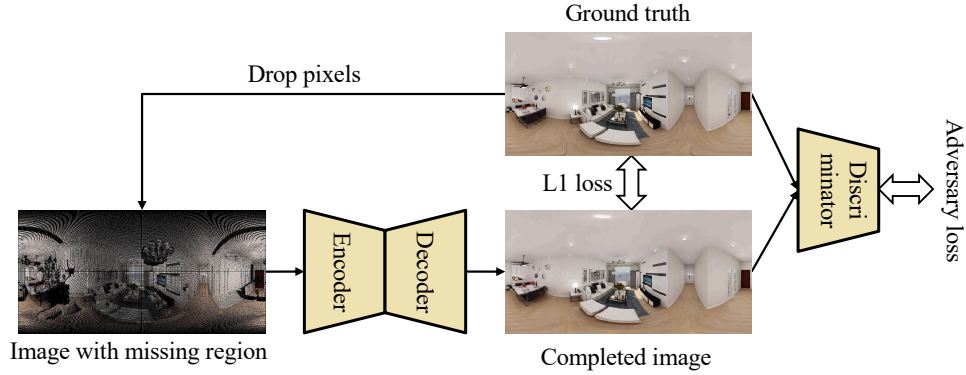


Figure 3: An image with randomly missing pixels from the ground truth image is utilized as input. Completed images are generated through the encoder and decoder, the L1 and adversary loss between the ground truth is obtained, and the weights of the network are learned to minimize the loss.

images, as usual. Simultaneously, the selector reselects the completed images based on the discrimination score of the synthesized images for real images. In the following sections, we describe the proposed method in detail.

4.1 Reprojection and completion

The input image is first reprojected onto a 360-degree image of the virtual camera position. This reprojection adopts the same approach as OmniNeRF, as described in Section 3. The reprojected image has missing regions due to occlusion and zooming, as illustrated in Fig. 1. We utilize a variation of the method in [18], which is an advanced 360 degree image completion. This method is based on the conditional VAE and GAN to estimate and control the latent symmetry in the image to complete the 360-degree image. The original method completes entire field of view from a normal-field-of-view image, and we modify it for a free-form completion. An overview of the image completion network, which comprises an encoder, decoder, and discriminator, is illustrated in Fig. 3. Each module is based on CNNs with self-attention [65]. Using an image with randomly missing pixels from the ground truth image as input, the weights of the network are learned in a self-supervised manner using L1 and adversarial losses. The details of network configurations are described in the Appendix A.

4.2 Selection of completion positions

Because the image completion is processed in the 360-degree field of view from a single point of view, the consistency between images observed from different positions cannot be guaranteed. Therefore, a sparse set of input images for training the NeRF is adaptively selected from the completed images. It is not easy to determine which combination of the completed images is 3D consistent. We can determine 3D consistency as a result of training NeRF, but it takes a long time to train NeRF once, and it is practically difficult to try all combinations. Therefore, we propose a method that simultaneously trains NeRF and selects completed images for training.

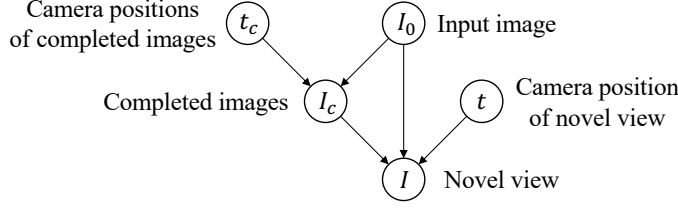


Figure 4: Graphical model of causal relationships for each variable in our novel view synthesis framework

4.2.1 Probabilistic framework

First, we consider a probabilistic framework for selecting the completed images. We denote the input image as I_0 , the selected set of M completion images as $I_c = \{I_c^{(i)}\}_{i=1}^M$, the camera positions for the selected completion images as $t_c = \{t_c^{(i)} \in \mathbb{R}^3\}_{i=1}^M$, generated novel views as I , and the camera position for the novel view as $t \in \mathbb{R}^3$. We assume the following joint distribution based on the causal relationship illustrated in Fig. 4:

$$p(I, t, I_0, I_c, t_c) = p(I|t, I_0, I_c)p(I_c|I_0, t_c)p(I_0)p(t)p(t_c). \quad (1)$$

The goal is to maximize the likelihood of the novel view with respect to the input image $p(I|I_0) = \mathbb{E}_{p(t)}[p(I|I_0, t)]$. Using Eq. (1) and Jensen's inequality, we can derive the variational lower bound of $\log p(I|I_0, t)$ as follows:

$$\begin{aligned} \log p(I|I_0, t) &= \log \iint p(I, I_c, t_c|I_0, t) dI_c dt_c \\ &= \log \iint p(I|t, I_0, I_c)p(I_c|I_0, t_c)p(t_c) dI_c dt_c \\ &\geq -\mathbb{KL}[q(t_c)||p(t_c)] \\ &\quad + \mathbb{E}_{q(t_c)}[\log p(I|t, I_0, \hat{I}_c(I_0, t_c)) + \log p(\hat{I}_c(I_0, t_c)|I_0, t_c)] \end{aligned} \quad (2)$$

where $q(t_c)$ is the variational posterior distribution of t_c and $\hat{I}_c(I_0, t_c)$ are completed images reprojected from I_0 at the camera positions t_c . A more detailed derivation of the equation is provided in the Appendix B. By determining $q(t_c)$ that maximizes this variational lower bound, it is possible to sample the camera positions of the completed images from $q(t_c)$.

4.2.2 Distribution setting

In the setting of this study, there is only one input image as the ground truth; hence, we cannot define $\log p(I|t, I_0, \hat{I}_c)$ as reconstruction errors with the ground truth, except for the input image position. Therefore, we represent $p(I|t, I_0, \hat{I}_c)$ by the output value $d \in [0, 1]$ of the discriminator in the image completion networks described in Section 4.1. $p(\hat{I}_c(I_0, t_c)|I_0, t_c)$ is also represented by the output value of the discriminator for $\hat{I}_c(I_0, t_c)$. Let $p(t_c)$ be the discrete uniform distribution of the predetermined completion positions. We assume that $q(t_c)$ can be decomposed into independent distributions for each camera position $t_c^{(i)}$, as $q(t_c) = \prod_{i=1}^M q(t_c^{(i)})$. Adjacency relations are defined between the discrete sampled completion positions in advance, and $q(t_c^{(i)})$ is defined as a distribution with the base position $\mu^{(i)} \in \mathbb{R}^3$ as a parameter, and the transition probability to the adjacent position $\epsilon \in [0, 1]$ as a hyperparameter. In the setting where L adjacent positions are defined, $q(t_c^{(i)}) = 1 - \epsilon L$ when $t_c^{(i)} = \mu^{(i)}$, and $q(t_c^{(i)}) = \epsilon$ when $t_c^{(i)}$ is the adjacent position. The details of the settings are described in Section 5.2.

4.2.3 Simultaneous training NeRF and selecting completion positions

We adopt the following heuristic algorithm to simultaneously train the NeRF and select the completed image positions.

Step 1 Initialize base positions $\mu = \{\mu^{(i)}\}_{i=1}^M$, and set μ to t_c .

Step 2 After K times of back propagation of NeRF, calculate the likelihood $p(I|I_0) = \mathbb{E}_{p(t)}[p(I|I_0, t)]$ using the approximate formula in Eq. (2). If the likelihood takes the maximum value in the past, then μ is replaced by t_c .

Step 3 Sample completed images positions t_c from $q(t_c)$.

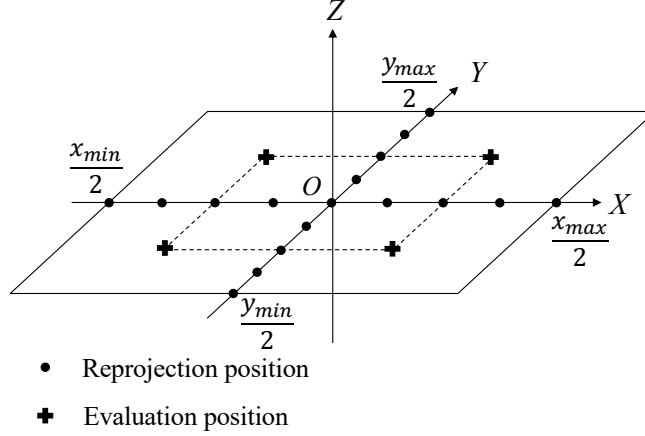


Figure 5: Settings of the reprojection and evaluation points. The origin O is the camera position of the input image, and the Z axis is parallel to the direction of gravity and orthogonal to the XY axis. One hundred re-projection points are equally spaced on $[\frac{x_{\min}}{2}, \frac{x_{\max}}{2}]$ on the x -axis and on the $[\frac{y_{\min}}{2}, \frac{y_{\max}}{2}]$ on the y -axis, where $x_{\min}, x_{\max}, y_{\min}, y_{\max}$ are the boundary point of the depth of the input image. The evaluation points for the likelihood are four points $(\frac{x_{\min}}{4}, \frac{y_{\min}}{4}), (\frac{x_{\min}}{4}, \frac{y_{\max}}{4}), (\frac{x_{\max}}{4}, \frac{y_{\min}}{4}), (\frac{x_{\max}}{4}, \frac{y_{\max}}{4})$.

