# Unconstrained Scene Generation with Locally Conditioned Radiance Fields

Terrance DeVries[†1]    Miguel Angel Bautista[1]    Nitish Srivastava[1]

Graham W. Taylor[2,3]    Joshua M. Susskind[1]

[1]Apple    [2]University of Guelph    [3]Vector Institute

Figure 1: Scenes sampled from our learned prior, rendered from freely moving camera paths containing various rotation, translation, and forward/backward motions.

## Abstract

*We tackle the challenge of learning a distribution over complex, realistic, indoor scenes. In this paper, we introduce **Generative Scene Networks (GSN)**, which learns to decompose scenes into a collection of many local radiance fields that can be rendered from a free moving camera. Our model can be used as a prior to generate new scenes, or to complete a scene given only sparse 2D observations. Recent work has shown that generative models of radiance fields can capture properties such as multi-view consistency and view-dependent lighting. However, these models are specialized for constrained viewing of single objects, such as cars or faces. Due to the size and complexity of realistic indoor environments, existing models lack the representational capacity to adequately capture them. Our decomposition scheme scales to larger and more complex scenes while preserving details and diversity, and the learned prior enables high-quality rendering from viewpoints that are significantly different from observed viewpoints. When compared to existing models, GSN produces quantitatively higher-quality scene renderings across several different scene datasets.*

## 1. Introduction

Spatial understanding entails the ability to infer the geometry and appearance of a scene when observed from any viewpoint or orientation given sparse or incomplete observations. Although a wide range of geometry and learning-based approaches for 3D view synthesis [8, 9, 33, 34, 43, 51, 62] can interpolate between observed views of a scene, they cannot extrapolate to infer unobserved parts of the scene. The fundamental limitation of these models is their inability to learn a prior over scenes. As a result, learning-based approaches have limited performance in the extrapolation regime, whether it be inpainting disocclusions or synthesizing views beyond the boundaries of the given observations. For example, the popular NeRF [34] approach represents a scene via its radiance field, enabling continuous view interpolation given a densely captured scene. However, since NeRF does not learn a scene prior, it cannot extrapolate views. On the other hand, conditional auto-encoder models for view synthesis [12, 63, 53, 11, 39, 54] are able to extrapolate views of simple objects from multiple viewpoints and orientations. Yet, they overfit to viewpoints seen during

---

† Work done as part of an internship at Apple.

| Model | Multiple Scenes/Objects | Generative | Latent | Radiance Field | Scene Level | Camera Placement |
|---|---|---|---|---|---|---|
| NeRF [34] | ✗ | ✗ | - | ✓ | ✓ | Sphere/Wide-baseline |
| NSVF [30] | ✗ | ✗ | 3D | ✓ | ✓ | Sphere/Wide-baseline |
| PixelNerf [58] | ✓ | ✗ | N×2D | ✓ | ✓ | Sphere/Wide-baseline |
| GRAF [49] | ✓ | ✓ | 1D | ✓ | ✗ | Sphere |
| HoloGAN [37] | ✓ | ✓ | 1D | ✗ | ✗ | Sphere |
| PlatonicGAN [16] | ✓ | ✗ | 1D | ✗ | ✗ | Sphere |
| ENR [11] | ✓ | ✗ | 3D | ✗ | ✗ | Sphere |
| GTM-SM [13] | ✓ | ✓ | 1D | ✗ | ✓ | Free moving |
| ISS [41] | ✓ | ✗ | 2D | ✗ | ✓ | Free moving |
| GSN (ours) | ✓ | ✓ | 2D | ✓ | ✓ | Free moving |

Table 1: Summary of contributions and comparison with relevant related work. (**Multiple Scene/Objects**): Ability to model multiple scenes/objects in the same network. (**Generative**) Whether the model is generative (*e.g.* allows for free sampling). (**Latent**) Latent code spatial dim. (**Radiance Field**) Whether the model predicts a radiance field. (**Scene level**) Results demonstrated in scene-level environments. (**Camera Placement**) What camera motion is permitted?

training (see [7] for a detailed analysis). Moreover, conditional auto-encoders tend to favor point estimates (*e.g.* the distribution mean) and produce blurry renderings when extrapolating far from observed views as a result.

A learned prior for scenes may be used for *unconditional* or *conditional* inference. A compelling use case for unconditional inference is to generate realistic scenes and freely move through them in the absence of input observations, relying on the prior distribution over scenes (see Fig. 1 for examples of trajectories of a freely moving camera on scenes sampled from our model). Likewise, conditional inference lends itself to different types of problems. For instance, plausible scene completions may be sampled by inverting scene observations back to the learned scene prior [57]. A generative model for scenes would be a practical tool for tackling a wide range of problems in machine learning and computer vision, including model-based reinforcement learning [15], SLAM [5, 6], content creation [24], and adaptation for AR/VR or immersive 3D photography.

In this paper we introduce *Generative Scene Networks* (GSN), a generative model of scenes that allows view synthesis of a freely moving camera in an open environment. Our contributions are the following. We: *(i)* introduce the first generative model for unconstrained scene-level radiance fields; *(ii)* demonstrate that decomposing a latent code into a grid of locally conditioned radiance fields results in an expressive and robust scene representation, which outperforming strong baselines; *(iii)* infer observations from arbitrary cameras given a sparse set of observations by inverting GSN (*i.e.*, fill in the scene); and *(iv)* show that GSN can be trained on multiple scenes to learn a rich scene prior, while rendered trajectories are smooth and consistent, maintaining scene coherence.

## 2. Related Work

An extensive body of literature has tackled view synthesis of simple synthetic objects [4] against a homogeneous background [11, 39, 42, 53, 63]. Relevant models typically rely on predicting 2D/3D flow fields to transform pixels in the input view to pixels in the predicted view. Moreover, in [50] a continuous representation is introduced in the form of the weights of a parametric vector field, which is used to infer appearance of 3D points expressed in a world coordinate system. At inference time the model is optimized to find latent codes that are predictive of the weights. The above approaches rely on the existence of training data of the form *input-transformation-output* tuples, where at training time the model has access to two views (*input* and *output*) in a viewing sphere and the corresponding 3D transformation. Sidestepping this requirement [16, 37] showed that similar models can be trained in an adversarial setting without access to *input-transformation-output* tuples.

Recent approaches [30, 34, 43, 45, 46] have been shown to encode a continuous radiance field in the parameters of a neural network by fitting it to a collection of high resolution views of realistic scenes. The encoded radiance field is rendered using the classical volume rendering equation [21] and a reconstruction loss is used to optimize the model parameters. While successful at modelling the radiance field at high resolutions, these approaches require optimizing a new model for every scene (where optimization usually takes days on commodity hardware). Thus, the fundamental limitation of these works is that they are unable to represent multiple scenes within the same model. As a result, these models cannot learn a prior distribution over multiple scenes.

Some newer works have explore generative radiance field models in the adversarial setting [3, 38, 49]. However, [3, 49] are restricted to modeling single objects where
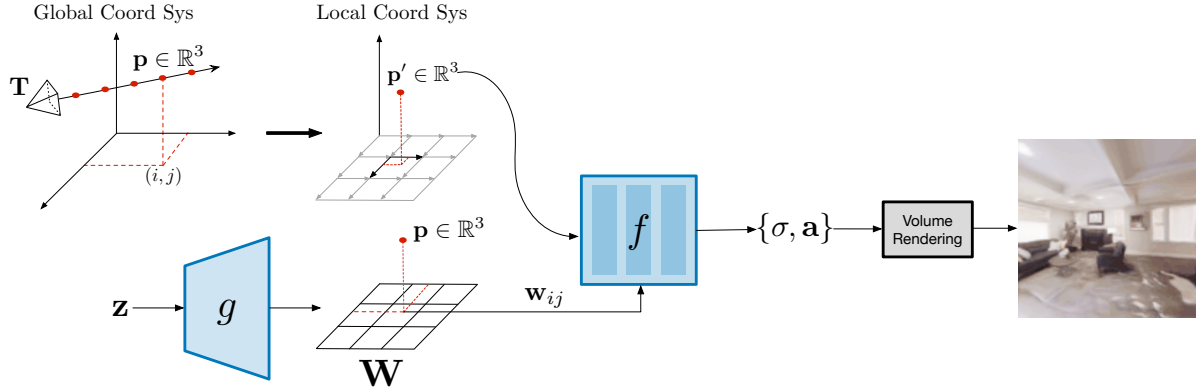
Figure 2: Architecture of the GSN generator. We sample a latent code $\mathbf{z} \sim p_z$ that is fed to our global generator $g$ producing a local latent grid $\mathbf{W}$. This local latent grid $\mathbf{W}$ conceptually represents a latent scene "floorplan" and is used for locally conditioning a radiance field $f$ from which images are rendered via volumetric rendering. For a given point $\mathbf{p}$ expressed in a global coordinate system to be rendered, we sample $\mathbf{W}$ at the location $(i, j)$, given by $\mathbf{p}$ resulting in $\mathbf{w}_{ij}$. In turn $f$ takes as input $\mathbf{p}'$ which results from expressing $\mathbf{p}$ relative to the local coordinate system of $\mathbf{w}_{ij}$.

a camera is positioned on a viewing sphere oriented towards the object, and they use a 1D latent representation, which prevents them from efficiently scaling to full scenes (c.f. § 4 for experimental evidence). Concurrent to our work, Niemeyer et al. [38] model a scene with multiple entities but report results on single-object datasets and simple synthetic scenes with 2-5 geometrical shapes in the CLEVR dataset [20]. The object-level problem may be solved by learning a pixel-aligned representation for conditioning the radiance field [55, 58, 44, 56]. However, [55, 58, 44, 56] are designed as conditional auto-encoders as opposed to a generative model, thereby ignoring the fundamentally stochastic nature of the view synthesis problem. In [28] a model is trained on long trajectories of nature scenarios. However, due to their iterative refinement approach the model has no persistent scene representation. While allowing for perpetual generation, the lack of a persistent scene representation limits its applicability for other downstream tasks that may require loop closure.

