

NeRF-Enhanced Outpainting for Faithful Field-of-View Extrapolation

Rui Yu^{1†}, Jiachen Liu^{2†}, Zihan Zhou³, Sharon X. Huang²

Abstract—In various applications, such as robotic navigation and remote visual assistance, expanding the field of view (FOV) of the camera proves beneficial for enhancing environmental perception. Unlike image outpainting techniques aimed solely at generating aesthetically pleasing visuals, these applications demand an extended view that faithfully represents the scene. To achieve this, we formulate a new problem of faithful FOV extrapolation that utilizes a set of pre-captured images as prior knowledge of the scene. To address this problem, we present a simple yet effective solution called NeRF-Enhanced Outpainting (NEO) that uses extended-FOV images generated through NeRF to train a scene-specific image outpainting model. To assess the performance of NEO, we conduct comprehensive evaluations on three photorealistic datasets and one real-world dataset. Extensive experiments on the benchmark datasets showcase the robustness and potential of our method in addressing this challenge. We believe our work lays a strong foundation for future exploration within the research community.

I. INTRODUCTION

The field of view (FOV) of a camera plays a pivotal role in the performance of vision-based navigation [1], [2]. A larger FOV enables robots to perceive more spatial elements and layouts (e.g., obstacles, doorways, etc.). This expanded perspective empowers them to make more informed and strategic decisions when planning their paths. A larger FOV also offers substantial benefits for remote sighted agents (RSAs) tasked with assisting visually impaired individuals in navigation [3], [4]. In light of this motivation, our work delves into the challenge of FOV extrapolation. Our goal is to enable robots and remote agents to perceive scene content beyond the immediate camera FOV, thereby enhancing their situational awareness.

In the computer vision domain, our task is closely related to image outpainting [5], [6], [7], [8], also referred to as image extrapolation [9], [10], [11] or image extension [12], which aims to extend the image boundaries with semantically consistent and visually appealing contents. The extrapolated image hallucinates visually plausible scene contents beyond the original FOV, delivering immersive viewing experience for applications such as virtual reality. The hallucination ability is acquired by, for example, training deep learning models on large-scale generic datasets [13] of realistic images. However, image outpainting models cannot be applied

to our problem because navigational applications necessitate the extended portions of the image maintain fidelity and coherence with the actual scene.

We define the problem of *faithful FOV extrapolation* as follows. As illustrated in Fig. 1 (right), a collection of training images that were captured in a given scene (e.g., images shown inside red boxes) serves as prior knowledge. We assume that the camera pose corresponding to each training image can be obtained, for example, via structure from motion (SfM) [14], [15], [16] methods. During the testing phase, our objective is to faithfully extrapolate the FOV of any newly captured image in the same scene to a specified FOV based on the prior knowledge of the scene (see example images inside blue boxes).

It is possible to adapt existing computer vision techniques, namely image stitching and video expansion, to tackle the faithful FOV extrapolation problem. Image stitching [17], [18], [19], [20] entails aligning overlapping portions of multiple images taken from various angles, whereas video expansion [21], or video extrapolation [22], [23], leverages adjacent frames to extend the FOV of a specific frame. Both methods require precise warping of source images to blend seamlessly with the target image. However, they often yield irregularly shaped non-overlapping areas, imposing limitations on the extent of FOV expansion. Therefore, these methods are not well suited for our goal, as the faithful FOV extrapolation task demands expanding the view to a desired rectangular size.

Given that there has been very limited prior research addressing the challenge of faithful FOV extrapolation for navigation, we propose a simple yet effective method called *NeRF-Enhanced Outpainting (NEO)*. Our method involves first training a neural radiance fields [24] (NeRF) model using training images of original FOV. We then densely sample a substantial number of camera poses within the same scene and, for each sampled pose, the trained NeRF model is applied to render an image with an expanded FOV. Finally, we leverage these rendered images to train an image outpainting model, which is subsequently employed to extrapolate the FOV of input images during the inference phase. We validate the proposed method on three photorealistic datasets and one real dataset. The NEO method excels at producing high-quality extrapolations tailored to the specified FOV and consistently surpasses the performance of three baseline methods.

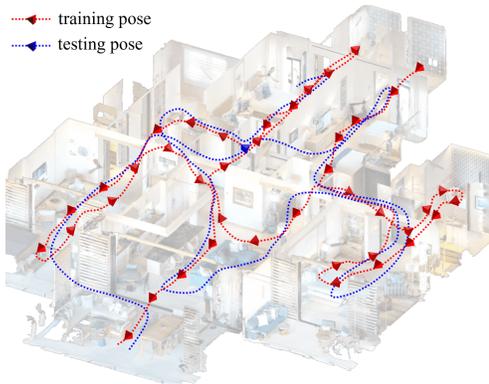
In summary, the contributions of this work are as follows. (1) We introduce a novel problem, namely faithful FOV extrapolation for navigation, which has been relatively underexplored in existing literature. (2) We propose the

[†]Equal contribution to this work.