Step 4 If it is repeated a specific number of times, the learning is completed; otherwise, it returns to step 2.

Note that in Step 2, Eq. (2) is computed approximately using t_c sampled at that time, and t is fixed to the evaluation positions as described in Section 5.2.

5 Experimental Results

Quantitative and qualitative experiments were conducted to verify the effectiveness of the proposed method for both synthetic and real-world datasets.

5.1 Dataset

5.1.1 Structured3D

Structured3D dataset [70] contains 3,500 synthetic departments (scenes) with 185,985 photorealistic panoramic renderings. As the original virtual environment is not publicly accessible, we utilized the rendered panoramas directly. The data were divided into 3,100 scenes for training, and 400 scenes for testing.

5.1.2 Matterport3D

Matterport3D dataset [10] is a large-scale indoor real-world 360 dataset, captured by Matterport’s Pro 3D camera in 90 furnished houses (scenes). The dataset provides a total of 10,800 RGB-D panorama images, where we find the RGB-D signals near the polar region are missing. The data were divided into 71 scenes for training, and 19 scenes for testing.

5.2 Implementation Details

We trained the image completion networks from scratch using the Adam optimizer [26] with a fixed learning rate of 1.0×10^{-4} and a mini-batch size of 8. For training NeRF, we also utilized the Adam optimizer while exponentially decreasing the learning rate from 5.0×10^{-4} to 5.0×10^{-5} . The NeRF model was trained with 200,000 iterations for each experiment with a batch size of 1,400. We set $N_c = 64$ and $N_f = 128$ for the coarse and refined networks. The NeRF settings were identical to those of OmniNeRF [21]. Learning one scene took approximately 16 hours on an NVIDIA Tesla V100 GPU. The image completion networks were trained on training data, while NeRF was trained on test data, according to the division defined in section 5.1. The input images, the completed images and novel views were all in equirectangular projection format with a resolution of $1,024 \times 512$.

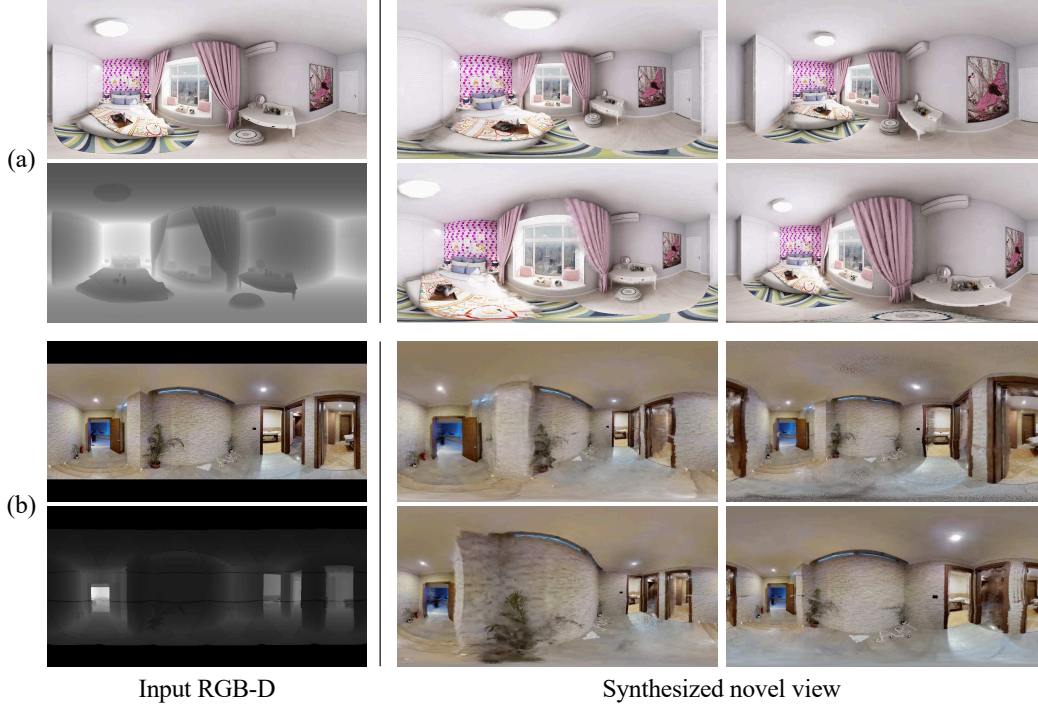


Figure 6: Visualization of novel view rendering in the proposed method at the evaluation positions defined in Fig. 5. (a) A sample in Structure3D dataset. (b) A sample in Matterport3D dataset. A plausible viewpoint image with 3D consistency is synthesized at a position different from the camera position of the input.

For the selection of completion positions, the transition probability is set as $\epsilon = 0.25$, and the update period of the base position is set as $K = 500$. $p(t_c)$ is a discrete uniform distribution that takes constant values at predetermined reprojection positions as illustrated in Fig. 5. The number of selected completion position is set as $M = 4$, and these positions are constrained on equally spaced 100 reprojection points on the X and Y axes. The reprojection points setting is identical to those of OmniNeRF [21]. The adjacent positions in $q(t_c)$ are adjacent reprojection points on each axis. The evaluation points for the likelihood are four points $(\frac{x_{\min}}{4}, \frac{y_{\min}}{4}), (\frac{x_{\min}}{4}, \frac{y_{\max}}{4}), (\frac{x_{\max}}{4}, \frac{y_{\min}}{4}), (\frac{x_{\max}}{4}, \frac{y_{\max}}{4})$, where $x_{\min}, x_{\max}, y_{\min}, y_{\max}$ are the boundary point of the depth of the input image.

5.3 Qualitative evaluation

First, we qualitatively validate the novel view synthesis using a single 360-degrees RGB-D image. Figure 6 illustrates examples of synthesized novel views by the proposed method. A plausible views with 3D consistency is synthesized at a position that is different from the camera position of the input. The depth data in Matterport3D dataset is subject to measurement error, which results in some distortion in the synthesized views. In Fig. 8 and Fig. 9, we compare images synthesized using OmniNeRF and the proposed method, including the case where all the completion images are utilized without selecting the completion positions. Although the occlusion regions are noisy in OmniNeRF, the proposed method produces a more plausible image. The proposed method, without selecting a completed image, generates blurry images. This is because there is no 3D consistency in the completed images at multiple camera positions; therefore, an average view is acquired that would explain each completed image. Additional results are available in the Appendix C.

5.4 Quantitative evaluation

We quantitatively evaluated each method using the following two evaluation metrics:

Peak-to-signal-noise ratio (PSNR) Because there is no ground truth image apart from the input image, we calculate the PSNR between the synthesized and the input images at the position of input image. This metric evaluates the performance of the input image reconstruction.

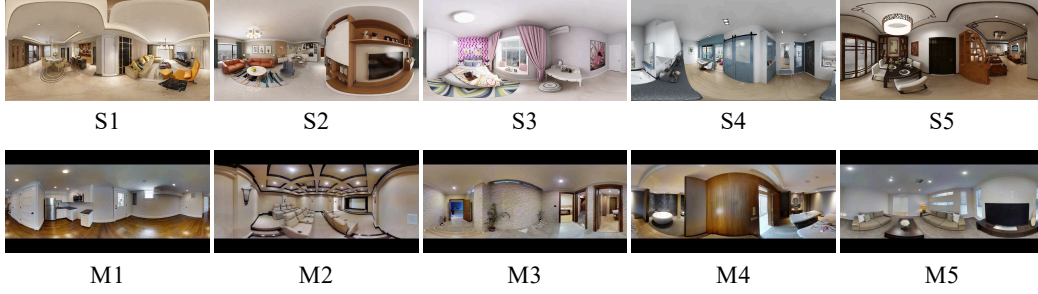


Figure 7: Input images used in quantitative evaluation. Images beginning with S are from the Structred3D dataset and images beginning with M are from the Matterport3D dataset.

Table 1: Quantitative evaluation of each novel view synthesis method on 10 images in Fig. 7. PSNR evaluates the performance of the input image reconstruction, and NLLF evaluates the plausibility of the synthesized views at different positions from the the input image. NLLF is described by dividing by 1,000.