Approaches tackling view synthesis for a freely moving camera in a scene offer the capability to fully explore a scene [12, 13, 41]. Irrespective of the output image resolution, this is strictly a more complex problem than a camera moving on a sphere oriented towards a single object [50, 63, 53, 58, 49, 38, 3]. For scenes with freely moving cameras the model needs to learn to represent a full scene that consists of a multitude of objects rather than a single object. In this setup, it is often useful to equip models with a memory mechanism to aggregate the set of incoming observations. In particular, it has been useful to employ a dictionary-based memory to represent an environment where keys are camera poses and values are latent observation representations [13]. Furthermore, the memory may be extended to a 2D feature map that represents

at top-down view of the environment [12, 41]. At inference time, given a query viewpoint the model queries the memory and predicts the expected observation. This dependency on stored observations significantly limits these models' ability to cope with unseen viewpoints. GSN can perform a similar task by matching observations to scenes in its latent space, but with the added benefit that the learned scene prior allows for extrapolation beyond the given observations. A summarized comparison of GSN with the most relevant related work is shown in Tab. 1.

## 3. Method

As in traditional GANs [14], our model is composed of a generator $G_\theta$ and a discriminator $D_\phi$, with latent codes $\mathbf{z} \sim p_z$, where $p_z$ denotes the prior distribution. The generator $G_\theta$ is composed of two subnetworks $G = g_{\theta_g} \cdot f_{\theta_f}$: a global generator $g_{\theta_g}$ and a locally conditioned radiance field network $f_{\theta_f}$ (in the rest of the text and figures we drop the subscripts for clarity). Unlike standard 2D GANs, which often only require samples $\mathbf{z} \sim p_z$ as input, GSN also takes a camera pose $\mathbf{T} \in SE(3)$ from an empirical pose distribution $p_T$, as well as camera intrinsic parameters $\mathbf{K} \in \mathbb{R}^{3 \times 3}$. These additional inputs allow for explicit control over the viewpoint and field of view of the output image, which is rendered via a radiance field. Fig. 2 outlines the GSN network architecture.

### 3.1. Global Generator

In GRAF [49] and $\pi$-GAN [3] a 1D global latent code $\mathbf{z} \sim p_z$ is used to condition an MLP which parameterizes a single radiance field. Although such 1D global latent code can be effective when representing individual object categories, such as cars or faces, it does not scale well to large,
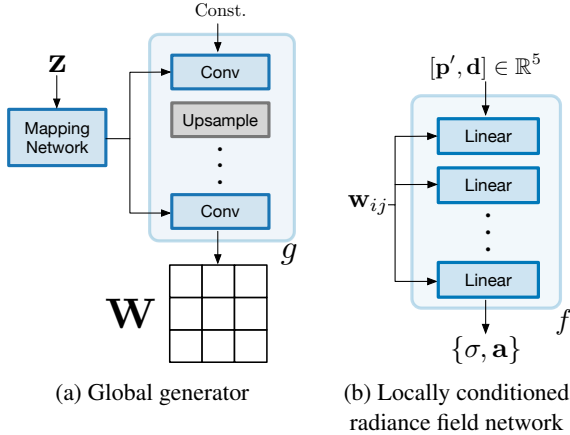
Figure 3: (a) Architecture of global generator $g$. We use a mapping network, modulated convolutional blocks, and a learned constant input as in StyleGAN2 [25]. (b) Architecture of the locally conditioned radiance field network $f$. Latent code $\mathbf{w}_{ij}$, sampled from $\mathbf{W}$, is used to modulate linear layers, similar to CIPS [2].

densely populated scenes. Simply increasing the capacity of a single radiance field network has diminishing returns with respect to render quality [43]. It is more effective to distribute a scene among several smaller networks, such that each can specialize on representing a local region [43]. Inspired by this insight we propose the addition of a global generator that learns to decompose large scenes into many smaller, more manageable pieces.

The global generator $g$ maps from the global latent code $\mathbf{z} \in \mathbb{R}^d$, which represents an entire scene, to a 2D grid of local latent codes $\mathbf{W} \in \mathbb{R}^{c \times s \times s}$, each of which represent a small portion of the larger scene (see Fig. 3a). Conceptually, the 2D grid of local latent codes can be interpreted as a latent *floorplan* representation of a scene, where each code is used to locally condition a radiance field network. We opt for a 2D representation produced with a convolutional generator for computational efficiency, but our framework can be easily extended to a 3D grid [30].

### 3.2. Locally Conditioned Radiance Field

To render an image $\hat{\mathbf{X}} \in \mathbb{R}^{3 \times w \times h}$ given a grid of latent codes $\mathbf{W} \in \mathbb{R}^{c \times s \times s}$, and a camera with pose $\mathbf{T} = [\mathbf{R}, \mathbf{t}] \in SE(3)$ and intrinsics $\mathbf{K} \in \mathbb{R}^{3 \times 3}$, we model scenes with a locally conditioned radiance field [34] $f : \mathbb{R}^5 \times \mathbb{R}^c \rightarrow [0,1] \times \mathbb{R}^3$. To predict radiance for a 3D point $\mathbf{p}$ and 2D camera direction $\mathbf{d}$, $f$ is locally conditioned on a vector $\mathbf{w}_{ij} \in \mathbb{R}^c$ sampled at a discrete location $(i, j)$ from the grid of latent codes $\mathbf{W}$ using bi-linear sampling [19]. The location $(i, j)$ is naturally given by the projection of $\mathbf{p}$ on the $zx$-plane (assuming the $y$-axis points up).

Radiance fields are usually defined over $\mathbb{R}^3$ considering a global coordinate system [34, 49, 3] (*i.e.* a coordinate system that spans the whole scene/object). This has been successful since these works do not learn a decomposition of scenes into multiple radiance fields. To decompose the scene into a grid of *independent* radiance fields, we enable our model to perform spatial sharing of local latent codes; the same latent code $w_{ij}$ can be used to represent the same scene part irrespective of its position $(i, j)$ on the grid (see Fig. 6 for empirical results). In order to achieve this spatial sharing, it is not sufficient to learn a grid of local latent codes, it is also necessary to decompose the global coordinate system into multiple local coordinate systems (one for each local latent $\mathbf{w}_{ij}$). The reason for this is to prevent $f$ from relying on the absolute coordinate value of input points to predict radiance. As a result, the input to $f$ is $\mathbf{p}'$, which results from expressing $\mathbf{p}$ relative to the local coordinate system of its corresponding local latent $\mathbf{w}_{ij}$ and applying positional encoding [34].

Finally, the locally conditioned radiance field $f$ outputs two variables: the occupancy $\sigma \in [0, \infty]$ and the appearance $\mathbf{a} \in \mathbb{R}^3$ (see Fig. 3b). To render $\hat{\mathbf{X}}$ given $f$ we leverage Eq. 1 for implicit volumetric rendering [34] and densely evaluate $f$ on points uniformly sampled along ray directions $\mathbf{r}$, where each pixel in $\hat{\mathbf{X}}$ corresponds to a ray and the color $\mathbf{c} \in \mathbb{R}^3$ of pixel/ray $\mathbf{r}$ is obtained by approximating the integral in Eq. 1. For computational efficiency, we use Eq. 1 to render a feature map as in [38]. The rendered feature map is then is upsampled with a small convolutional refinement network to achieve the desired output resolution, resulting in final rendered output $\hat{\mathbf{X}}$.

$$\mathbf{c}(\mathbf{r}, \mathbf{W}) = \int_{u_n}^{u_f} Tr(u)\sigma\left(\mathbf{r}(u), \mathbf{w}_{ij}\right) \mathbf{a}\left(\mathbf{r}(u), \mathbf{d}, \mathbf{w}_{ij}\right) du$$

$$Tr(u) = \exp\left(-\int_{u_n}^{u} \sigma(\mathbf{r}(u), \mathbf{w}_{ij}) du\right) \quad (1)$$

### 3.3. Sampling Camera Poses

Unlike standard 2D generative models which have no explicit concept of viewpoint, radiance fields require a camera pose for rendering an image. Therefore, camera poses $\mathbf{T} = [\mathbf{R}|\mathbf{t}] \in SE(3)$ need to be sampled from pose distribution $p_T$ in addition to the latent code $\mathbf{z} \sim p_z$, which is challenging for the case of realistic scenes and a freely moving camera. GRAF [49] and $\pi$-GAN [3] avoid this issue by training on datasets containing objects placed at the origin, where the camera is constrained to move on a viewing sphere around the object and oriented towards the origin.[1] For GSN, sampling camera poses becomes more challenging due to i) the absence of such constraints (*i.e.* need to

---

[1] Camera poses $\mathbf{T}$ on a sphere oriented towards the origin are constrained to $SO(3)$ as opposed to free cameras in $SE(3)$.

sample $\mathbf{T} \in SE(3)$), ii) and the possibility of sampling invalid locations, such as inside walls or other solid objects that sporadically populate the scene.