¹Rui Yu is with the Department of Computer Science and Engineering, University of Louisville, Louisville, KY 40208, USA rui.yu@louisville.edu

²Jiachen Liu and Sharon X. Huang are with College of Information Sciences and Technology, Pennsylvania State University, University Park, PA 16802, USA {jz16493, suh972}@psu.edu

³Zihan Zhou is with Manycore Tech Inc., Hangzhou, Zhejiang, China shuer@qunhemail.com



- Given: Images captured in a scene as prior knowledge



- Goal: Faithfully extrapolate FOV of a new image in the same scene

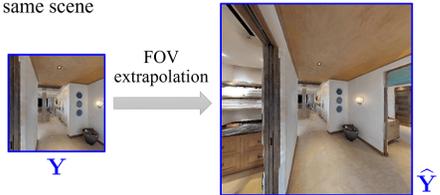


Fig. 1: Problem formulation of faithful FOV extrapolation. A collection of training images \mathbf{X}_i (red box) taken within a specific scene serves as prior knowledge. In the testing phase, our objective is to faithfully extrapolate the FOV of newly captured image \mathbf{Y} (blue box) to a specified FOV by leveraging the prior knowledge of the scene.

NeRF-Enhanced Outpainting (NEO) pipeline as a solution for faithful FOV extrapolation. (3) Comprehensive empirical investigations on both photorealistic and real-world datasets consistently validate the effectiveness of the proposed NEO method compared to the baseline counterparts.

II. RELATED WORK

A. Image Outpainting

Image outpainting, often referred to as image extrapolation or extension, is a task that seeks to expand image boundaries while maintaining semantically coherent content. Typically, the ability to infer such contents is acquired through learning from large-scale datasets of real images. Image outpainting approaches can be broadly categorized as either non-parametric or parametric. Non-parametric methods [25], [26] are restricted to basic pattern outpainting, and they become increasingly fragile as the extrapolation range grows. The emergence of GAN-based models has resulted in significant advancements in image outpainting. Some notable works [27], [9], [12], [5], [10] utilize a single GAN model for extrapolating the input image. More recently, Khurana *et al.* [11] propose an image outpainting framework that extends the image within the semantic label space, thereby produce new objects within the extrapolated area. Li *et al.* [6] introduce CTO-GAN, which deduces the potential semantic layout based on foreground elements and subsequently generates the corresponding background content with the guidance of the predicted semantics. Yao *et al.* [8] formulate this problem as a sequence-to-sequence autoregression task based on image patches, and present a query-based encoder-decoder transformer model to perform extrapolation. In addition to dedicated outpainting methods, certain image inpainting models [28], [29], [30], which have the capability to fill large masks, can also be adapted for image outpainting. Inspired by the pioneering work [31] on diffusion models, there have been endeavors [32], [33], [34], [35] that tackle the image outpainting problem via a diffusion-then-denoising process. One limitation of these pretrained image outpainting models is that the extrapolated parts lack geometric consistency and interpretability for a specific scene. This characteristic

renders them unsuitable for real-world application scenarios such as navigation. Therefore, we propose a new problem, faithful FOV extrapolation, the solution of which enables and facilitates navigational applications by ensuring that the extrapolated content remains faithful and relevant to the scene at hand.

B. NeRF and Data Augmentation

Neural radiance fields (NeRF) [24] enables novel view synthesis by representing the density and color of 3D spatial points of a specific scene through a neural network. With a novel camera pose as input, NeRF can render an image of specified FOV by performing ray marching from the camera’s central viewpoint, querying the corresponding color and density fields and conducting volume rendering. Follow-up works on NeRF have explored improving generalizability [36], [37], scene editing [38], [39], [40], neural scene reconstruction [41], [42], training and inference acceleration [43], [44], among others. Unlike other image outpainting techniques which rely on scene distribution priors to extrapolate large FoV, NeRF has its unique advantage in that it implicitly encodes the entire 3D scene, enabling the rendering of novel views in a manner that is both geometrically and semantically coherent. This positions NeRF as a potential approach to achieving faithful FOV extrapolation. Furthermore, NeRF has been utilized as a data augmentation tool to generate synthetic images for training deep neural networks. Moreau *et al.* [45] employ NeRF model to create a fresh dataset of synthetic images for training a camera pose regression model. Ge *et al.* [46] propose an online data augmentation pipeline based on NeRF synthesis for real-world object detection. In this paper, we present a NeRF-enhanced outpainting pipeline that leverages NeRF to generate sufficient synthetic images for training an FOV extrapolation model.

III. PROPOSED METHOD

A. Problem Formulation

As shown in Fig. 1, a camera captured N training images $\{\mathbf{X}_i \in \mathbb{R}^{h \times w \times 3} | i = 1, 2, \dots, N\}$ in a specific scene, which

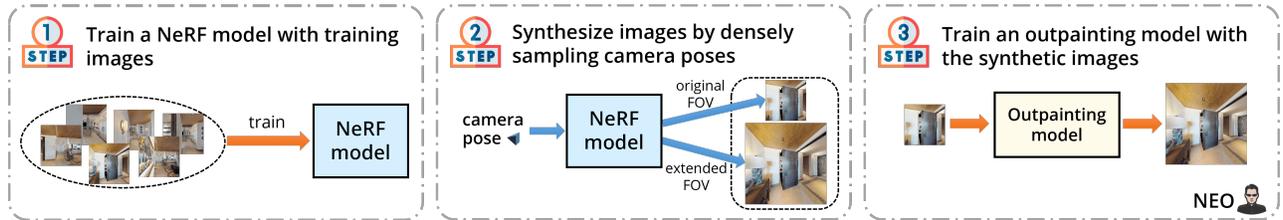


Fig. 2: **NEO pipeline:** (1) training a NeRF model with captured training images of original small FOV; (2) using the trained NeRF to synthesize images of extended FOV by densely sampling camera poses in the scene; (3) training an outpainting model with the synthetic images of extended FOV. During inference, we use the trained outpainting model for faithful FOV extrapolation.