	S1		S2		S3		S4		S5	
	PSNR↑	NLLF↓	PSNR↑	NLLF↓	PSNR↑	NLLF↓	PSNR↑	NLLF↓	PSNR↑	NLLF↓
OmniNeRF	29.26	4.270	30.59	4.445	32.85	3.876	29.59	3.298	26.72	4.864
Ours w/o selection	12.22	4.384	21.68	4.814	24.06	3.923	22.36	3.372	22.62	4.905
Ours	28.75	4.251	30.23	4.421	31.78	3.855	29.81	3.312	27.17	4.809
	M1		M2		M3		M4		M5	
	PSNR↑	NLLF↓	PSNR↑	NLLF↓	PSNR↑	NLLF↓	PSNR↑	NLLF↓	PSNR↑	NLLF↓
OmniNeRF	19.72	3.027	24.29	4.011	23.89	3.620	22.86	4.114	21.00	3.614
Ours w/o selection	17.64	3.240	21.84	4.454	21.10	3.916	19.06	4.333	20.71	3.755
Ours	19.83	3.010	24.07	4.094	23.71	3.515	22.93	4.092	21.67	3.586

Negative log likelihood of features (NLLF) We calculate the negative log likelihood of features in the synthesized image by modeling the feature distribution for the test data with a normal distribution, using the 2048 dimensional values of the last pooling layer of inception-v3 [54] as features. A total of 20,000 perspective images are randomly cropped from the test data, and the mean $\mu_f \in \mathbb{R}^{2048}$ and covariance matrix $\Sigma_f \in \mathbb{R}^{2048 \times 2048}$ of the 2048 dimensional features are calculated. The features $\{x_i \in \mathbb{R}^{2048}\}_{i=1}^N$ are calculated by randomly synthesizing $N = 512$ perspective projection images from the learned 3D scene, and the NLLF is calculated using the following formula:

$$\text{NLLF} = \frac{1}{N} \sum_{i=1}^N (x_i - \mu_f)^T \Sigma_f^{-1} (x_i - \mu_f) \quad (3)$$

This metric evaluates the plausibility of the synthesized views at different positions from the the input image.

We extract five images each from the test data of Structure3D and Matterport3D as shown in Fig. 7. The evaluation results for each image are presented in Table 1. In most images, the proposed method outperforms OmniNeRF in terms of the NLLF, which indicates that the proposed method is able to synthesize plausible views that have features close to the test dataset. This is thought to be because the image completion network trained on the training data generalized image features and was able to plausibly complete the missing regions in the scene. In contrast, the superiority in terms of PSNR depends on the input images, in some cases, the image completion reduces the performance of reconstruction by adding information that is different from the original image. However, note that the PSNR for the input images is employed as one of the characterization analyses of the methods and not to measure the performance of the novel view synthesis. Without the completion positions selection, both PSNR and NLLF are significantly degraded, indicating the importance of the selection module.

5.5 Limitations

Although the performance of the proposed method is promising, it has several limitations. First, if there are large missing regions that exceed the image completion capabilities, it is difficult to synthesize plausible views. Second, the reprojection process is highly dependent on the depth accuracy, and geometric distortion occurs in the synthesized image when the depth accuracy is low.

6 Conclusions

In this paper, we proposed a method for synthesizing novel views by learning the neural radiance field (NeRF) from a single 360-degree image. In the proposed method, the input image is reprojected to 360-degree images at other camera positions, the missing regions of the reprojected images were completed by a self-supervised trained generative model, and the completed images were utilized for training the NeRF. Because multiple completed images contain inconsistencies in 3D, we introduced a method to train NeRF while dynamically selecting a sparse set of the completed images to reduce the discrimination error of the synthesized views with real images. Experiments indicated that the proposed method can synthesize plausible novel views while preserving the features of the scene, both for artificial and real-world data. These results confirm the effectiveness of employing image completion and dynamic selection of the image completion positions for novel view synthesis.

Acknowledgements

This work was partially supported by JST AIP Acceleration Research JPMJCR20U3, Moonshot R&D Grant Number JPMJPS2011, JSPS KAKENHI Grant Number JP19H01115, and JP20H05556 and Basic Research Grant (Super AI) of Institute for AI and Beyond of the University of Tokyo.

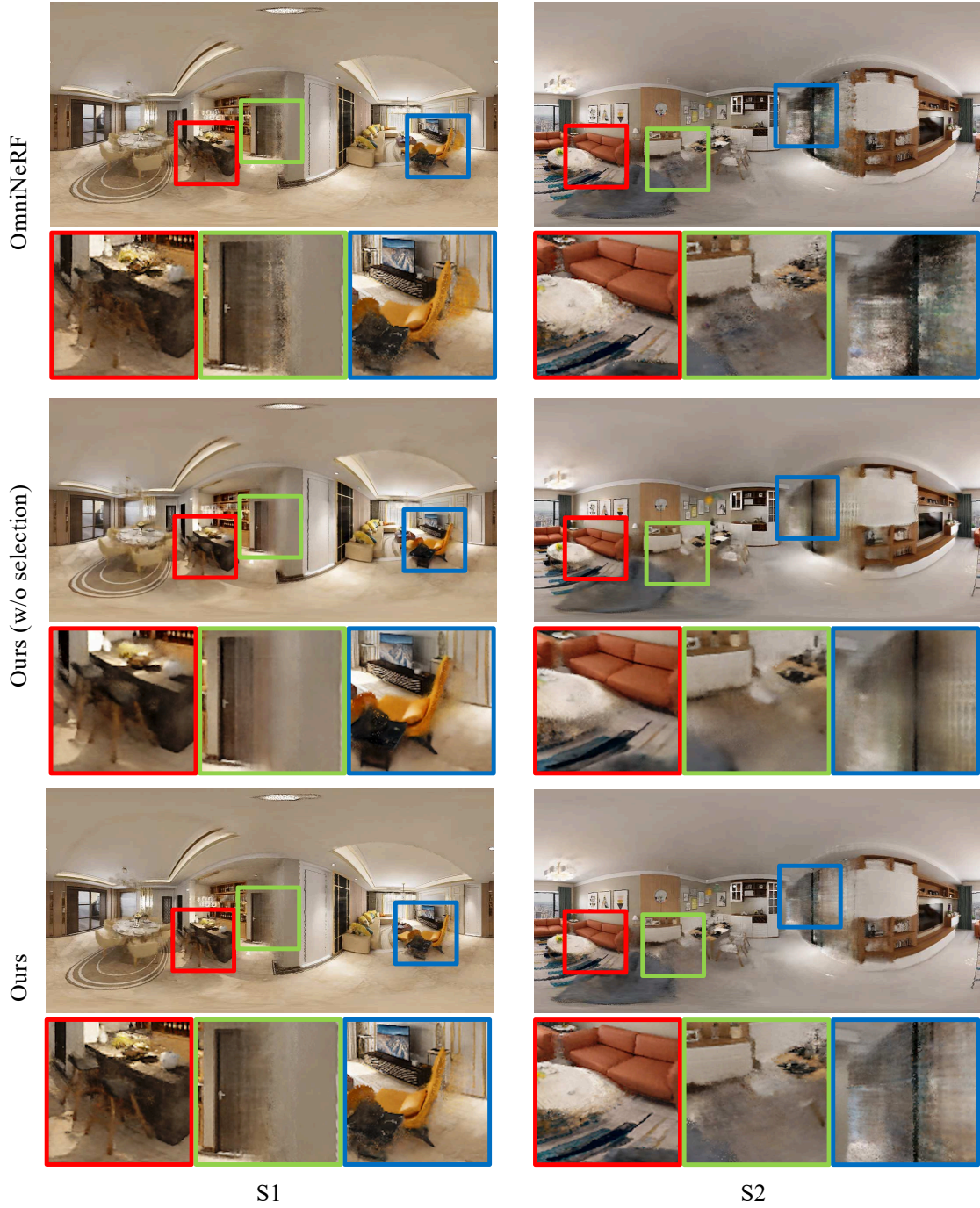


Figure 8: Qualitative comparison of OmniNeRF and the proposed method on Structure3D dataset. OmniNeRF produces artifacts in occlusion regions such as behind the yellow chair in scene S1 and on the wall in scene S2, while the proposed method reduces these artifacts. The backless chair in scene S1 has a collapsed image in OmniNeRF due to changes in the resolution and the viewing angle, while the proposed method keeps its shape. Without the completed images selection, synthesized views tend to be blurred.