To overcome the issue of sampling invalid locations we perform stochastic weighted sampling over a an empirical pose distribution $p_T$ composed by a set of candidate poses, where each pose is weighted by the occupancy (i.e., the $\sigma$ value predicted by the model) at that location. When sampling the candidate poses, we query the generator at each candidate location to retrieve the corresponding occupancy. Then the occupancy values are normalized with a softmin to produce sampling weights for a multinomial distribution. This sampling strategy reduces the likelihood of sampling invalid camera locations while retaining stochasticity required for scene exploration and sample diversity.

### 3.4. Discriminator

Our model adopts the convolutional discriminator architecture from StyleGAN2 [25].The discriminator takes as input an image, concatenated with corresponding depth map normalized between $[0, 1]$ and predicts whether the input comes from the true distribution or the generated distribution. We found it critical to add a small decoder $C$ to the discriminator and to enforce a reconstruction penalty on real images, similar to the self-supervised discriminator proposed by Lui et al. [29]. The additional regularization term restricts the discriminator from overfitting, while encouraging it to learn relevant features about the input that will provide a useful training signal for the generator.

### 3.5. Training

Let $\mathcal{X} = \{\mathbf{X}\} : \mathbf{X} \in \mathbb{R}^{4 \times w \times h}$ denote the set of i.i.d RGB-D samples obtained by recording camera trajectories on the true distribution of scenes $p_S$. Generator $G$, discriminator $D$, and decoder $C$ are parameterized by $\theta$, $\phi$, and $\omega$, respectively. We train our model by optimizing the non-saturating GAN loss [14] with R1 gradient penalty [32], as well as the discriminator reconstruction objective [29]:

$$
\begin{aligned}
\mathcal{L}(\theta, \phi, \omega) = \; & \mathbf{E}_{\mathbf{z} \sim p_z, \mathbf{T} \sim p_T}[h(D_\phi(G_\theta(\mathbf{z}, \mathbf{T})))] \\
& + \mathbf{E}_{\mathbf{X} \sim p_S}[h(-D_\phi(\mathbf{X})) + \lambda_{R1}|\nabla D_\phi(\mathbf{X})|^2 \\
& + \lambda_{Recon}|\mathbf{X} - C_\omega(D_\phi(\mathbf{X}))|^2], \\
& \text{where} \quad h(u) = -\log(1 + \exp(-u)).
\end{aligned}
\tag{2}
$$

## 4. Experiments

In this section we report both quantitative and qualitative experimental results on different settings. First, we compare the generative performance of GSN with recent state-of-the-art approaches. Second, we provide extensive ablation experiments that show the quantitative improvement obtained by the individual components of our model. Finally, we report results on view synthesis by inverting GSN and query-

ing the model to predict views of a scene given a set of input observations.

### 4.1. Generation Performance

We evaluate the generative performance of our model on three datasets: i) the VizDoom environment [26], a *synthetic* computer-generated world, ii) the Replica dataset [52] containing 18 scans of *realistic* scenes that we render using the Habitat environment [48], and iii) the Active Vision Dataset (AVD) [1] consisting of 20k images with noisy depth measurements from 9 *real world* scenes.[2] Images are resized to $64 \times 64$ resolution for all generation experiments. To generate data for the VizDoom and Replica experiments we render 100 sequences of 100 steps each, where an interactive agent explores a scene collecting RGB and depth observations as well as camera poses. Sequences for AVD are defined by an agent performing a random walk through the scenes according to rules defined in [41]. At training time we sample sub-sequences and express camera poses relative the middle frame of the sub-sequence. This normalization enforces an egocentric coordinate system whose origin is placed at the center of $\mathbf{W}$ (see supplementary material for details).

We compare GSN to two recent approaches for generative modelling of radiance fields: GRAF [49] and $\pi$-GAN [3]. For fair comparison all models use the same base architecture and training pipeline, with two main differences between models. The first difference is that GSN uses locally conditioned radiance fields, while GRAF and $\pi$-GAN use global conditioning. The second difference is the type of layer used in the radiance field generator network: GRAF utilizes linear layers, $\pi$-GAN employs modulated sine layers, and GSN uses modulated linear layers. Quantitative performance is evaluated with the Fréchet Inception Distance (FID) [17] and SwAV-FID [35] metrics, which measure the distance between the distributions of real and generated images in pretrained image embedding spaces. We sample 5,000 real and 5,000 generated images when calculating either of these metrics.

Tab. 2 shows the results of our generative modelling comparison. Despite our best attempts to tune $\pi$-GAN's hyperparameters, we find that it struggles to learn detailed depth in this setting, which significantly impedes render quality and leads to early collapse of the model. GSN achieves much better performance than GRAF or $\pi$-GAN across all datasets, obtaining an improvement of 10-14 absolute points on FID. We attribute GSN's higher performance to the increased expressiveness afforded by the locally conditioned radiance field, and not the specific layer type used (compare GRAF in Tab. 2 to GSN + global latents in Tab. 3). As a measure of qualitative results we show scenes sampled from GSN trained on the Replica dataset

---

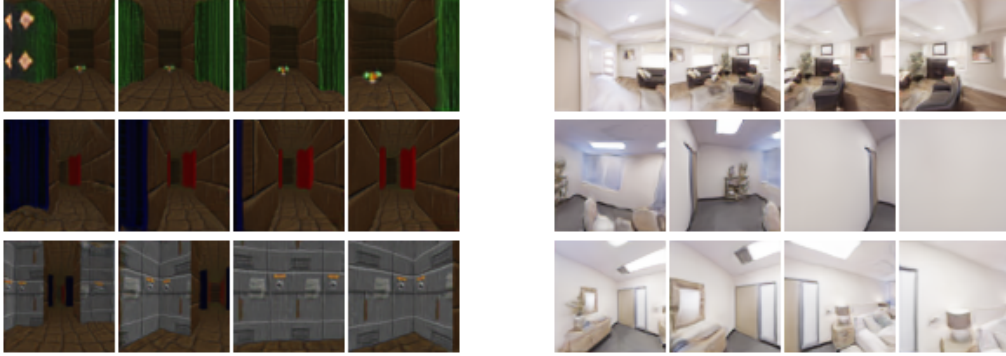[2]We will release code for reproducibility.

Figure 4: Random trajectories through scenes generated by GSN. Models are trained on VizDoom [26] (left), Replica [52] (right) at $64 \times 64$ resolution. We omit qualitative results for AVD [1] due to unclear licensing terms regarding reproduction of figures for this dataset.

| | VizDoom [26] | | Replica [52] | | AVD [1] | |
|---|---|---|---|---|---|---|
| | FID ↓ | SwAV-FID ↓ | FID ↓ | SwAV-FID ↓ | FID ↓ | SwAV-FID ↓ |
| GRAF [49] | $47.50 \pm 2.13$ | $5.44 \pm 0.43$ | $65.37 \pm 1.64$ | $5.76 \pm 0.14$ | $62.59 \pm 2.41$ | $6.95 \pm 0.15$ |
| $\pi$-GAN[3] | $143.55 \pm 4.81$ | $15.26 \pm 0.15$ | $166.55 \pm 3.61$ | $13.17 \pm 0.20$ | $98.76 \pm 1.49$ | $9.54 \pm 0.29$ |
| GSN (Ours) | $\mathbf{37.21 \pm 1.17}$ | $\mathbf{4.56 \pm 0.19}$ | $\mathbf{41.75 \pm 1.33}$ | $\mathbf{4.14 \pm 0.02}$ | $\mathbf{51.11 \pm 1.37}$ | $\mathbf{6.59 \pm 0.03}$ |

Table 2: Generative performance of state-of-the-art approaches for generative modelling of radiance fields on 3 scene-level datasets: Vizdoom [26], Replica [52] and Active Vision (AVD) [1], according to FID [17] and SwAV-FID [35] metrics.

at $128 \times 128$ resolution in Fig. 1, and on the VizDoom, Replica, and AVD datasets at $64 \times 64$ resolution in Fig. 4.

**Latent Space Interpolation** To confirm that GSN is learning a useful scene prior we demonstrate in Fig. 5 some examples of interpolation in the global latent space. The model aligns both geometry and appearance features such that traversing the latent space transitions smoothly between scenes without producing unrealistic off-manifold samples.

## 4.2. Ablation

To further analyze our contributions, we report extensive ablation experiments where we show how our design choices affect the generative performance of GSN. For ablation, we perform our analysis on the Replica dataset [52] at $64 \times 64$ resolution and analyze the following factors: *(i)* the choice of latent code representation (global vs local); *(ii)* the effects of global and local coordinate systems on the learned representation; *(iii)* the effect of sampled trajectory length; *(iv)* the depth resolution needed in order to successfully learn a scene prior; *(v)* the regularization applied to the discriminator.

**How useful are locally conditioned radiance fields?** We compare the effect of replacing our local latent codes $\mathbf{W} \in \mathbb{R}^{c \times s \times s}$ with a global latent code $\mathbf{w} \in \mathbb{R}^c$ that globally conditions the radiance field $f_{\theta_f}$. Generative performance can be improved by decomposing the scene into many inde-

pendent locally conditioned radiance fields (Tab. 3).

| Model | FID ↓ | SwAV-FID ↓ |
|---|---|---|
| GSN + global latents | $68.42 \pm 3.88$ | $6.14 \pm 0.24$ |
| GSN + local latents | $41.75 \pm 1.33$ | $4.14 \pm 0.02$ |

Table 3: Decomposition of the global latent code into a local latent grid has a drastic impact in performance.