may have been taken sparsely. The FOV of each image is (α_x, α_y) , indicating the horizontal and vertical FOV angles, respectively. We assume that the camera pose $\mathbf{P}_i \in \mathbf{SE}(3)$ for each training image \mathbf{X}_i can be acquired through structure from motion (SfM) pipelines such as COLMAP [16]. During testing, we seek to extrapolate a testing image $\mathbf{Y} \in \mathbb{R}^{h \times w \times 3}$ with the same FOV (α_x, α_y) captured in the same scene to a new image $\hat{\mathbf{Y}} \in \mathbb{R}^{H \times W \times 3}$ with larger FOV (γ_x, γ_y) while keeping the camera’s focal length constant. The extrapolated portions of the image must be consistent with the real scene.

B. NeRF-Enhanced Outpainting (NEO)

To address this problem, we propose a simple yet effective method, dubbed NeRF-Enhanced Outpainting (NEO) with three steps in the training stage, as illustrated in Fig. 2.

Step 1: Training a NeRF model. We start by training a NeRF model using training images $\{\mathbf{X}_i | i = 1, 2, \dots, N\}$. NeRF learns an implicit representation of the given 3D scene with a multilayer perceptron (MLP). With the trained NeRF, we can emit a ray from any direction and sample points on the ray to obtain their density and radiance, then render a novel view through volume rendering. Therefore, NeRF can generate an image of a specified FOV with a new camera pose. Using the trained NeRF, we can render an arbitrary number of images with specified poses and desired FOV.

Step 2: Synthesizing images. Next, we sample a multitude of new camera poses within the scene. For each pose, we leverage the trained NeRF model to render a pair of images with FOVs of (α_x, α_y) and (γ_x, γ_y) , respectively. We sample new camera poses by ensuring that they cover all walkable areas in the training trajectories. Moreover, we manage to sample poses with different degrees of freedom (DoF) following their DoF distribution in the specific environment. The details will be presented in Sec. IV.

Step 3: Training an outpainting model. Finally, we utilize the NeRF-rendered images as training data to train an image outpainting model, which takes the small-FOV images as input and extrapolates the large-FOV ones. Various image outpainting models [7] can be employed in this step. Image inpainting models [30], [29] that fill in large empty spaces can also be applied to outpainting. However, it is worth noting that some outpainting or inpainting models are designed to generate diverse results with randomness.

This property does not align with the objective of faithful FOV extrapolation. Thus, we modify such models during our implementation for the purpose of training an outpainting model that performs deterministic extrapolations faithful to the environment.

During inference, a real small-FOV image is given to the trained outpainting model as input and the model performs faithful FOV extrapolation to obtain the large-FOV image.

C. Discussions

1) *Why not directly train an outpainting model using the training images $\{\mathbf{X}_i | i = 1, 2, \dots, N\}$?*

The outpainting model takes $\mathbf{X}_i \in \mathbb{R}^{h \times w \times 3}$ as input and produces $\hat{\mathbf{X}}_i \in \mathbb{R}^{H \times W \times 3}$ as output. However, there are no $\hat{\mathbf{X}}_i$ in the training data. To address this issue, we could simply resize $\mathbf{X}_i \in \mathbb{R}^{h \times w \times 3}$ to $\mathbf{X}'_i \in \mathbb{R}^{H \times W \times 3}$ and then crop the central part $\mathbf{X}''_i \in \mathbb{R}^{h \times w \times 3}$ to use as input. Since the training is simply achieved through resizing and cropping the original small-FOV training images, we refer to this method as *naive outpainting*. However, this approach encounters two main challenges. First, the quantity and coverage of the training data prove inadequate for hallucinating from a new viewpoint. Second, cropping the central part reduces the FOV of the training image, leading to a mismatch of the FOVs during training and testing stages.

Our proposed NEO pipeline addresses the above two challenges by leveraging NeRF-based synthesis. First, there is no longer a concern about limited amount of training data, as NEO can theoretically generate an unlimited number of synthetic images by sampling arbitrary camera poses across walkable areas. Second, the issue of training-testing FOV mismatch is resolved, because the NEO pipeline trains the outpainting model using synthetic images with the same FOV as the target of the testing stage. The underlying principle of NEO is that the resulting outpainting model learns to extrapolate faithfully in the given scene by processing an extensive volume of training images that comprehensively cover the entire scene.