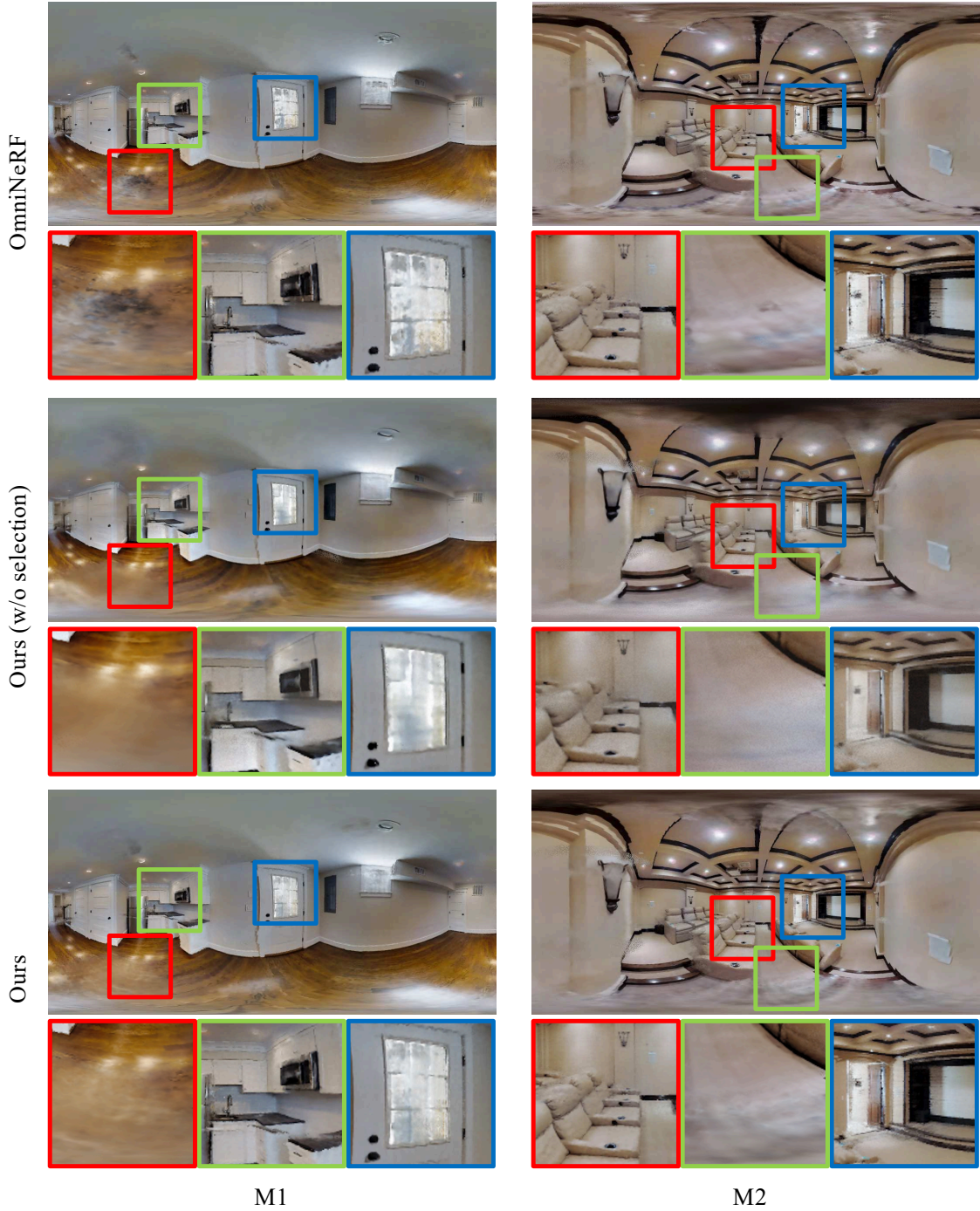


Figure 9: Qualitative comparison of each novel view synthesis method on Matterport3D dataset. The Matterport3D dataset contains missing regions at the top and bottom regions of the input image, and the proposed method completes the missing regions naturally, as in the floor of scene M1. The screen in scene M2 has some artifacts in OmniNeRF due to resolution change, but the proposed method reduces these artifacts. Without the completed images selection, synthesized views tend to be blurred.

References

- [1] Akimoto, N., Aoki, Y.: Image completion of 360-degree images by cgan with residual multi-scale dilated convolution. *IIEEJ Trans. Image Electronics and Vis. Comput.* **8**(1), 35–43 (2020)
- [2] Akimoto, N., Kasai, S., Hayashi, M., Aoki, Y.: 360-degree image completion by two-stage conditional gans. In: *Proc. IEEE Int. Conf. Image Processing*. pp. 4704–4708 (2019)
- [3] Attal, B., Ling, S., Gokaslan, A., Richardt, C., Tompkin, J.: Matryodshka: Real-time 6dof video view synthesis using multi-sphere images. In: *Europ. Conf. Comput. Vis.* (2020)
- [4] Badki, A., Gallo, O., Kautz, J., Sen, P.: Meshlet priors for 3d mesh reconstruction. In: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* p. 2846–2855 (2020)
- [5] Ballester, C., Bertalmio, M., Caselles, V., Sapiro, G., Verdera, J.: Filling-in by joint interpolation of vector fields and gray levels. *IEEE Trans. Image Processing* **10**(8), 1200–1211 (2001)
- [6] Barnes, C., Shechtman, E., Finkelstein, A., Goldman, D.B.: Patchmatch: A randomized correspondence algorithm for structural image editing. *ACM Trans. Graph.* **28**(3) (2009)
- [7] Barron, J.T., Mildenhall, B., Verbin, D., Srinivasan, P.P., Hedman, P.: Mip-nerf 360: Unbounded anti-aliased neural radiance fields. *arXiv:2111.12077* (2021)
- [8] Bertalmio, M., Sapiro, G., Caselles, V., Ballester, C.: Image inpainting. In: *Proc. Conf. Comput. Graph. Interact. Techniq.* pp. 417–424 (2000)
- [9] Chan, E.R., Lin, C.Z., Chan, M.A., Nagano, K., Pan, B., Mello, S.D., Gallo, O., Guibas, L., Tremblay, J., Khamis, S., Karras, T., Wetzstein, G.: Efficient geometry-aware 3d generative adversarial networks. *arXiv:2112.07945* (2020)
- [10] Chang, A., Dai, A., Funkhouser, T., Halber, M., Nießner, M., Savva, M., Song, S., Zeng, A., Zhang, Y.: Matterport3d: Learning from rgb-d data in indoor environments. In: *Int. Conf. 3D Vis.* (2017)
- [11] Criminisi, A., Perez, P., Toyama, K.: Object removal by exemplar-based inpainting. In: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* pp. II–II (2003)
- [12] DeVries, T., Bautista, M.A., Srivastava, N., Taylor, G.W., Susskind, J.M.: Unconstrained scene generation with locally conditioned radiance fields. In: *Proc. IEEE Int. Conf. Comput. Vis.* (2021)
- [13] Fukushima, K., Miyake, S.: Neocognitron: A new algorithm for pattern recognition tolerant of deformations and shifts in position. *Pattern Recognition* **15**(6), 455–469 (1982)
- [14] Garbin, S.J., Kowalski, M., Johnson, M., Shotton, J., Valentin, J.: Fastnerf: High-fidelity neural rendering at 200fps. In: *International Conference on Computer Vision* (2021)
- [15] Gardner, M.A., Sunkavalli, K., Yumer, E., Shen, X., Gambaretto, E., Christian, G., Lalonde, J.F.: Learning to predict indoor illumination from a single image. *ACM Trans. Graph.* **9**(4) (2017)
- [16] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: *Proc. Int. Conf. Neural Inf. Process. Syst.* pp. 2672–2680. (2014)
- [17] Han, S.W., Suh, D.Y.: A 360-degree panoramic image inpainting network using a cube map. *CMC-Computers, Materials and Continua* **66**(1), 213–228 (2021)
- [18] Hara, T., Mukuta, Y., Harada, T.: Spherical image generation from a single image by considering scene symmetry. In: *Proc. AAAI Conf. Artif. Intell.* pp. 1513–1521 (2021)
- [19] Hartley, R., Zisserman, A.: *Multiple view geometry in computer vision*. Cambridge university press (2003)
- [20] Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. In: *Proc. Int. Conf. Neural Inf. Process. Syst.* (2020)
- [21] Hsu, C.Y., Sun, C., Chen, H.T.: Moving in a 360 world: Synthesizing panoramic parallaxes from a single panorama. *arXiv:2106.10859* (2021)
- [22] Iizuka, S., Simo-Serra, E., Ishikawa, H.: Globally and locally consistent image completion. *ACM Trans. Graph.* **36**(4) (2017)
- [23] Jain, A., Tancik, M., Abbeel, P.: Putting nerf on a diet: Semantically consistent few-shot view synthesis. In: *Proc. IEEE Int. Conf. Comput. Vis.* (2021)
- [24] Jeong, Y., Ahn, S., Choy, C., Anandkumar, A., Cho, M., Park, J.: Self-calibrating neural radiance fields. In: *International Conference on Computer Vision* (2021)