**What are the benefits of a local coordinate system?** Enabling feature sharing in local coordinate systems to (*e.g.* a latent $w_{ij}$ can be used to represent the same scene part irrespective of its position $(i, j)$ on the grid). To empirically validate this hypothesis we train GSN with both *global* and *local* coordinate systems. We then sample from the prior obtaining 5k latent grids $\mathbf{W}$ for each model. Next, we perform a simple rigid re-arrangement of latent codes in $\mathbf{W}$ by applying a 2D rotation at different angles and measuring FID by using the resulting grids to predict a radiance field. A local coordinate system is significantly more robust to rearranging local latent codes than a global one (Fig. 6; see supplementary material for qualitative results).

**How long do camera trajectories need to be?** The length of trajectories that collect the camera poses affects the representation of large scenes. Given that we normal-

Figure 5: Two example latent interpolations between global latent codes **z**. Scenes transition smoothly by aligning geometry features such as walls (top) and appearance features such as the picture frame and doorway (bottom). Views are rendered from a fixed camera pose.
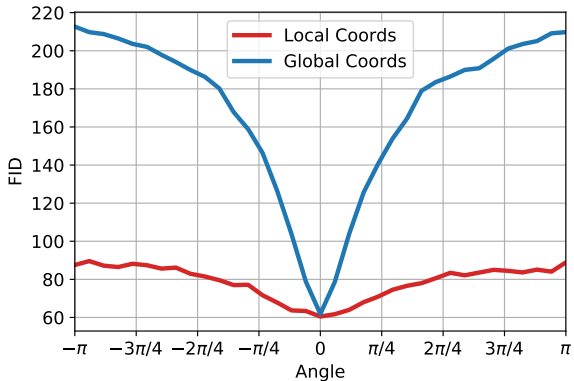


Figure 6: Robustness of generation quality w.r.t. to rotations of the local latent grid **W**. Using a local coordinate system for each local radiance field limits degradation due to rigid re-arrangement of the local latent codes.
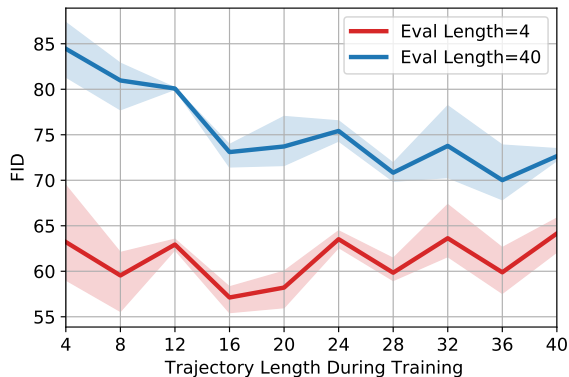


Figure 7: Effect of trajectory length (x axis) on FID scores (y axis) when evaluating short 4-step trajectories (red) vs. 40-step trajectories (blue) during training. Models trained only with short trajectories achieve good local render quality, but cannot move far from the origin without encountering quality degradation.

ize camera poses w.r.t. the middle step in a trajectory, if the trajectories are too short the model will only explore small local neighbourhoods and will not be able to deal with long displacements in camera poses during test time. Models trained on short trajectories fail to generalize to long trajectories during evaluation (Fig. 7). However, models trained with long trajectories do not struggle when evaluated on short trajectories, as evidenced by the stability of the FID when evaluated on short trajectories.

**How much depth information is required during training?** The amount of depth information used during training affects GSN. To test the sensitivity of GSN to depth information, we down-sample both real and generated depth maps to the target resolution, then up-sample back to the original resolution before concatenating with the corresponding image and feeding them to the discriminator. Without depth information GSN fails to train (Tab. 4).

However, we demonstrate that the depth resolution can be reduced to *a single pixel* without finding statistically significant reduction in the generated image quality. This is a significant result, as it enables GSN to be trained in settings with sparse depth information in future work. We speculate that in early stages of training the added depth information guides the model to learn some aspect of depth, after which it can learn without additional supervision.

**Does $D_\phi$ need to be regularized?** Different forms of regularizing the discriminator affect the quality of generated samples. We discover that $D_\phi$ greatly benefits from regularization (Tab. 5). In particular, data augmentation [61] and the discriminator reconstruction penalty [29] are both crucial; training to rapidly diverge without either of these components. The R1 gradient penalty [32] offers an addi-

| Depth Resolution | FID ↓ | SwAV-FID ↓ |
|---|---|---|
| No depth | $311.35 \pm 24.45$ | $28.75 \pm 1.68$ |
| 1 | $41.57 \pm 2.05$ | $4.18 \pm 0.03$ |
| 2 | $40.64 \pm 0.28$ | $4.07 \pm 0.09$ |
| 4 | $44.18 \pm 2.87$ | $4.44 \pm 0.21$ |
| 8 | $42.51 \pm 1.18$ | $4.28 \pm 0.13$ |
| 16 | $42.75 \pm 2.34$ | $4.22 \pm 0.17$ |
| 32 (original) | $41.75 \pm 1.33$ | $4.14 \pm 0.02$ |

Table 4: Generative performance of GSN models trained with different depth resolutions. Remarkably, we show that we can down-sample depth resolution to *a single pixel* without degrading the quality of the generative model.

tional improvement of training stability helping adversarial learning to converge.

| Model | FID ↓ | SwAV-FID ↓ |
|---|---|---|
| GSN (full model) | $41.75 \pm 1.33$ | $4.14 \pm 0.02$ |
| – R1 gradient penalty | $52.9 \pm 19.9$ | $4.56 \pm 1.50$ |
| – Reconstruction penalty | $274.3 \pm 41.4$ | $23.58 \pm 4.11$ |
| – Data augmentation | $412.25 \pm 14.34$ | $35.71 \pm 9.51$ |

Table 5: Regularizing $D_\phi$ shows great benefits for GSN, especially when using a reconstruction penalty and data augmentation.

### 4.3. View Synthesis

We now turn to the task of view synthesis, where we show how GSN performs in comparison with two approaches for scene-level view synthesis of free moving cameras: Generative Temporal models with Spatial Memory (GTM-SM) [13] and Incremental Scene Synthesis (ISS) [41]. Taking the definition in [41], the view synthesis task is defined as follows: for a given step $t$ in a sequence we want the model to predict the target $t + 5$ views $\mathcal{T} = \{(\mathbf{X}, \mathbf{T})_i\}_{i=t:t+5}$ conditioned on the source $t - 5$ views $\mathcal{S} = \{(\mathbf{X}, \mathbf{T})_i\}_{i=t-5:t}$ along the camera trajectory. Note that this view synthesis problem is unrelated to video prediction, since the scenes are static and camera poses for source and target views are given. To tackle this task both GTM-SM [13] and ISS [41] rely on auto-encoder architectures augmented with memory mechanisms that directly adapt to the conditional nature of the task. For GSN to deal with this task we follow standard practices for GAN inversion [57] (see supplementary material for details on the encoder architecture and inversion algorithm). We invert source views $\mathcal{S}$ into the prior and use the resulting latent to locally condition the radiance field and render observations using the camera poses of the target views $\mathcal{T}$.

|  | Memorization | | Hallucination | |
|---|---|---|---|---|
|  | L1 ↓ | SSIM ↑ | L1 ↓ | SSIM ↑ |
| GTM-SM[13] | 0.09 | 0.52 | 0.13 | 0.49 |
| ISS [41] | 0.09 | 0.56 | **0.11** | 0.54 |
| GSN | **0.07** | **0.66** | **0.11** | **0.57** |

(a) View synthesis results on Vizdoom [26]

|  | Memorization | | Hallucination | |
|---|---|---|---|---|
|  | L1 ↓ | SSIM ↑ | L1 ↓ | SSIM ↑ |
| GTM-SM[13] | 0.37 | 0.12 | 0.43 | 0.1 |
| ISS [41] | 0.22 | 0.31 | **0.25** | 0.23 |
| GSN | **0.19** | **0.54** | 0.35 | **0.35** |

(b) View synthesis results on AVD[1]

Table 6: Results of view synthesis on Vizdoom [26] (a) and AVD [1] (b). GSN improves view synthesis quality as a result of modeling a powerful prior over scenes.

Following [41] we report results on the Vizdoom [26] and AVD [1] datasets in Tab. 6.[3] We report two different aspects of view synthesis: the capacity to fit the source views or *Memorization* (*e.g.* reconstruction), and the ability to predict the target views or *Hallucination* (*e.g.* scene completion) using L1 and SSIM metrics.[4] GSN outperforms both GTM-SM [13] and ISS [41] for nearly all tasks (Tab. 6), even though it was not trained to explicitly learn a mapping from $\mathcal{S}$ to $\mathcal{T}$ (see supplementary material for qualitative results). We attribute this success to the powerful scene prior learned by GSN.

### 5. Conclusions

We have made great strides towards generative modeling of unconstrained, complex and realistic indoor scenes. We introduced GSN, a generative model that decomposes a scene into many local radiance fields, which can be rendered by a free moving camera. Our decomposition scheme scales efficiently to big scenes while preserving details and distribution coverage. We show that GSN can be trained on multiple scenes to learn a rich scene prior, while rendered trajectories on scenes sampled from the prior are smooth and consistent, maintaining scene coherence. The prior distribution learned by GSN can also be used to infer observations from arbitrary cameras given a sparse set of observations. A multitude of avenues for future work are

---

[3]There is no code or data release for [41]. In private communication with the authors of [41] we discussed the data splits and settings used for generating trajectories and follow them as closely as possible. We thank them for their help.

[4]We note that while these metrics are a good proxy for reconstruction quality, they do not tell a complete picture in the hallucination setting due to the stochastic nature of the view synthesis problem.

now enabled by GSN, from improving the rendering performance and training on large-scale datasets to exploring the wide range of down-stream tasks that benefit from a learned scene prior like model-based reinforcement learning [27], SLAM [6] or 3D completion for immersive photography [40].