2) *Why not directly synthesize the target images using the trained NeRF model?*

NeRF can generate a specific extended-FOV image given a camera pose. However, during the testing phase, the camera pose of the testing image \mathbf{Y} is unknown. We could employ camera relocalization (also known as visual localization)

paradigms [15], [16] to estimate the camera pose of \mathbf{Y} . Yet a main issue with combining relocalization and NeRF is that, the estimated pose may not be precise enough due to errors in feature detection, matching, as well as perspective-n-point (PnP) [47] estimation. Consequently, the resulting extended-FOV image rendered by NeRF may not be aligned well with the testing image \mathbf{Y} . In contrast, the NEO pipeline circumvents the need for highly accurate estimation of the camera poses for testing images. Instead, it benefits from NeRF by training an outpainting model with the underlying scene priors from NeRF renderings.

IV. EXPERIMENTS

A. Baseline Methods

To demonstrate the effectiveness of our proposed NEO, we first introduce three baseline methods to address this problem, as we have discussed earlier. The first method is the “naive outpainting” mentioned in Sec. III-C, where an outpainting model is trained only using the original small-FOV images. The original image is resized larger and its cropped central part is used as the small-FOV input, while the original image is used as the corresponding large-FOV target. The second baseline approach, dubbed “warping & fusion” (B2), goes through an image stitching pipeline. To be specific, for a testing image, we first retrieve the nearest images from the training set, then build correspondence between the testing image and the retrieved (source) images. Finally, image warping is employed to get a larger FOV image by warping and fusing the contents from source images. The last baseline is called “relocalized NeRF”. We first train a NeRF model on the training images. For a testing image, we relocalize its camera pose by employing camera relocalization methods. Finally we render a large-FOV image with the trained NeRF and the relocalized pose.

B. Datasets and Metrics

We first evaluate the FOV extrapolation performance on three scenes from three photorealistic datasets: 1) **Replica** dataset [51] includes 18 realistic indoor scenes. We adopt the first floor of the `apartment0` scene for evaluation. 2) **Gibson** dataset [52] includes realistic scans of 572 full buildings. We adopt the `Bonesteel` building to verify our method. 3) **Habitat-Matterport 3D (HM3D)** dataset [53] includes realistic scans of 1,000 buildings. We adopt scene `00065` for evaluation. To verify our method on real scenes, we further demonstrate our method on **ScanNet** [54] which samples posed RGB images from 1513 real indoor scans. We adopt scene `00000_00` as our training set and scene `00000_01` whose data is sampled at the same scene but with different trajectories as our testing set.

For photorealistic datasets, we consider a simplified scenario: a robot with a fixed height and a fixed front camera navigates in an indoor environment. In such a setting, the camera has a constant height and can only rotate horizontally, so the motion of the camera has only 3 DoFs. For each scene, we use the Habitat environment [55] to render 1,000 training images with 256×256 resolution and 90° FOV from random

camera poses at a fixed height of $1.5m$. We then render 2,000 testing images with the same resolution and FOV, which contains 40 random walking paths at the same height with 50 images on each path. The target resolution during testing is 512×512 (126.87° FOV) by uniformly extrapolating in four directions. For ScanNet, we simulate a challenging but more realistic scenario to extrapolate an input central image with resolution 240×160 in horizontal directions (left and right sides), whose target resolution is 240×320 . We uniformly sample the original training trajectory by image IDs with an interval 10, leading to 558 training images. For testing, we randomly sample three pieces of continuous trajectory, each with 200 images, leading to 600 testing views.

Since we require the extrapolated regions to be consistent with the scene, we adopt PSNR, SSIM and LPIPS [56] as the evaluation metrics for faithful FOV extrapolation.

C. Implementation Details

Photorealistic Datasets. For Replica, Gibson, and HM3D datasets, we use MAT [30] as our default outpainting model but remove its style manipulation module for deterministic extrapolation. For warping & fusion baseline (B2), we employ the pipeline proposed in [48] as our implementation. For camera relocalization, we first apply an image retrieval method NetVLAD [57] to retrieve the nearest images from the training set w.r.t. a testing image, then run COLMAP [16] which generalizes well to different environments to get the relative pose between the testing image and the retrieved training image. Finally we transform the relative pose to the absolute pose using the known training pose. For NeRF model, we employ DirectVoxGO [49] which leverages 3D voxel representation to accelerate training and inference.

For new data generation, we first derive the walkable areas from the 2D floor plans of the datasets. We sample new camera poses on a 2D horizontal grid with a default interval of $0.05m$. At each position, we uniformly sample 72 yaw angles for the horizontal rotation. As a result, we sampled about 1.63 million, 1.43 million, 1.56 million camera poses for Replica, Gibson, and HM3D, respectively.

ScanNet. For ScanNet [54], we instead use LaMa [29] for outpainting, to satisfy our need on different resolution and aspect ratios. For NeRF model, we replace DirectVoxGO with a state-of-the-art NeRF method [50] on ScanNet to achieve better rendering performance. Other baselines are evaluated in a similar manner as on the photorealistic ones.

For new data generation, we aim to generate 6-DoF novel poses whose distribution is similar to the training trajectory. Specifically, on x - y plane, we sample new camera poses on a uniform 2D horizontal grid with an interval of $0.2m$, and ensure the sampled trajectories are roughly covered by the training set, *i.e.*, the Euclidean distance from the new pose to the nearest training pose on x - y plane should be limited within a threshold of $0.3m$. For horizontal rotation (yaw angle), we sample from a uniform distribution whose upper and lower bounds come from the training poses. For other DoFs (vertical translation, pitch and roll rotation), we empirically discover that the training poses

TABLE I: **Quantitative evaluation on four datasets.** The backbone for computing LPIPS metrics is VGG network.