- [25] Kato, H., Ushiku, Y., Harada, T.: Neural 3d mesh renderer. In: Proc. IEEE Conf. Comput. Vis. Pattern Recognit. p. 3907–3916 (2018)
- [26] Kingma, D.P., Ba, J.L.: Adam: A method for stochastic optimization. arXiv:1412.6980 (2014)
- [27] Kingma, D.P., Welling, M.: Auto-encoding variational bayes. arXiv:1312.6114 (2013)
- [28] Koh, J.Y., Lee, H., Yang, Y., Baldridge, J., Anderson, P.: Pathdreamer: A world model for indoor navigation. In: Proc. IEEE Int. Conf. Comput. Vis. (2021)
- [29] LeCun, Y., Boser, B., Denker, J.S., Henderson, D., Howard, R.E., Hubbard, W., Jackel, L.D.: Backpropagation applied to handwritten zip code recognition. *Neural computation* **1**(4), 541–551 (1989)
- [30] Levoy, M., Hanrahan, P.: Light field rendering. In: SIGGRAPH (2020)
- [31] Li, Y., Liu, S., Yang, J., Yang, M.H.: Generative face completion. In: Proc. IEEE Conf. Comput. Vis. Pattern Recognit. pp. 3911–3919 (2017)
- [32] Lin, C.H., Ma, W.C., Torralba, A., Lucey, S.: Barf : Bundle-adjusting neural radiance fields. In: International Conference on Computer Vision (2021)
- [33] Liu, G., Reda, F.A., Shih, K.J., Wang, T.C., Tao, A., Catanzaro, B.: Image inpainting for irregular holes using partial convolutions. In: Proc. Eur. Conf. Comput. Vis. pp. 85–100 (2018)
- [34] Liu, H., Wan, Z., Huang, W., Song, Y., Han, X., Liao, J.: Pd-gan: Probabilistic diverse gan for image inpainting. In: Proc. IEEE Conf. Comput. Vis. Pattern Recognit. pp. 9371–9381 (2021)
- [35] Lugmayr, A., Danelljan, M., Romero, A., Yu, F., Timofte, R., Gool, L.V.: Repaint: Inpainting using denoising diffusion probabilistic models. arXiv:2201.09865 (2022)
- [36] Martin-Brualla, R., Radwan, N., Sajjadi, M.S.M., Barron, J.T., Dosovitskiy, A., Duckworth, D.: Nerf in the wild: Neural radiance fields for unconstrained photo collections. In: Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (2021)
- [37] Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: Europ. Conf. Comput. Vis. (2020)
- [38] Müller, T., Evans, A., Schied, C., Keller, A.: Instant neural graphics primitives with a multiresolution hash encoding. arXiv:2201.05989 (Jan 2022)
- [39] Park, T., Liu, M.Y., Wang, T.C., Zhu, J.Y.: Semantic image synthesis with spatially-adaptive normalization. In: Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (2019)
- [40] Peng, J., Liu, D., Xu, S., Li, H.: Generating diverse structure for image inpainting with hierarchical vq-vae. In: Proc. IEEE Conf. Comput. Vis. Pattern Recognit. pp. 10775–10784 (2021)
- [41] Reiser, C., Peng, S., Liao, Y., Geiger, A.: Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In: Proc. IEEE Int. Conf. Comput. Vis. (Jan 2021)
- [42] Rematas, K., Liu, A., Srinivasan, P.P., Barron, J.T., Tagliasacchi, A., Funkhouser, T., Ferrari, V.: Urban radiance fields. arXiv:2111.14643 (2021)
- [43] Rezende, D.J., Mohamed, S., Wierstra, D.: Stochastic backpropagation and approximate inference in deep generative models. In: Proc. Int. Conf. Mach. Learn. pp. 1278–1286 (2014)
- [44] Shih, M.L., Su, S.Y., Kopf, J., Huang, J.B.: 3d photography using context-aware layered depth inpainting. In: Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (2020)
- [45] Shum, H.Y., Chan, S.C., Kang, S.B.: Image-based rendering. Springer Science and Business Media (2008)
- [46] Sohl-Dickstein, J., Weiss, E.A., Maheswaranathan, N., Ganguli, S.: Deep unsupervised learning using nonequilibrium thermodynamics. In: Proc. Int. Conf. Mach. Learn. (2015)
- [47] Song, S., Funkhouser, T.: Neural illumination: Lighting prediction for indoor environments. In: Proc. IEEE Conf. Comput. Vis. Pattern Recognit. pp. 6918–6926 (2019)
- [48] Song, S., Zeng, A., Chang, A.X., Savva, M., Savarese, S., Funkhouser, T.: Im2pano3d: Extrapolating 360° structure and semantics beyond the field of view. In: Proc. IEEE Conf. Comput. Vis. Pattern Recognit. pp. 3847–3856 (2018)
- [49] Song, W., Shi, C., Xiao, Z., Duan, Z., Xu, Y., Zhang, M., Tang, J.: Autoint: Automatic feature interaction learning via self-attentive neural networks. In: Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (2021)
- [50] Srinivasan, P.P., Mildenhall, B., Tancik, M., Barron, J.T., Tucker, R., Snavely, N.: Lighthouse: Predicting lighting volumes for spatially-coherent illumination. In: Proc. IEEE Conf. Comput. Vis. Pattern Recognit. pp. 8080–8089 (2020)

- [51] Sumantri, J.S., Park, I.K.: 360 panorama synthesis from a sparse set of images on a low-power device. *IEEE Trans. on Comput. Imaging* **6**, 1179–1193 (2020)
- [52] Sumantri, J.S., Park, I.K.: 360 panorama synthesis from a sparse set of images with unknown fov. In: *IEEE Winter Conf. Applications of Comput. Vis.* pp. 2386–2395 (2020)
- [53] Suvorov, R., Logacheva, E., Mashikhin, A., Remizova, A., Ashukha, A., Silvestrov, A., Kong, N., Goka, H., Park, K., Lempitsky, V.: Resolution-robust large mask inpainting with fourier convolutions. In: *IEEE Winter Conf. Applications of Comput. Vis.* (2022)
- [54] Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* (2016)
- [55] Tancik, M., Casser, V., Yan, X., Pradhan, S., Mildenhall, B., Srinivasan, P.P., Barron, J.T., Kretzschmar, H.: Block-nerf: Scalable large scene neural view synthesis. *arXiv:2202.05263* (2022)
- [56] Tewari, A., Christian, T., Goldman, D.B., Shechtman, E., Wetzstein, G., Saragih, J., Zhu, J.Y., Thies, J., Sunkavalli, K., Agrawala, M., Niessner, M., Zollhofer, M., Fried, O., Brualla, R.M., Pandey, R.K., Fanello, S., Lombardi, S., Simon, T., Sitzmann, V.: State of the art on neural rendering. *Computer Graphics Forum* (2020)
- [57] Trevithick, A., Yang, B.: Grf: Learning a general radiance field for 3d scene representation and rendering. In: *International Conference on Computer Vision* (2021)
- [58] Wang, Q., Wang, Z., Genova, K., Srinivasan, P., Zhou, H., Barron, J.T., Martin-Brualla, R., Snavely, N., Funkhouser, T.: Ibrnet: Learning multi-view image-based rendering. In: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* (2021)
- [59] Wang, Z., Wu, S., Xie, W., Chen, M., Prisacariu, V.A.: Nerf-: Neural radiance fields without known camera parameters. *arXiv:2102.07064* (2021)
- [60] Yan, Z., Li, X., Li, M., Zuo, W., Shan, S.: Shift-net: Image inpainting via deep feature rearrangement. In: *Proc. Eur. Conf. Comput. Vis.* pp. 1–17 (2018)
- [61] Yu, A., Li, R., Tancik, M., Li, H., Ng, R., Kanazawa, A.: Plenotrees for real-time rendering of neural radiance fields. In: *International Conference on Computer Vision* (2021)
- [62] Yu, A., Ye, V., Tancik, M., Kanazawa, A.: Pixelnerf: Neural radiance fields from one or few images. In: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* (2021)
- [63] Yu, J., Lin, Z., Yang, J., Shen, X., Lu, X., Huang, T.S.: Generative image inpainting with contextual attention. in *proceedings of the ieee conference on computer vision and pattern recognition*. In: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* pp. 5505–5514 (2018)
- [64] Yu, J., Lin, Z., Yang, J., Shen, X., Lu, X., Huang, T.S.: Free-form image inpainting with gated convolution. In: *Proc. IEEE Int. Conf. Comput. Vis.* pp. 4471–4480 (2019)
- [65] Zhang, H., Goodfellow, I., Metaxas, D., Odena, A.: Self-attention generative adversarial networks. *arXiv:1805.08318* (2018)
- [66] Zhang, K., Riegler, G., Snavely, N., Koltun, V.: Nerf++: Analyzing and improving neural radiance fields. *arXiv:2010.07492* (2020)
- [67] Zhao, L., Mo, Q., Lin, S., Wang, Z., Zuo, Z., Chen, H., Xing, W., Lu, D.: Uctgan: Diverse image inpainting based on unsupervised cross-space translation. In: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* pp. 5741–5750 (2020)
- [68] Zhao, S., Cui, J., Sheng, Y., Dong, Y., Liang, X., Chang, E.I., Xu, Y.: Large scale image completion via co-modulated generative adversarial networks. In: *Proc. Int. Conf. Learn. Representations* (2021)
- [69] Zheng, C., Cham, T.J., Cai, J.: Pluralistic image completion. In: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* pp. 1438–1447 (2020)
- [70] Zheng, J., Zhang, J., Li, J., Tang, R., Gao, S., Zhou, Z.: Structured3d: A large photo-realistic dataset for structured 3d modeling. In: *Europ. Conf. Comput. Vis.* (2020)
- [71] Zhou, T., Tucker, R., Flynn, J., Fyffe, G., Snavely, N.: Stereo magnification: Learning view synthesis using multiplane images. In: *ACM Trans. Graph* (2018)

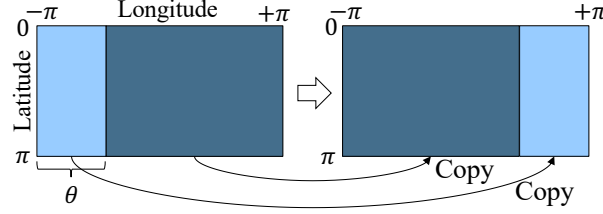


Figure 10: Circular shift [18]. A horizontal circular shift in the equirectangular projection corresponds to a rotation around the gravity axis.