# References

[1] Phil Ammirato, Patrick Poirson, Eunbyung Park, Jana Košecká, and Alexander C Berg. A dataset for developing and benchmarking active vision. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1378–1385. IEEE, 2017. 5, 6, 8

[2] Ivan Anokhin, Kirill Demochkin, Taras Khakhulin, Gleb Sterkin, Victor Lempitsky, and Denis Korzhenkov. Image generators with conditionally-independent pixel synthesis. *arXiv preprint arXiv:2011.13775*, 2020. 4, 12

[3] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-GAN: Periodic implicit generative adversarial networks for 3d-aware image synthesis. *arXiv preprint arXiv:2012.00926*, 2020. 2, 3, 4, 5, 6

[4] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 2

[5] Devendra Singh Chaplot, Dhiraj Prakashchand Gandhi, Abhinav Gupta, and Russ R Salakhutdinov. Object goal navigation using goal-oriented semantic exploration. *Advances in Neural Information Processing Systems*, 33, 2020. 2

[6] Devendra Singh Chaplot, Ruslan Salakhutdinov, Abhinav Gupta, and Saurabh Gupta. Neural topological slam for visual navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12875–12884, 2020. 2, 9

[7] Xu Chen, Jie Song, and Otmar Hilliges. Monocular neural image based rendering with continuous view control. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4090–4100, 2019. 2

[8] Inchang Choi, Orazio Gallo, Alejandro Troccoli, Min H Kim, and Jan Kautz. Extreme view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7781–7790, 2019. 1

[9] Peng Dai, Yinda Zhang, Zhuwen Li, Shuaicheng Liu, and Bing Zeng. Neural point cloud rendering via multi-plane projection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7830–7839, 2020. 1

[10] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017. 14

[11] Emilien Dupont, Miguel Angel Bautista, Alex Colburn, Aditya Sankar, Carlos Guestrin, Josh Susskind, and Qi Shan. Equivariant neural rendering. *arXiv preprint arXiv:2006.07630*, 2020. 1, 2

[12] SM Ali Eslami, Danilo Jimenez Rezende, Frederic Besse, Fabio Viola, Ari S Morcos, Marta Garnelo, Avraham Ruderman, Andrei A Rusu, Ivo Danihelka, Karol Gregor, et al. Neural scene representation and rendering. *Science*, 360(6394):1204–1210, 2018. 1, 3

[13] Marco Fraccaro, Danilo Jimenez Rezende, Yori Zwols, Alexander Pritzel, SM Eslami, and Fabio Viola. Generative temporal models with spatial memory for partially observed environments. *arXiv preprint arXiv:1804.09401*, 2018. 2, 3, 8

[14] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *arXiv preprint arXiv:1406.2661*, 2014. 3, 5

[15] David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018. 2

[16] Philipp Henzler, Niloy J Mitra, and Tobias Ritschel. Escaping plato's cave: 3d shape from adversarial rendering. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9984–9993, 2019. 2

[17] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *arXiv preprint arXiv:1706.08500*, 2017. 5, 6

[18] Geoffrey Hinton, Nitish Srivastava, and Kevin Swersky. Neural networks for machine learning lecture 6a overview of mini-batch gradient descent. 2012. 13

[19] Max Jaderberg, Karen Simonyan, Andrew Zisserman, and Koray Kavukcuoglu. Spatial transformer networks. *arXiv preprint arXiv:1506.02025*, 2015. 4

[20] Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2901–2910, 2017. 3

[21] James T Kajiya and Brian P Von Herzen. Ray tracing volume densities. *ACM SIGGRAPH computer graphics*, 18(3):165–174, 1984. 2

[22] Abhishek Kar, Christian Häne, and Jitendra Malik. Learning a multi-view stereo machine. *arXiv preprint arXiv:1708.05375*, 2017. 14

[23] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017. 12, 14

[24] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019. 2, 12, 13

[25] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8110–8119, 2020. 4, 5, 12, 14

[26] Michał Kempka, Marek Wydmuch, Grzegorz Runc, Jakub Toczek, and Wojciech Jaśkowski. Vizdoom: A doom-based

ai research platform for visual reinforcement learning. In *2016 IEEE Conference on Computational Intelligence and Games (CIG)*, pages 1–8. IEEE, 2016. 5, 6, 8

[27] Seung Wook Kim, Yuhao Zhou, Jonah Philion, Antonio Torralba, and Sanja Fidler. Learning to simulate dynamic environments with gamegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1231–1240, 2020. 9

[28] Andrew Liu, Richard Tucker, Varun Jampani, Ameesh Makadia, Noah Snavely, and Angjoo Kanazawa. Infinite nature: Perpetual view generation of natural scenes from a single image. *arXiv preprint arXiv:2012.09855*, 2020. 3

[29] Bingchen Liu, Yizhe Zhu, Kunpeng Song, and Ahmed Elgammal. Towards faster and stabilized GAN training for high-fidelity few-shot image synthesis. *arXiv e-prints*, pages arXiv–2101, 2021. 5, 7

[30] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *arXiv preprint arXiv:2007.11571*, 2020. 2, 4

[31] Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml*, volume 30, page 3. Citeseer, 2013. 12

[32] Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for GANs do actually converge? In *International conference on machine learning*, pages 3481–3490. PMLR, 2018. 5, 7

[33] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 38(4):1–14, 2019. 1

[34] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *arXiv preprint arXiv:2003.08934*, 2020. 1, 2, 4, 12

[35] Stanislav Morozov, Andrey Voynov, and Artem Babenko. On self-supervised image representations for GAN evaluation. In *International Conference on Learning Representations*, 2021. 5, 6

[36] Zak Murez, Tarrence van As, James Bartolozzi, Ayan Sinha, Vijay Badrinarayanan, and Andrew Rabinovich. Atlas: End-to-end 3d scene reconstruction from posed images. *arXiv preprint arXiv:2003.10432*, 2020. 14

[37] Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian Richardt, and Yong-Liang Yang. Hologan: Unsupervised learning of 3d representations from natural images. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7588–7597, 2019. 2

[38] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. *arXiv preprint arXiv:2011.12100*, 2020. 2, 3, 4, 12

[39] Kyle Olszewski, Sergey Tulyakov, Oliver Woodford, Hao Li, and Linjie Luo. Transformable bottleneck networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7648–7657, 2019. 1, 2

[40] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo-Martin Brualla. Deformable neural radiance fields. *arXiv preprint arXiv:2011.12948*, 2020. 9, 12

[41] Benjamin Planche, Xuejian Rong, Ziyan Wu, Srikrishna Karanam, Harald Kosch, YingLi Tian, Jan Ernst, and Andreas Hutter. Incremental scene synthesis. In *Advances in Neural Information Processing Systems*, pages 1668–1678, 2019. 2, 3, 5, 8

[42] Muhammad Usman Rafique, Hunter Blanton, Noah Snavely, Nathan Jacobs, and Cornell Tech. Generative appearance flow: A hybrid approach for outdoor view synthesis. In *British Machine Vision Conference (BMVC)*, volume 9, 2020. 2

[43] Daniel Rebain, Wei Jiang, Soroosh Yazdani, Ke Li, Kwang Moo Yi, and Andrea Tagliasacchi. Derf: Decomposed radiance fields. *arXiv preprint arXiv:2011.12490*, 2020. 1, 2, 4

[44] Konstantinos Rematas, Ricardo Martin-Brualla, and Vittorio Ferrari. Sharf: Shape-conditioned radiance fields from a single view. *arXiv preprint arXiv:2102.08860*, 2021. 3

[45] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *European Conference on Computer Vision*, 2020. 2

[46] Gernot Riegler and Vladlen Koltun. Stable view synthesis. *arXiv preprint arXiv:2011.07233*, 2020. 2

[47] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 14

[48] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A platform for embodied ai research. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9339–9347, 2019. 5

[49] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. GRAF: Generative radiance fields for 3d-aware image synthesis. *arXiv preprint arXiv:2007.02442*, 2020. 2, 3, 4, 5, 6

[50] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *Advances in Neural Information Processing Systems*, pages 1121–1132, 2019. 2, 3

[51] Pratul P Srinivasan, Richard Tucker, Jonathan T Barron, Ravi Ramamoorthi, Ren Ng, and Noah Snavely. Pushing the boundaries of view extrapolation with multiplane images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 175–184, 2019. 1

[52] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, et al. The replica dataset: A digital replica of indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019. 5, 6, 14

[53] Shao-Hua Sun, Minyoung Huh, Yuan-Hong Liao, Ning Zhang, and Joseph J Lim. Multi-view to novel view: Synthesizing novel views with self-learned confidence. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 155–171, 2018. 1, 2, 3

[54] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. Multi-view 3d models from single images with a convolutional network. In *European Conference on Computer Vision*, pages 322–337. Springer, 2016. 1

[55] Alex Trevithick and Bo Yang. Grf: Learning a general radiance field for 3d scene representation and rendering. *arXiv preprint arXiv:2010.04595*, 2020. 3

[56] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. *arXiv preprint arXiv:2102.13090*, 2021. 3

[57] Weihao Xia, Yulun Zhang, Yujiu Yang, Jing-Hao Xue, Bolei Zhou, and Ming-Hsuan Yang. GAN inversion: A survey. *arXiv preprint arXiv:2101.05278*, 2021. 2, 8, 14

[58] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. *arXiv preprint arXiv:2012.02190*, 2020. 2, 3

[59] Richard Zhang. Making convolutional networks shift-invariant again. In *International Conference on Machine Learning*, pages 7324–7334. PMLR, 2019. 12

[60] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 14

[61] Shengyu Zhao, Zhijian Liu, Ji Lin, Jun-Yan Zhu, and Song Han. Differentiable augmentation for data-efficient gan training. *arXiv preprint arXiv:2006.10738*, 2020. 7, 14

[62] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817*, 2018. 1

[63] Tinghui Zhou, Shubham Tulsiani, Weilun Sun, Jitendra Malik, and Alexei A Efros. View synthesis by appearance flow. In *European conference on computer vision*, pages 286–301. Springer, 2016. 1, 2, 3

# A. Model Architectures and Training Details

In this section we summarize the model architectures, hyperparameter settings, and other training details used for producing the GSN models presented in this paper.