Method	Replica (photorealistic)			Gibson (photorealistic)			HM3D (photorealistic)			ScanNet (real)		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
(B1) Naive Outpainting [30]	20.59	0.781	0.348	18.05	0.705	0.404	17.92	0.630	0.427	19.98	0.755	0.188
(B2) Warping & Fusion [48]	19.03	0.745	0.397	16.61	0.688	0.496	15.94	0.582	0.512	21.02	0.755	0.238
(B3) Relocalized NeRF [49], [50], [16]	16.78	0.724	0.386	14.90	0.641	0.474	14.43	0.602	0.484	18.88	0.695	0.188
Oracle NeRF [49], [50]	32.90	0.936	0.174	32.16	0.928	0.188	27.03	0.824	0.299	23.80	0.805	0.108
NEO	25.94	0.868	0.217	23.53	0.822	0.263	21.54	0.731	0.338	22.40	0.793	0.168



Fig. 3: Qualitative results on three photorealistic datasets: **Replica** (1st row), **Gibson** (2nd row), and **HM3D** (3rd row).



Fig. 4: Qualitative results on **ScanNet** dataset.

roughly follow a Gaussian distribution, thus we sample from a Gaussian distribution whose mean and standard deviation are calculated from the training poses. Then we employ NeRF to render novel views using these poses. However, practically we found that a non-negligible portion of the rendered images on ScanNet are too blurry to provide useful information for training the outpainting model. To address this problem, we apply a blur detection algorithm based on Laplacian variance to compute the blurry scores for filtering out blurry images. Eventually, we generated 0.75 million outpainting-trainable images on ScanNet.

D. Results

Quantitative Results. Table I illustrates the extrapolation results of the proposed NEO approach and three baseline

methods on four datasets. In addition, for validation purposes, the performance of “oracle NeRF” is reported, which employs the groundtruth camera pose of the test image for NeRF rendering. “Oracle NeRF” reflects the quality of the trained NeRF model and sets an upper bound for (B3) relocalized NeRF. Since NEO learns the outpainting from NeRF-augmented images, its performance is expected to be lower than that of “oracle NeRF”. As shown in Table I, NEO significantly outperforms the three baseline methods. As anticipated, “oracle NeRF” achieves the best results although it is not for practical use since groundtruth camera poses for testing images are unknown. Although the three baseline methods produce reasonably good results, they encounter non-trivial problems with faithfulness, which we will demonstrate in qualitative results below.

TABLE II: Effect of (a) *pose sampling density* and (b) *FOV of training images* in NEO pipeline on Replica dataset.

	Interval	# Pose	FOV	PSNR \uparrow	SSIM \uparrow
(a)	0.05	1,629,864	extended	25.94	0.868
	0.10	406,224	extended	24.64	0.851
	0.20	101,304	extended	24.52	0.849
	0.50	16,056	extended	23.37	0.833
	1.00	3,744	extended	22.72	0.824
(b)	0.05	1,629,864	extended	25.94	0.868
	0.05	1,629,864	original	20.61	0.802

Qualitative Results. We show some qualitative results on Replica, Gibson, and HM3D datasets in Fig. 3. “Oracle NeRF” learns a geometrically consistent 3D representation thus demonstrates appealing rendering in a coherent way on the three datasets. The areas extrapolated by NEO are also much more accurate and faithful to the scene compared to the three baselines. NEO sometimes suffers from slight misalignment around small objects (*e.g.*, the painting on the left wall on the third row of Fig. 3). The extrapolated regions of (B1) naive outpainting tend to be blurry, which is mainly caused by the limited number of training images. The extended areas by (B2) warping & fusion are limited by the non-overlapping regions of neighboring images. As for (B3) relocalized NeRF, the input region (central part) always misaligns with the extrapolated regions due to evident errors in pose estimation. Comparison of Results on the ScanNet dataset, as shown in Fig. 4, leads to similar observations. Surprisingly, we found that NEO can avoid some issues encountered by “oracle NeRF” and achieve better visual quality in some regions, such as the blurry floor region near the border of extrapolated image by “oracle NeRF” in the second row of Fig. 4. We suspect the reason is that, NeRF is trained on sparse views from the original training set, thereby its rendering quality in specific areas may be dependent on the availability of informative, overlapping training images. In contrast, NEO trains the outpainting model on sufficient, dense novel views rendered from NeRF, so it is more capable of learning a semantically and geometrically coherent color field of the scene, effectively reducing the impact of insufficient information in some areas of the original NeRF.

E. Discussions

Pose Sampling. It is important in the NEO approach to cover as many views as possible in the scene in Step 2 of the training process. Thus, the distribution and number of sampled poses are crucial for training a highly effective outpainting model. In this study, we vary the interval of the 2D grid to control the sampling density on the Replica dataset. As shown in Table II (a), the generative performance naturally improves when increasing the pose sampling density. The improvement is significant (+1.30 in PSNR) when reducing the interval from $0.1m$ to $0.05m$, where the number of sampled poses increases from $0.4M$ to $1.6M$.