A Details of Image Completion

In this section, we describe the network configuration and training method in our image completion network.

A.1 Network configuration

We utilize a variation of the spherical image generation network (SIGN) [18] as the image completion network. SIGN generates 360-degree images by combining symmetrical patterns of features in the input normal-field-of-view (NFOV) image. However, this method cannot be employed in its original form, because the focus of this study is on the task of free-form completion. A rotation around the gravity axis corresponds to a horizontal cyclic shift in the equirectangular projection, as illustrated in Fig. 10. We assume that the combination of rotational symmetric patterns in SIGN improves completion performance by expanding the receptive field of the CNN. Based on this assumption, we design a network that cyclic shifts some channels of the features. Specifically, cyclic shifts corresponding to rotations of 90, 180, and 270 degrees are applied to 37.5% of the total channels.

The image completion network comprises an encoder F , decoder G , and discriminator D , and each module have the same combination of ResBlocks [69] with circular padding (RBCP) [18], as illustrated in Fig. 11. The RBCP has four modes as well as ResBlock: (s) start, (n) normal, (d) down-sampling, and (u) up-sampling, and the behaviors of each of the four modes are different. The structure of each function is as follows:

- F : one-layer RBCP(s) and four-layer RBCP(d).
- G : two-layer RBCP(u), self-attention layer, and three-layer RBCP(u).
- D : one-layer RBCP(s), two-layer RBCP(d), self-attention layer, two-layer RBCP(d), one-layer RBCP(n), and four-layer RBCP(d). The last four-layer RBCP(d) corresponds to the global discriminator.

The encoder F outputs f_e and f_l which are the outputs of the third and final (fifth) layers, respectively. f_l is input to the first layer of the decoder G and f_e is concatenated to the third layer of G in the channel.

A.2 Training Method

We trained the network with common weights progressively from low resolution to high resolution. First, we trained the network for up to 200,000 iterations with a mini-batch size of 8 in a resolution of 512×256 , then 90,000 iterations with a mini-batch size of 4 in a resolution of 1024×512 .

B Probabilistic Framework

We describe the derivation of Eq. (2). We denote the input image as I_0 , the selected set of M completion images as $I_c = \{I_c^{(i)}\}_{i=1}^M$, the camera positions for the selected completion images as $t_c = \{t_c^{(i)} \in \mathbb{R}^3\}_{i=1}^M$, generated novel views as I , and the camera position for the novel view as $t \in \mathbb{R}^3$. We assume the following joint distribution based on the causal relationship illustrated in Fig. 4.

$$p(I, t, I_0, I_c, t_c) = p(I|t, I_0, I_c)p(I_c|I_0, t_c)p(I_0)p(t)p(t_c). \quad (4)$$

Therefore, as the conditional probability for t, I_0 , we obtain the following:

$$p(I, I_c, t_c|t, I_0) = \frac{p(I, t, I_0, I_c, t_c)}{\iiint p(I, t, I_0, I_c, t_c) dI dI_c dt_c}$$

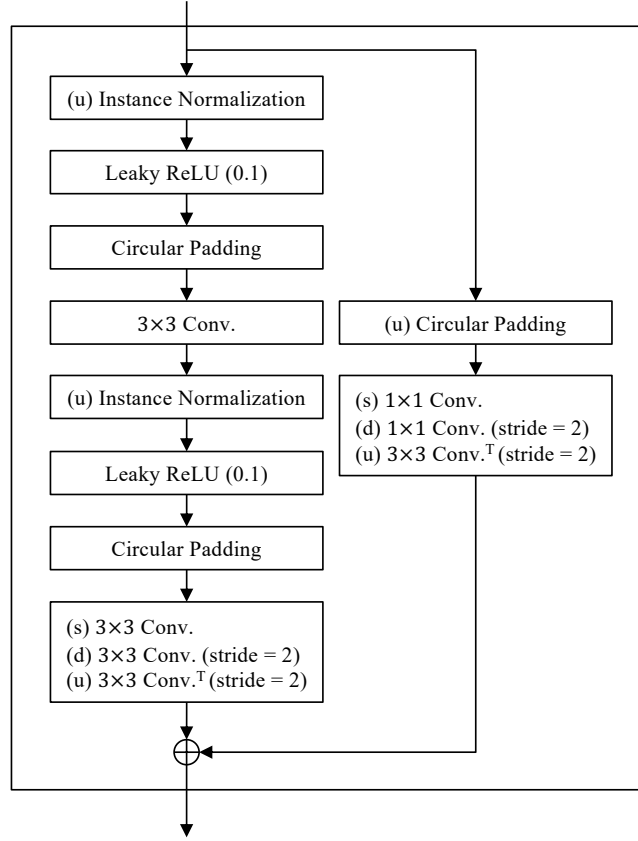


Figure 11: Residual block with circular padding [18]. There are four modes: (s) start, (n) normal, (d) down-sampling, and (u) up-sampling. For example, during the start mode, the modules marked with (s) and unmarked modules are used.

$$\begin{aligned}
 &= \frac{p(I, t, I_0, I_c, t_c)}{p(t)p(I_0)} \\
 &= p(I|t, I_0, I_c)p(I_c|I_0, t_c)p(t_c)
 \end{aligned} \tag{5}$$

Next, we introduce the variational lower bound of $\log p(I|t, I_0)$ using Jensen's inequality as follows:

$$\log p(I|t, I_0) = \log \iint q(I_c, t_c) \frac{p(I, I_c, t_c|t, I_0)}{q(I_c, t_c)} dI_c dt_c \tag{6}$$

$$\begin{aligned}
 &\geq \iint q(I_c, t_c) \log \frac{p(I, I_c, t_c|t, I_0)}{q(I_c, t_c)} dI_c dt_c, \\
 &:= \mathcal{B}
 \end{aligned} \tag{7}$$

where $q(I_c, t_c)$ is the variational posterior distribution of I_c, t_c and can be decomposed as $q(I_c, t_c) = q(I_c|t_c)q(t_c)$. Since I_c is uniquely determined to be \hat{I}_c with respect to t_c by the image completion, we set $q(I_c|t_c) = \delta(I_c - \hat{I}_c)$, where δ is the Dirac delta function. Therefore, the lower bound \mathcal{B} can be expressed as:

$$\begin{aligned}
 \mathcal{B} &= \iint q(I_c|t_c)q(t_c) \log \frac{p(I, I_c, t_c|t, I_0)}{q(I_c|t_c)q(t_c)} dI_c dt_c \\
 &= \int q(t_c) \log \frac{p(I, \hat{I}_c, t_c|t, I_0)}{q(t_c)} dt_c.
 \end{aligned} \tag{8}$$

Using equation Eq. (5), we finally obtain the following equation:

$$\begin{aligned}
 \mathcal{B} &= \int q(t_c) \log \frac{p(I|t, I_0, \hat{I}_c)p(\hat{I}_c|I_0, t_c)p(t_c)}{q(t_c)} dt_c, \\
 &= -\mathbb{KL}[q(t_c)||p(t_c)] + \mathbb{E}_{q(t_c)}[\log p(I|t, I_0, \hat{I}_c) + \log p(\hat{I}_c|I_0, t_c)]
 \end{aligned} \tag{9}$$

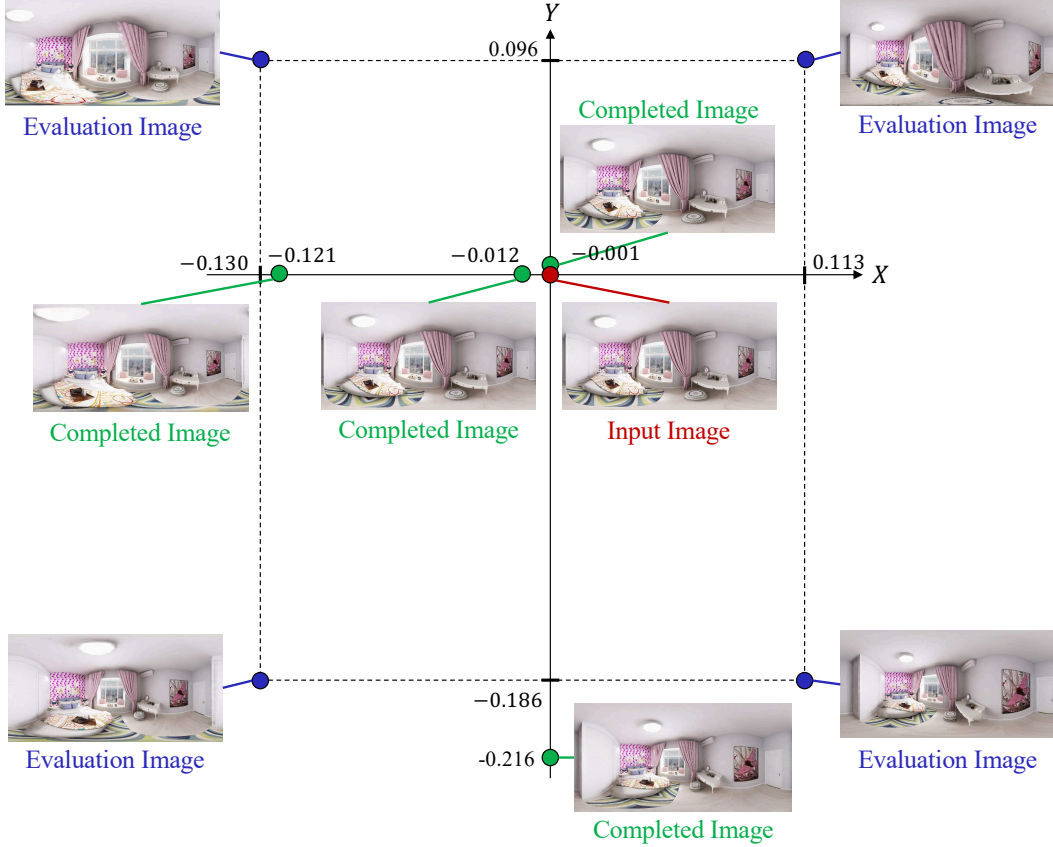


Figure 12: An example of the content and location of the input image, the completed images selected to train NeRF, and the evaluation images. The X -axis and Y -axis lie on a plane that is orthogonal to the gravity axis Z , and the maximum depth in the input image is scaled to 1.0.