## A.1. Mapping Network

The mapping network maps the global latent code $\mathbf{z}$ to an intermediate non-linear latent space [24]. All models in our experiments, including those that do not use the global generator, use a mapping network which consists of a normalization step followed by three linear layers with LeakyReLU activations [31], as shown in Tab. 7.

|  | Activation | Output Shape |
|---|---|---|
| Input $\mathbf{z}$ | – | 128 |
| Normalize | – | 128 |
| Linear | LeakyReLU (0.2) | 128 |
| Linear | LeakyReLU (0.2) | 128 |
| Linear | LeakyReLU (0.2) | 128 |

Table 7: Mapping network architecture.

## A.2. Global Generator

The purpose of the global generator (Tab. 8) is to map from a single global latent code $\mathbf{z}$ to a 2D grid of local latent codes $\mathbf{W}$ which represent the spatial layout of the scene. The global generator is composed of successive modulated convolutional layers [25] that are conditioned on the global latent code. Following StyleGAN [24], the model learns a constant input for the first layer. The first layers in every pair of modulated convolutional layers thereafter upsamples the feature map resolution by $2\times$, which is implemented as a transposed convolution with a stride of 2, followed by bilinear filtering [59].

The output resolution of the global generator (which is also the spatial resolution of $\mathbf{W}$) is a hyperparameter which effectively controls the size of the spatial region represented by each individual local latent code. We set the global generator output resolution to $32 \times 32$ for all experiments.

|  | Activation | Output Shape |
|---|---|---|
| Constant Input | – | $256 \times 4 \times 4$ |
| ModulatedConv ($3 \times 3$) | LeakyReLU (0.2) | $256 \times 8 \times 8$ |
| ModulatedConv ($3 \times 3$) | LeakyReLU (0.2) | $256 \times 8 \times 8$ |
| ModulatedConv ($3 \times 3$) | LeakyReLU (0.2) | $256 \times 16 \times 16$ |
| ModulatedConv ($3 \times 3$) | LeakyReLU (0.2) | $256 \times 16 \times 16$ |
| ModulatedConv ($3 \times 3$) | LeakyReLU (0.2) | $256 \times 32 \times 32$ |
| ModulatedConv ($3 \times 3$) | LeakyReLU (0.2) | $256 \times 32 \times 32$ |
| ModulatedConv ($3 \times 3$) | – | $32 \times 32 \times 32$ |

Table 8: Global generator architecture.

## A.3. Local Generator

The local generator is composed of a locally conditioned radiance field network which maps coordinates and view direction to appearance $a$ and occupancy $\sigma$, a volumetric rendering step which accumulates along sampled rays to convert $a$ and $\sigma$ values to feature vectors, and a refinement network which upsamples feature maps to higher resolution RGB images.

The locally conditioned radiance field network (Fig. 8) mimics the architecture of the the original NeRF network [34]. To condition the network such that it can represent many different radiance fields we swap out the fixed linear layers for modulated linear layers similar to those used in CIPS [2], where each modulated linear layer is conditioned on $\mathbf{w_{ij}}$. Each modulated linear layer has 128 channels.

When performing volumetric rendering we threshold occupancy values $\sigma$ with a softplus as in D-NeRF [40] as opposed to the standard ReLU, as we find it leads to more stable training. For all experiments we sample 64 samples per ray. When generating $64 \times 64$ images we sample the radiance field network to produce feature maps at $32 \times 32$ resolution, and when generating $128 \times 128$ resolution images we sample feature maps at $64 \times 64$ resolution.

Once volumetric rendering has been performed we upsample the resulting feature map with refinement blocks (Fig. 10) until the desired resolution is achieved, then apply a sigmoid to bound the final output, as in GIRAFFE [38]. In general, we found that sampling higher resolution feature maps directly from the radiance field produced higher quality results compared to sampling at low resolution and applying many refinement blocks, but the computational cost is significantly higher.

## A.4. Discriminator

The discriminator is based on the architecture used in StyleGAN2[25] (Tab. 9), including residual blocks (Fig. 10) and minibatch standard deviation layer [23]. When including depth information as input to the discriminator we normalize it to $[0, 1]$. In the case that the radiance field network is sampled at a resolution lower than the final output resolution (such as when using the refinement network), then resulting depth maps will have lower resolution than real examples. To prevent the discriminator from using this difference in detail to differentiate real and fake samples we downsample all real depth maps to match the resolution of the generated depth maps, then upsample them both back to full resolution.

The decoder (Tab. 10) takes as input $4 \times 4$ resolution feature maps from the discriminator (before the minibatch standard deviation layer), and applies successive transposed convolutions with a stride of 2 and bilinear filtering [59] to upsample the input until the original resolution is recovered.
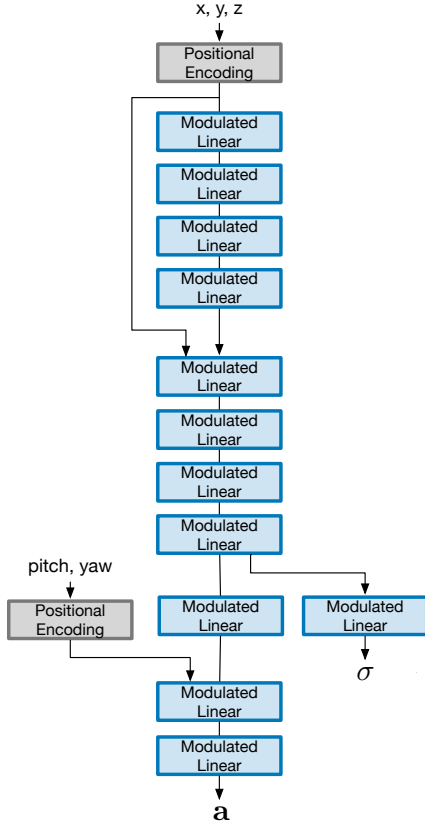
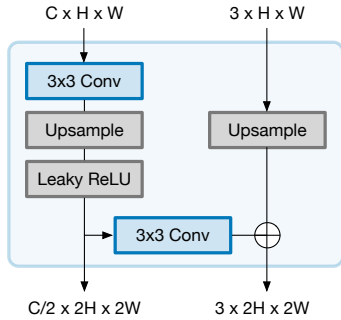Figure 8: Locally conditioned radiance field network.



Figure 9: Refinement block. Feature maps pass through the left path, while RGB images pass through the right path. The input to the right path is not applied for the first refinement block after the volumetric rendering step.



Figure 10: Residual block.

|  | Activation | Output Shape |
|---|---|---|
| Input RGB-D | – | $4 \times 64 \times 64$ |
| Conv $(3 \times 3)$ | LeakyReLU (0.2) | $64 \times 64 \times 64$ |
| Residual Block | LeakyReLU (0.2) | $128 \times 32 \times 32$ |
| Residual Block | LeakyReLU (0.2) | $256 \times 16 \times 16$ |
| Residual Block | LeakyReLU (0.2) | $512 \times 8 \times 8$ |
| Residual Block | LeakyReLU (0.2) | $512 \times 4 \times 4$ |
| Minibatch stdev | – | $513 \times 4 \times 4$ |
| Conv $(3 \times 3)$ | LeakyReLU (0.2) | $512 \times 4 \times 4$ |
| Flatten | – | 8192 |
| Linear | LeakyReLU (0.2) | 512 |
| Linear | – | 1 |

Table 9: Discriminator architecture.

|  | Activation | Output Shape |
|---|---|---|
| Input Feature Map | – | $512 \times 4 \times 4$ |
| ConvTranspose $(3 \times 3)$ | LeakyReLU (0.2) | $256 \times 8 \times 8$ |
| ConvTranspose $(3 \times 3)$ | LeakyReLU (0.2) | $128 \times 16 \times 16$ |
| ConvTranspose $(3 \times 3)$ | LeakyReLU (0.2) | $64 \times 32 \times 32$ |
| ConvTranspose $(3 \times 3)$ | LeakyReLU (0.2) | $32 \times 64 \times 64$ |
| Conv $(3 \times 3)$ | – | $4 \times 64 \times 64$ |

Table 10: Decoder architecture

## A.5. Sampling Camera Poses

We use poses from camera trajectories in the training set as candidate poses during sampling, as real camera poses better reflect the true distribution of occupable locations compared to uniformly sampling over the entire scene region.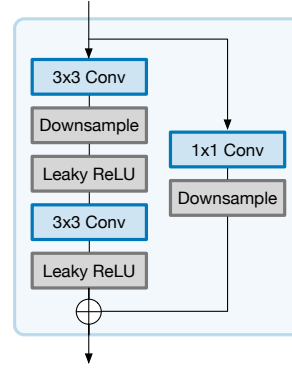 Sampled camera poses are normalized and expressed relative to the camera pose in the middle of the trajectory. This normalization enforces an egocentric coordinate system whose origin is placed at the center of $\mathbf{W}$. Note that despite working with trajectories of multiple camera poses, we still only sample a single camera pose per generated scene during training.