FOV of Training Images. A key issue of training an outpainting model for FOV extrapolation is the consistency in

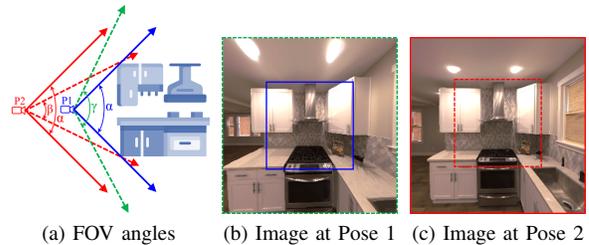


Fig. 5: FOV analysis. The discrepancy in the FOV between the (c) training image and (b) testing image may lead to false extrapolation behaviors.

FOV between training and testing images. Figure 5 demonstrates the FOV mismatch problem in the naive outpainting method. The solid red and blue arrows in Fig. 5(a) represents the camera’s inherent FOV α . From a resized training image captured at Pose 2 (Fig. 5(c)), naive outpainting learns to extrapolate the cropped FOV β (dashed red) to α . However, for a testing input image (solid blue) at Pose 1 (Fig. 5(b)), the goal is to extrapolate the inherent FOV α to a larger FOV γ (dashed green). Though the central parts (inputs) of Fig. 5(b) and Fig. 5(c) are similar, their extrapolated parts are totally different. The FOV mismatch issue of naive outpainting can also be observed in the qualitative results (*e.g.*, the painting on the second row in Fig. 3). To further examine the effect of FOV, we evaluate a variant of NEO, which uses the original-FOV synthetic images to train the same outpainting model. As seen in Table II (b), the performance greatly decreases (-5.33 in PSNR), indicating the significance of training FOV.

V. LIMITATIONS AND CONCLUSIONS

As an initial exploration of the faithful FOV extrapolation task, this paper focuses on tackling the problem for only static scenes. However, real-world navigation usually entails dynamic objects or people. Moreover, scenes are rarely static over time, *e.g.*, furniture may be rearranged. This study can serve as a probe to more comprehensive research in this area. In the future, we envision to explore solutions that can accommodate the complexities of more realistic scenarios. One way may leverage dynamic NeRF [58] that better handles dynamic scenarios.

To conclude, in this paper, we formulate a new problem named *faithful image extrapolation* to increase FOV of a given image. It requires the expanded area to adhere to the real environment. To address this problem, inspired by the recent surge of NeRF-based rendering approaches, we propose a novel pipeline dubbed NEO, to train a NeRF-enhanced image outpainting model. Our key insight is to obtain sufficient and interpretable training data to aid the training of outpainting model from the novel views rendered by NeRF on a specific scene. Compared with competing baselines, our model has showcased superior generative performance. Our synthesized views are geometrically and semantically consistent with the 3D environment, thereby achieving faithful extrapolation that opens up potential applications such as AR-based navigation.