C Additional Results

C.1 Selected completion images

Figure 12 shows an example of the content and location of the input image, the completed images selected to train NeRF, and the evaluation images. The completed images correspond to the base positions μ of $q(t_c)$. The evaluation images correspond to the synthesized image at evaluation points described in Section 5.2. Note that, although not indicated in this figure, 100 non-completed reprojection images are also utilized to train NeRF, as in OmniNeRF [21]. In this example, the completed images are selected to have high completion accuracy around the input image and to complete large missing regions away from the input position.

C.2 Synthesized novel views

The results of novel view synthesis using OmniNeRF [21] and the proposed method (with and without the completed images selection) in scenes S3, S4, S5, M3, M4 and M5 are illustrated in Fig. 13, Fig. 14 and Fig. 15. These figures show synthesized novel views in equirectangular and perspective projections with different camera positions.

C.3 Rendering from a free moving camera

Examples of rendering from a freely moving camera are shown in Fig. 16. In this figure, novel views are rendered in perspective projection with a horizontal field of view of 90 degrees.

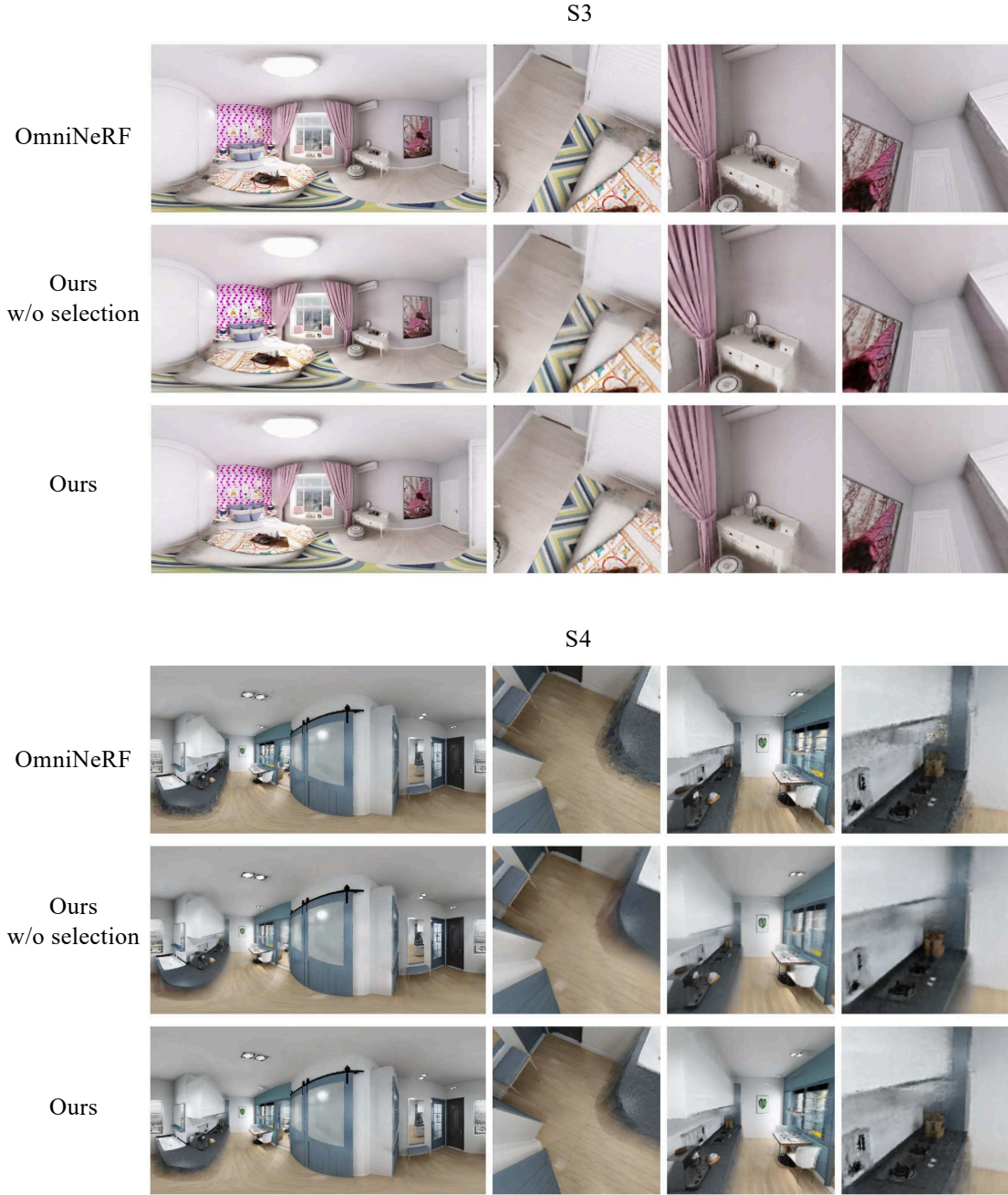


Figure 13: Novel views synthesized in equirectangular and perspective projections with different camera positions for scenes S3 and S4.

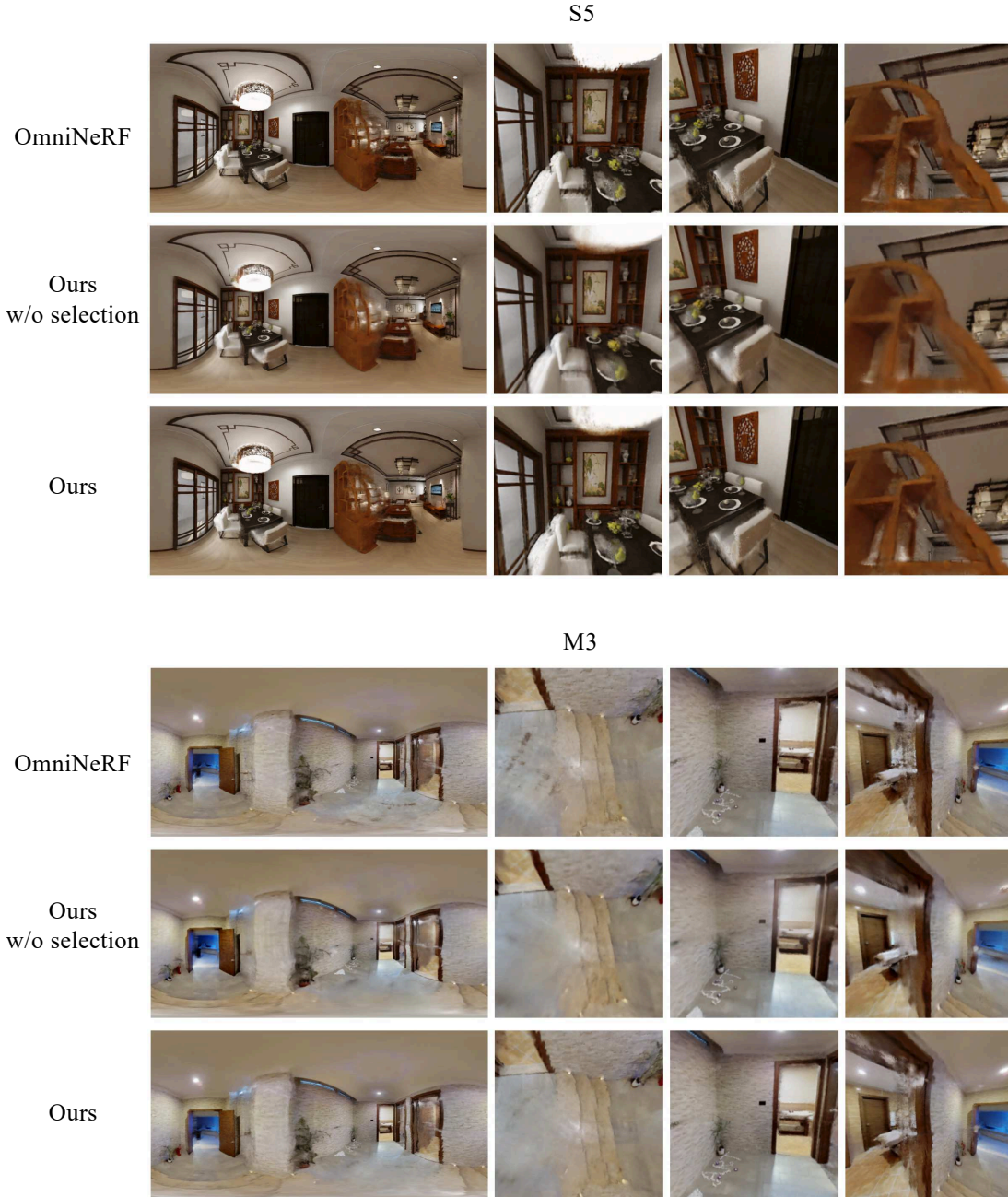


Figure 14: Novel views synthesized in equirectangular and perspective projections with different camera positions for scenes S5 and M3.