## A.6. Training Details

We use the RMSprop [18] optimizer with $\alpha = 0.99$, $\epsilon = 10^{-8}$, and a learning rate of 0.002 for both the generator and discriminator. Following StyleGAN [24], we set the learning rate of the mapping network $100\times$ less than the rest of the network for improved training stability. Equalized

learning rate [23] is used for all learnable parameters, and an exponential moving average of the generator weights [23] with a decay of 0.999 is used to stabilize test-time performance. Differentiable data augmentations [61] such as random translation, colour jitter, and Cutout [10] are applied to all inputs to the discriminator in order to combat overfitting. To save compute, the R1 gradient penalty is applied using a lazy regularization strategy [25] by applying it every 16 iterations. We set $\lambda_{R1}$ to 0.01 and $\lambda_{Recon}$ to 1000 for all experiments.

All $64 \times 64$ resolution models used for the generation performance evaluation (GSN and otherwise) were trained for 500k iterations with a batch size of 32. Training takes 4 days on two NVIDIA A100 GPUs with 40GB of memory each. Mixed precision training is applied to the generator for a small reduction in memory cost and training time. We do not apply mixed precision training to the discriminator as training stability decreases in this case.

## B. Inverting GSN for View Synthesis

In order for GSN to deal with the view synthesis problem we follow common practices for GAN inversion [57], adopting a hybrid inversion approach where we first train an encoder $E_{\theta_E} : \mathbb{R}^{3 \times w \times h} \times SE(3) \longrightarrow \mathbb{R}^{c \times s \times s \times s}$ on $\{(\hat{\mathbf{X}}, \mathbf{T}, \mathbf{W})_i\}_{i=1:n}$ tuples sampled from a trained GSN (trained on the training set of the same dataset). The goal of this encoder is to predict an initial grid of latent codes $\mathbf{W}_0$ given a set of posed views. Our encoder is conceptually similar to [22, 36] where views $\hat{\mathbf{X}}$ are first processed with a backbone (UNet [47] with a ResNet-50 encoder in our case) and the resulting feature maps are back-projected using camera poses $\mathbf{T}$ into a shared feature volume $\mathbf{V} \in \mathbb{R}^{c \times s \times s \times s}$. Finally, we perform average pooling over the height dimension of $\mathbf{V}$ to get $\mathbf{W}_0 \in \mathbb{R}^{c \times s \times s}$. We train the encoder by minimizing the following reconstruction loss:

$$\mathcal{L}(\hat{\mathbf{X}}, \mathbf{T}, \mathbf{W}; \theta_E) = \|\mathbf{W} - E_{\theta_E}(\hat{\mathbf{X}}, \mathbf{T})\|_2 + \qquad (3)$$
$$+ \|\hat{\mathbf{X}} - f(E_{\theta_E}(\hat{\mathbf{X}}, \mathbf{T}), \mathbf{T})\|_2, \qquad (4)$$

where the first term encourages the reconstruction of the local latent grid, and the second term encourages samples from locally conditioned radiance field $f$ to match the original input views.

At inference time, given a trained encoder $E_{\theta_E}$, we feed the source views $\mathcal{S} = \{(\mathbf{X}, \mathbf{T})_i\}_{i=t-5:t}$ through our encoder to predict an initialization latent code grid $\mathbf{W}_0$ that we then optimize via SGD for 1000 iterations to get $\hat{\mathbf{W}}$. Given that scenes do not share a canonical orientation, we predict $\mathbf{W}_0$ at multiple rotation angles of $\{\mathbf{T}_i\}_{i=t-5:t}$ about the $y-$axis to find the generator's preferred orientation and use this orientation during optimization (note that relative transformations between camera poses do not change with this global

transformation). We define the preferred orientation as the one that minimizes an auto-encoding LPIPS loss [60]. The optimization process is performed by freezing the weights of $f_{\theta_f}$ and computing a reconstruction loss w.r.t. $\mathcal{S}$. We then use $\hat{\mathbf{W}}$ in the locally conditioned radiance field and render observations using the camera poses of $\mathcal{S}$ to produce $\hat{\mathcal{S}}$ (i.e. to auto-encode source views), while also rendering from the camera poses of the target views $\mathcal{T}$ to produce $\hat{\mathcal{T}}$ (i.e. to predict unseen parts of the scene). Future work will explore in depth how to improve the quality and efficiency of the inversion approach for GSN-based models where the generative model tends to prefer a certain orientation. In Fig. 11 we show qualitative results for view synthesis on Vizdoom on held out sequences not seen during training. We can see how GSN learns a robust prior that is able to fill in the scene with plausible completions (e.g. row 5 and 8), even if those completions do not strictly minimize the L1 reconstruction loss.

In addition, Fig. 12 shows qualitative view synthesis results on the Replica dataset [52], showing the applicability of GSN for view synthesis on realistic data. In this experiment we follow the settings described for Vizdoom in terms of $\mathcal{S}$ and $\mathcal{T}$. In Fig. 12 the top 3 rows show results on data from the training set (e.g. scenes that were observed during training) and the bottom 3 rows show test set results (e.g. results on unseen scenes). We can see how GSN successfully uses the prior learned from training data to find a plausible scene completion for unseen scenes that respects the global scene geometry.

## C. Qualitative Results on Local vs. Global Coordinate Systems

In this section we demonstrate the robustness of GSN w.r.t. re-arrangement of the latent codes in $\mathbf{W}$. In order to do so we sample different scenes from our learned prior and apply a rigid transformation to their corresponding $\mathbf{W}$ (a 2D rotation). In principle, this rotation should amount to a rotation of the scene represented by $\mathbf{W}$ that does not change the radiance field prediction. To qualitatively evaluate this effect we sample different scenes and rotate their corresponding $\mathbf{W}$ by $\{0, 90, 180, 270\}$ degrees while (i) rotating the camera by the same amount about the $y-$axis so that the rendered image should remain constant and (ii) keeping the camera fixed so that the scene should rotate. In Fig. 13-14 we show the result of the setting in (i) for a local and global coordinate system respectively. In these results we see how a local coordinate system is drastically more robust to re-arrangements of the latent codes than a global coordinate system. In addition, we show results for the setting in (ii) in Fig. 15-16 for local and global coordinate systems respectively. In this case, we see how given a fixed camera, a rotation of $\mathbf{W}$ amounts to rotating the scene. In this case we
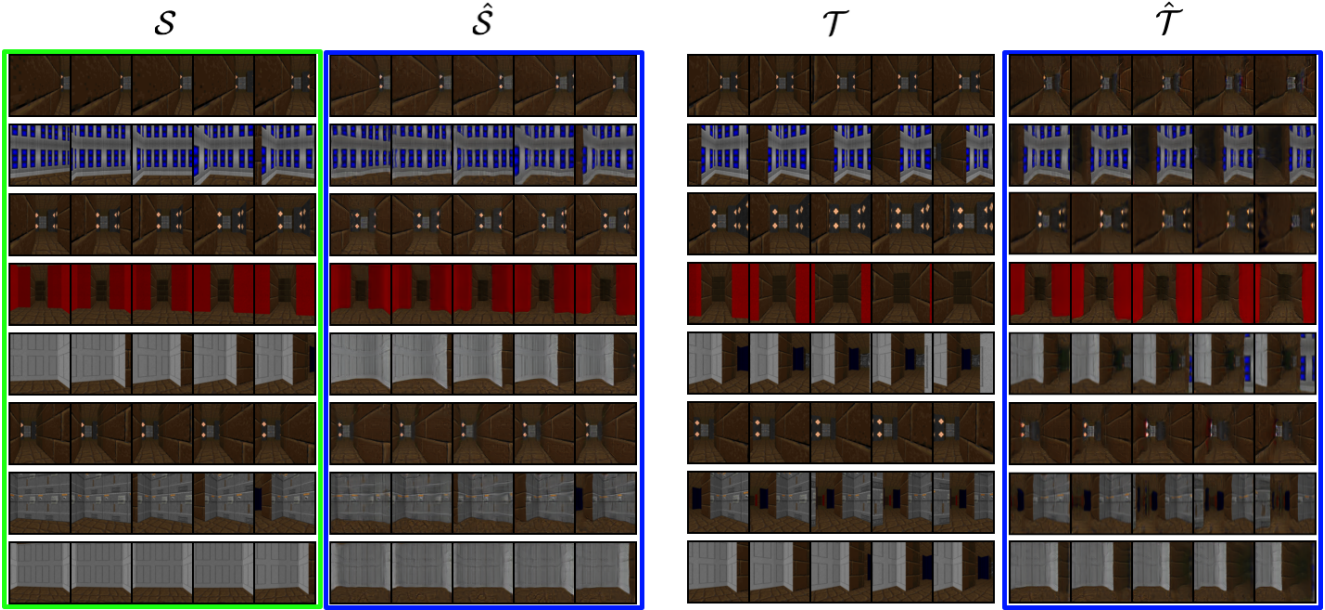
Figure 11: Qualitative view synthesis results on Vizdoom sequences not seen during training. Given source views $\mathcal{S}$ we invert GSN to obtain a local latent code grid $\hat{\mathbf{W}}$, which is then use both to reconstruct $\mathcal{S}$, denoted as $\hat{\mathcal{S}}$, and also to predict target views $\mathcal{T}$ (given their camera poses) which are denoted as $\hat{\mathcal{T}}$. Each row corresponds to a different set of source views $\mathcal{S}$. Frames highlighted in green are input to GSN, frames highlighted in blue are predictions.
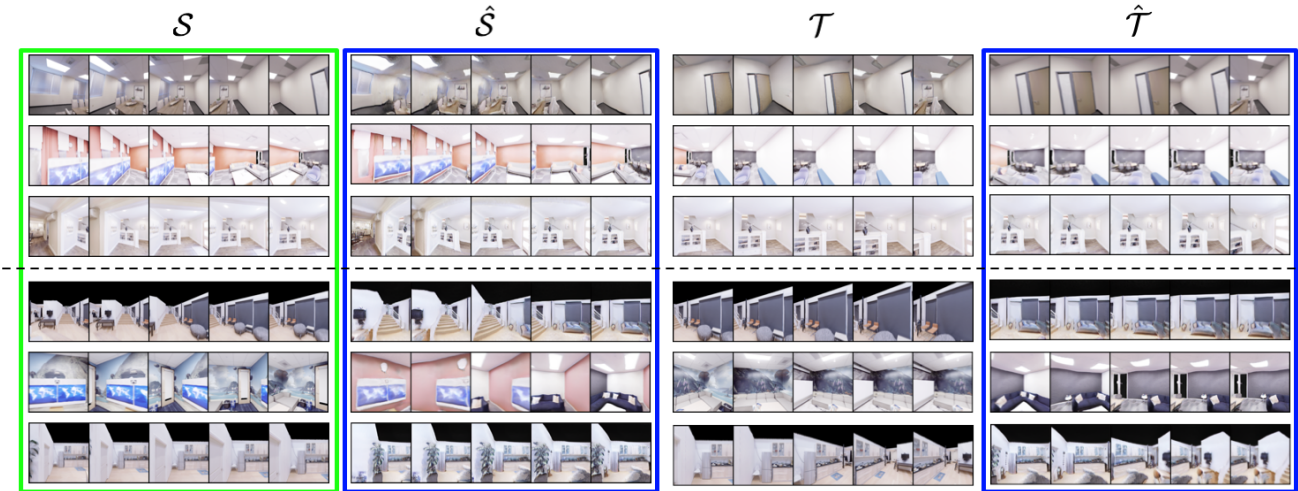