REFERENCES

- [1] Y. Zhu, R. Mottaghi, E. Kolve, J. J. Lim, A. Gupta, L. Fei-Fei, and A. Farhadi, "Target-driven visual navigation in indoor scenes using deep reinforcement learning," in *IEEE International Conference on Robotics and Automation, ICRA*, 2017, pp. 3357–3364.
- [2] P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson, N. Sünderhauf, I. D. Reid, S. Gould, and A. van den Hengel, "Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2018, pp. 3674–3683.
- [3] S. Lee, R. Yu, J. Xie, S. M. Billah, and J. M. Carroll, "Opportunities for human-AI collaboration in remote sighted assistance," in *ACM International Conference on Intelligent User Interfaces, IUI*, 2022, pp. 63–78.
- [4] J. Xie, R. Yu, S. Lee, Y. Lyu, S. M. Billah, and J. M. Carroll, "Helping helpers: Supporting volunteers in remote sighted assistance with augmented reality maps," in *ACM Conference on Designing Interactive Systems, DIS*, 2022, pp. 881–897.
- [5] Z. Yang, J. Dong, P. Liu, Y. Yang, and S. Yan, "Very long natural scenery image prediction by outpainting," in *IEEE/CVF International Conference on Computer Vision, ICCV*, 2019, pp. 10560–10569.
- [6] J. Li, C. Chen, and Z. Xiong, "Contextual outpainting with object-level contrastive learning," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2022, pp. 11441–11450.
- [7] Y. Cheng, C. H. Lin, H. Lee, J. Ren, S. Tulyakov, and M. Yang, "Inout: Diverse image outpainting via GAN inversion," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, pp. 11421–11430.
- [8] K. Yao, P. Gao, X. Yang, J. Sun, R. Zhang, and K. Huang, "Outpainting by queries," in *European Conference on Computer Vision, ECCV*, 2022, pp. 153–169.
- [9] Y. Wang, X. Tao, X. Shen, and J. Jia, "Wide-context semantic image extrapolation," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2019, pp. 1399–1408.
- [10] D. Guo, H. Liu, H. Zhao, Y. Cheng, Q. Song, Z. Gu, H. Zheng, and B. Zheng, "Spiral generative network for image extrapolation," in *European Conference on Computer Vision, ECCV*, 2020, pp. 701–717.
- [11] B. Khurana, S. R. Dash, A. Bhatia, A. Mahapatra, H. Singh, and K. Kulkarni, "Semie: Semantically-aware image extrapolation," in *IEEE/CVF International Conference on Computer Vision, ICCV*, 2021, pp. 14880–14889.
- [12] D. Krishnan, P. Teterwak, A. Sarna, A. Maschinot, C. Liu, D. Belanger, and W. T. Freeman, "Boundless: Generative adversarial networks for image extension," in *IEEE/CVF International Conference on Computer Vision, ICCV*, 2019, pp. 10520–10529.
- [13] B. Zhou, À. Lapedriza, A. Khosla, A. Oliva, and A. Torralba, "Places: A 10 million image database for scene recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 6, pp. 1452–1464, 2018.
- [14] P. Lindenberger, P.-E. Sarlin, V. Larsson, and M. Pollefeys, "Pixel-perfect structure-from-motion with featuremetric refinement," in *IEEE/CVF International Conference on Computer Vision, ICCV*, 2021, pp. 5987–5997.
- [15] P. Sarlin, A. Unagar, M. Larsson, H. Germain, C. Toft, V. Larsson, M. Pollefeys, V. Lepetit, L. Hammarstrand, F. Kahl, and T. Sattler, "Back to the feature: Learning robust camera localization from pixels to pose," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2021, pp. 3247–3257.
- [16] J. L. Schonberger and J.-M. Frahm, "Structure-from-motion revisited," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2016, pp. 4104–4113.
- [17] R. Szeliski, "Image alignment and stitching: A tutorial," *Found. Trends Comput. Graph. Vis.*, vol. 2, no. 1, 2006.
- [18] N. Li, Y. Xu, and C. Wang, "Quasi-homography warps in image stitching," *IEEE Trans. Multimed.*, vol. 20, no. 6, pp. 1365–1375, 2018.
- [19] T. Liao and N. Li, "Single-perspective warps in natural image stitching," *IEEE Trans. Image Process.*, vol. 29, pp. 724–735, 2020.
- [20] K. Lee and J. Sim, "Warping residual based image stitching for large parallax," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2020, pp. 8195–8203.
- [21] M. Wang, A. Shamir, G.-Y. Yang, J.-K. Lin, G.-W. Yang, S.-P. Lu, and S.-M. Hu, "Biggerselfie: Selfie video expansion with hand-held camera," *IEEE Trans. Image Process.*, vol. 27, no. 12, pp. 5854–5865, 2018.
- [22] S. Lee, J. Lee, B. Kim, K. Kim, and J. Noh, "Video extrapolation using neighboring frames," *ACM Trans. Graph.*, vol. 38, no. 3, pp. 20:1–20:13, 2019.
- [23] L. Ma, S. Georgoulis, X. Jia, and L. Van Gool, "Fov-net: Field-of-view extrapolation using self-attention and uncertainty," *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 4321–4328, 2021.
- [24] 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," in *European Conference on Computer Vision, ECCV*, 2020, pp. 405–421.
- [25] A. A. Efros and T. K. Leung, "Texture synthesis by non-parametric sampling," in *IEEE/CVF International Conference on Computer Vision, ICCV*, 1999, pp. 1033–1038.
- [26] A. A. Efros and W. T. Freeman, "Image quilting for texture synthesis and transfer," in *SIGGRAPH*, L. P. Poock, Ed., 2001, pp. 341–346.
- [27] R. A. Yeh, C. Chen, T. Lim, A. G. Schwing, M. Hasegawa-Johnson, and M. N. Do, "Semantic image inpainting with deep generative models," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2017, pp. 6882–6890.
- [28] S. Zhao, J. Cui, Y. Sheng, Y. Dong, X. Liang, E. I. Chang, and Y. Xu, "Large scale image completion via co-modulated generative adversarial networks," in *International Conference on Learning Representations, ICLR*, 2021.
- [29] R. Suvorov, E. Logacheva, A. Mashikhin, A. Remizova, A. Ashukha, A. Silvestrov, N. Kong, H. Goka, K. Park, and V. Lempitsky, "Resolution-robust large mask inpainting with fourier convolutions," in *IEEE/CVF Winter Conference on Applications of Computer Vision, WACV*, 2022, pp. 3172–3182.
- [30] W. Li, Z. Lin, K. Zhou, L. Qi, Y. Wang, and J. Jia, "MAT: mask-aware transformer for large hole image inpainting," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2022, pp. 10748–10758.
- [31] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020.
- [32] C. Saharia, W. Chan, H. Chang, C. Lee, J. Ho, T. Salimans, D. Fleet, and M. Norouzi, "Palette: Image-to-image diffusion models," in *ACM SIGGRAPH*, 2022, pp. 1–10.
- [33] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2022, pp. 10684–10695.
- [34] A. Lugmayr, M. Danelljan, A. Romero, F. Yu, R. Timofte, and L. Van Gool, "Repaint: Inpainting using denoising diffusion probabilistic models," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2022, pp. 11461–11471.
- [35] S. Cai, E. R. Chan, S. Peng, M. Shahbazi, A. Obukhov, L. Van Gool, and G. Wetzstein, "Diffdreamer: Towards consistent unsupervised single-view scene extrapolation with conditional diffusion models."
- [36] K. Schwarz, Y. Liao, M. Niemeyer, and A. Geiger, "Graf: Generative radiance fields for 3d-aware image synthesis," *Advances in Neural Information Processing Systems*, vol. 33, pp. 20154–20166, 2020.
- [37] A. Yu, V. Ye, M. Tancik, and A. Kanazawa, "pixelnerf: Neural radiance fields from one or few images," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2021, pp. 4578–4587.
- [38] S. Liu, X. Zhang, Z. Zhang, R. Zhang, J.-Y. Zhu, and B. Russell, "Editing conditional radiance fields," in *IEEE/CVF International Conference on Computer Vision, ICCV*, 2021, pp. 5773–5783.
- [39] J. Zhang, X. Liu, X. Ye, F. Zhao, Y. Zhang, M. Wu, Y. Zhang, L. Xu, and J. Yu, "Editable free-viewpoint video using a layered neural representation," *ACM Transactions on Graphics (TOG)*, vol. 40, no. 4, pp. 1–18, 2021.
- [40] C. Wang, M. Chai, M. He, D. Chen, and J. Liao, "Clip-nerf: Text-and-image driven manipulation of neural radiance fields," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2022, pp. 3835–3844.
- [41] P. Wang, L. Liu, Y. Liu, C. Theobalt, T. Komura, and W. Wang, "Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction," *arXiv preprint arXiv:2106.10689*, 2021.
- [42] L. Yariv, J. Gu, Y. Kasten, and Y. Lipman, "Volume rendering of neural implicit surfaces," *Advances in Neural Information Processing Systems*, vol. 34, pp. 4805–4815, 2021.
- [43] T. Müller, A. Evans, C. Schied, and A. Keller, "Instant neural graphics primitives with a multiresolution hash encoding," *ACM Trans. Graph.*, vol. 41, no. 4, pp. 102:1–102:15, 2022.