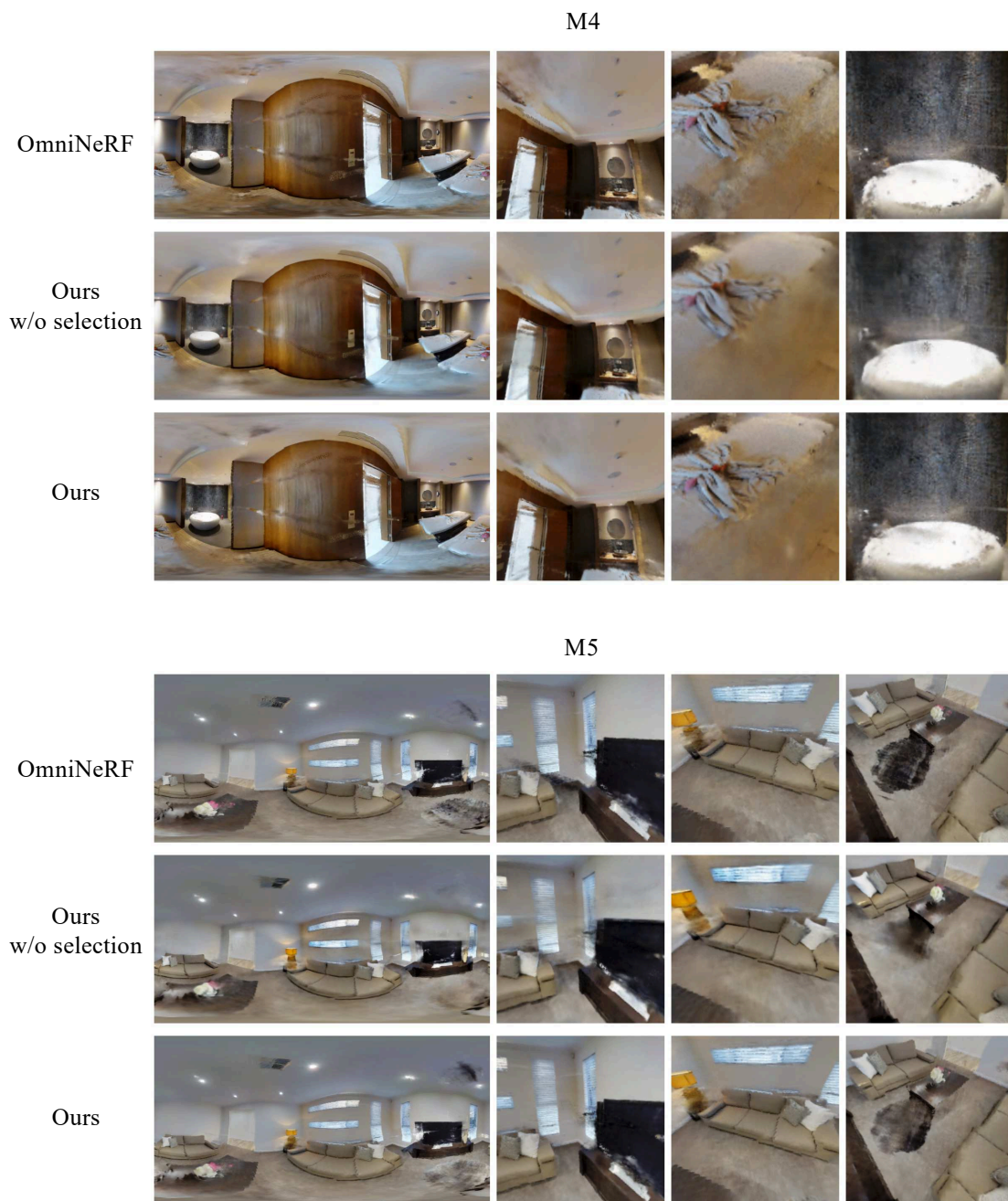


Figure 15: Novel views synthesized in equirectangular and perspective projections with different camera positions for scenes M4 and M5.

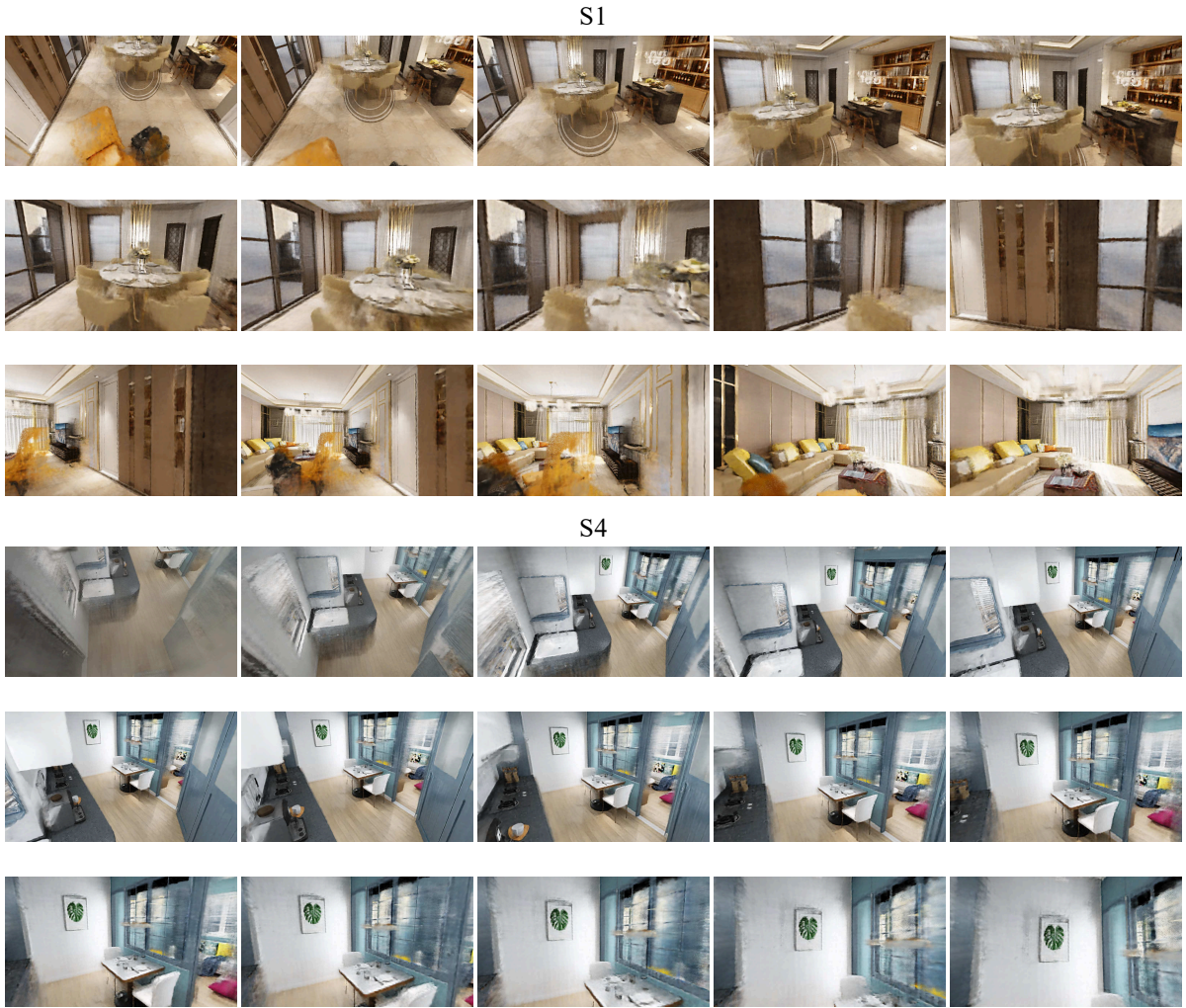


Figure 16: Examples of rendering from a freely moving camera. The novel views are rendered in perspective projection with a horizontal field of view of 90 degrees.