Figure 12: Qualitative view synthesis results on Replica. Given source views $\mathcal{S}$ we invert GSN to obtain a local latent code grid $\hat{\mathbf{W}}$, which is then use both to reconstruct $\mathcal{S}$, denoted as $\hat{\mathcal{S}}$, and also to predict target views $\mathcal{T}$ (given their camera poses) which are denoted as $\hat{\mathcal{T}}$. Each row corresponds to a different set of source views $\mathcal{S}$ (top 3 rows are scenes from the training set, bottom 3 rows are scenes in a heldout test set). Frames highlighted in green are input to GSN, frames highlighted in blue are predictions.

can also see how a local coordinate system results in higher rendering quality compared to that of a global coordinate system, which suffers from degradation as the rotation angle increases.

## D. Scene Editing

A nice property of the local latent grid produced by GSN is that it can be used to perform scene editing by directly altering $\mathbf{W}$. This property allows us a degree of manual control for scene synthesis beyond what we get from ran-
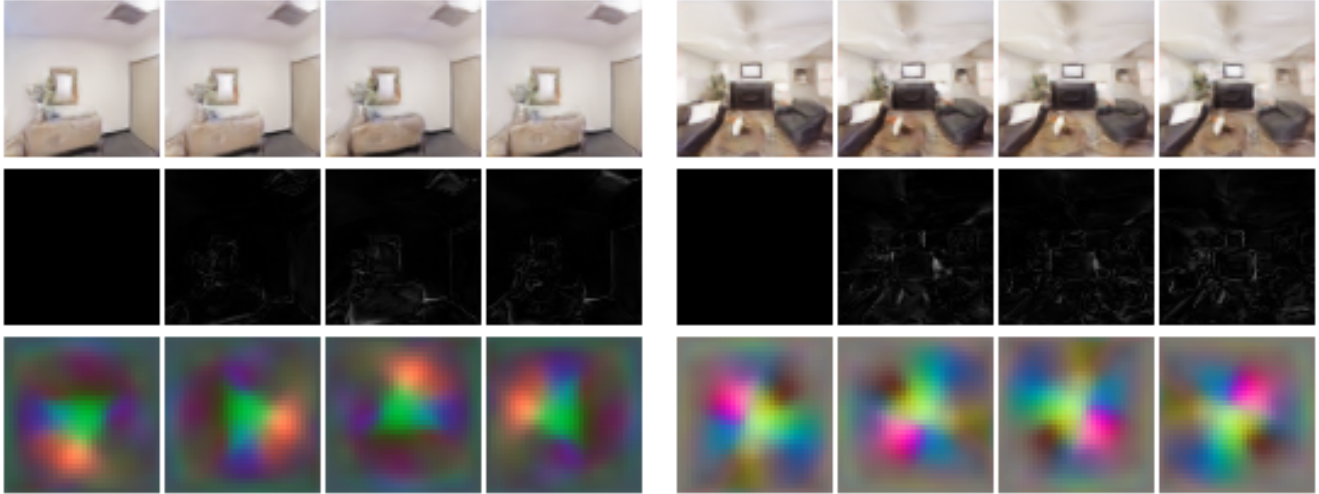
Figure 13: Change in generation output as local latent codes are rotated with a *local* coordinate system for two different scenes. (Top) Rendered image. (Middle) Residual w.r.t. 0 degree rotation. (Bottom) Visualization of **W**. Each column corresponds to a rotation of the camera and **W** in $\{0, 90, 180, 270\}$.
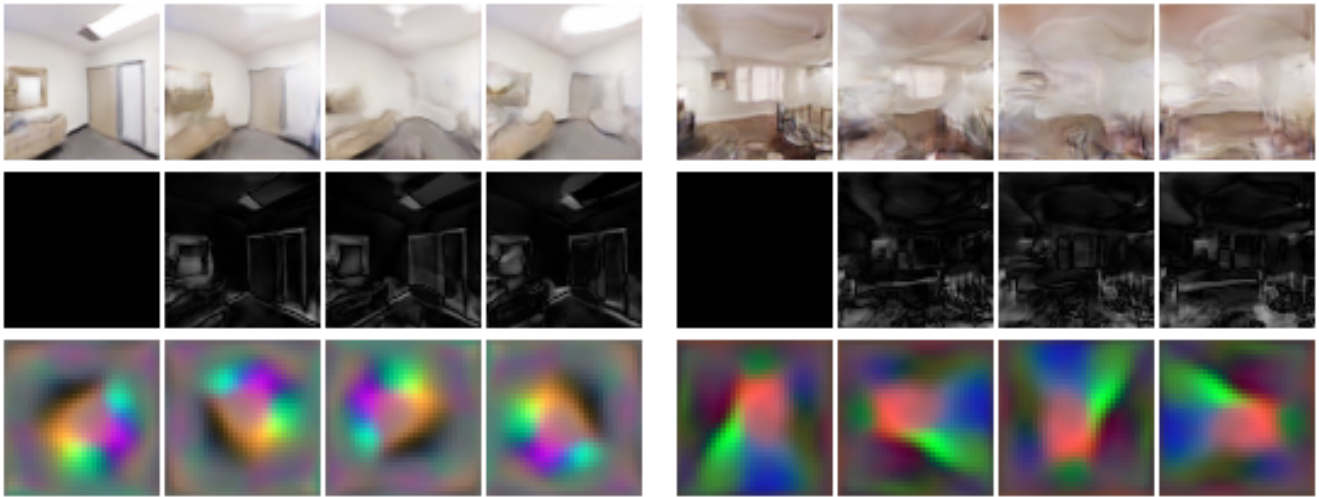


Figure 14: Change in generation output as local latent codes are rotated with a *global* coordinate system. (Top) Rendered image. (Middle) Residual w.r.t. 0 degree rotation. (Bottom) Visualization of **W**. Each column corresponds to a rotation of the camera and **W** in $\{0, 90, 180, 270\}$.

domly sampling the generator. While the low resolution of **W** used in current models currently limits us to high level scene modifications, training with larger local latent grids could allow for more fine-grained control over scenes, such as rearrangement of furniture.

We find that, as with most image composition operations, the results of scene editing appear most convincing when the inputs are well aligned in terms of appearance and geometry. We demonstrate editing operations by manipulating the codes from single scenes, since we don't need to worry about matching appearance and geometry, but multi-

ple scenes could be combined if they were similar enough. In Fig. 17-18 we manipulate the local latent codes by mirroring them along the horizontal axis to produce unique scenes.

Figure 15: Change in generation output for a *fixed camera* as local latent codes are rotated with a *local* coordinate system for two different scenes. (Top) Rendered image. (Bottom) Visualization of **W**. Each column corresponds to a rotation of **W** in $\{0, 90, 180, 270\}$.



Figure 16: Change in generation output for a *fixed camera* as local latent codes are rotated with a *global* coordinate system for two different scenes. (Top) Rendered image. (Bottom) Visualization of **W**. Each column corresponds to a rotation of **W** in $\{0, 90, 180, 270\}$.



Figure 17: Panoramas and corresponding local latent codes for scenes produced by GSN. Mirroring the local latent code from a single room (top row) produces a new room (bottom row).



Figure 18: Panoramas and corresponding local latent codes for scenes produced by GSN. Mirroring the local latent code from a single room (top row) produces a new room (bottom row).