- [44] A. Chen, Z. Xu, A. Geiger, J. Yu, and H. Su, "Tensorf: Tensorial radiance fields," in *European Conference on Computer Vision, ECCV*, vol. 13692, 2022, pp. 333–350.
- [45] A. Moreau, N. Piasco, D. Tsishkou, B. Stanciulescu, and A. de La Fortelle, "Lens: Localization enhanced by nerf synthesis," in *Conference on Robot Learning*, 2022, pp. 1347–1356.
- [46] Y. Ge, H. Behl, J. Xu, S. Gunasekar, N. Joshi, Y. Song, X. Wang, L. Itti, and V. Vineet, "Neural-sim: Learning to generate training data with nerf," in *European Conference on Computer Vision, ECCV*, 2022, pp. 477–493.
- [47] M. A. Fischler and R. C. Bolles, "Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography," *Communications of the ACM*, vol. 24, no. 6, pp. 381–395, 1981.
- [48] Y. Liu, W. Lai, M. Yang, Y. Chuang, and J. Huang, "Hybrid neural fusion for full-frame video stabilization," in *IEEE/CVF International Conference on Computer Vision, ICCV*, 2021, pp. 2279–2288.
- [49] C. Sun, M. Sun, and H. Chen, "Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2022, pp. 5449–5459.
- [50] P. Dai, Y. Zhang, X. Yu, X. Lyu, and X. Qi, "Hybrid neural rendering for large-scale scenes with motion blur," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2023, pp. 154–164.
- [51] J. Straub, T. Whelan, L. Ma, Y. Chen, E. Wijmans, S. Green, J. J. Engel, R. Mur-Artal, C. Ren, S. Verma, *et al.*, "The replica dataset: A digital replica of indoor spaces," *arXiv:1906.05797*, 2019.
- [52] F. Xia, A. R. Zamir, Z. He, A. Sax, J. Malik, and S. Savarese, "Gibson env: Real-world perception for embodied agents," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2018, pp. 9068–9079.
- [53] S. K. Ramakrishnan, A. Gokaslan, E. Wijmans, O. Maksymets, A. Clegg, J. Turner, E. Undersander, W. Galuba, A. Westbury, A. X. Chang, *et al.*, "Habitat-matterport 3d dataset (HM3d): 1000 large-scale 3d environments for embodied AI," *arXiv:2109.08238*, 2021.
- [54] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Nießner, "ScanNet: Richly-annotated 3d reconstructions of indoor scenes," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2017, pp. 5828–5839.
- [55] M. Savva, A. Kadian, O. Maksymets, Y. Zhao, E. Wijmans, B. Jain, J. Straub, J. Liu, V. Koltun, J. Malik, *et al.*, "Habitat: A platform for embodied AI research," in *IEEE/CVF International Conference on Computer Vision, ICCV*, 2019, pp. 9339–9347.
- [56] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2018, pp. 586–595.
- [57] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "Netvlad: Cnn architecture for weakly supervised place recognition," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2016, pp. 5297–5307.
- [58] Z. Li, Q. Wang, F. Cole, R. Tucker, and N. Snavely, "Dynibar: Neural dynamic image-based rendering," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2023, pp. 4273–4284.