

Persistent Nature: A Generative Model of Unbounded 3D Worlds

Lucy Chai¹, Richard Tucker², Zhengqi Li², Phillip Isola¹, Noah Snavely^{2,3}

¹MIT ²Google Research ³Cornell Tech

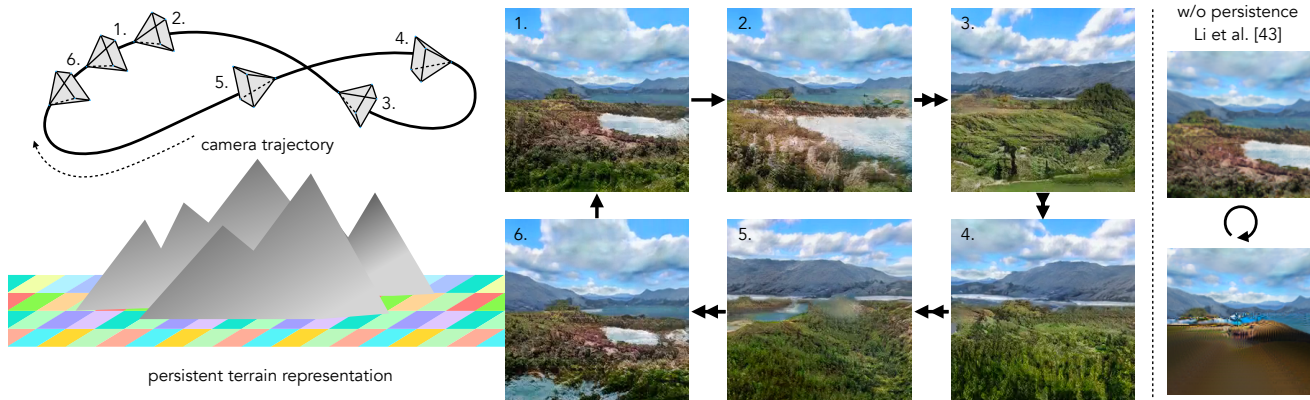


Figure 1. Our approach enables unconditional synthesis of unbounded 3D nature scenes with a persistent scene representation (**left**), using a scene layout grid representing a large-scale terrain model (depicted above as the checkered ground plane). This representation enables us to generate arbitrary camera trajectories, such as the six numbered views shown along a cyclic camera path (**center**). The *persistence* inherent to our representation stands in contrast to prior auto-regressive methods [42] that do not preserve consistency under circular camera trajectories (**right**); while the two images shown on the right are at the start and end of a cyclic path, the terrain depicted is completely different. Our method is trained solely from unposed, single-view landscape photos.

Abstract

Despite increasingly realistic image quality, recent 3D image generative models often operate on 3D volumes of fixed extent with limited camera motions. We investigate the task of unconditionally synthesizing unbounded nature scenes, enabling arbitrarily large camera motion while maintaining a persistent 3D world model. Our scene representation consists of an extendable, planar scene layout grid, which can be rendered from arbitrary camera poses via a 3D decoder and volume rendering, and a panoramic skydome. Based on this representation, we learn a generative world model solely from single-view internet photos. Our method enables simulating long flights through 3D landscapes, while maintaining global scene consistency—for instance, returning to the starting point yields the same view of the scene. Our approach enables scene extrapolation beyond the fixed bounds of current 3D generative models, while also supporting a persistent, camera-independent world representation that stands in contrast to auto-regressive 3D prediction models. Our project page: <https://chai1.github.io/persistent-nature/>.

1. Introduction

Generative image and video models have achieved remarkable levels of realism, but are still far from presenting a convincing, explorable world. Moving a virtual camera through these models—either in their latent space [3, 23, 29, 72] or via explicit conditioning [35]—is not like walking about in the real world. Movement is either very limited (for example, in object-centric models [5]), or else camera motion is unlimited but quickly reveals the lack of a persistent world model. Auto-regressive 3D synthesis methods exemplify this lack of persistence [42, 45]; parts of the scene may change unexpectedly as the camera moves, and you may find that the scene is entirely different when returning to previous positions. The lack of spatial and temporal consistency can give the output of these models a strange, dream-like quality. In contrast, machines that can generate unbounded, persistent 3D worlds could be used to develop agents that plan within a world model [21], or to build virtual reality experiences that feel closer to the natural world, rather than appearing as ephemeral hallucinations [42].

We therefore aim to develop a unconditional generative model capable of generating unbounded 3D scenes with a

persistent underlying world representation. We want synthesized content to move in a way that is consistent with camera motion, yet we should also be able to move arbitrarily far and still generate the same scene upon returning to a previous camera location, regardless of the camera trajectory.

To achieve this goal, we model a 3D world as a *terrain* plus a *skydome*. The terrain is represented by a *scene layout grid*—an extendable 2D array of feature vectors that acts as a map of the landscape. We ‘lift’ these features into 3D and decode them with an MLP into a radiance field for volume rendering. The rendered terrain images are super-resolved and composited with renderings from the skydome model to synthesize final images. We train using a layout grid of limited size, but can extend the scene layout grid by any desired amount during inference, enabling unbounded camera trajectories. Since our underlying representation is persistent over space and time, we can fly around 3D landscapes in a consistent manner. Our method does not require multiview data; each part of our system is trained from an unposed collection of single-view images using GAN objectives.

Our work builds upon two prior threads of research that tackle generating immersive worlds: 1) generative models of 3D data, and 2) generative models of infinite videos. Along the first direction are generators of meshes, volumes, radiance fields, etc (e.g., [5, 54, 59]). These models represent a consistent 3D world by construction, and excel at rendering isolated objects and bounded indoor scenes. Our work, in contrast, tackles the challenging problem of generating large-scale *unbounded* nature scenes. Along the second direction are methods like InfiniteNature [42, 45], which can indeed simulate visual worlds of infinite extent. These methods enable unbounded scene synthesis by predicting new viewpoints auto-regressively from a starting view. However, they do not ensure a persistent world representation; content may change when revisited.

Our method aims to combine the best of both worlds, generating boundless scenes (unlike prior 3D generators) while still representing a persistent 3D world (unlike prior video generative models). In summary:

- We present an unconditional 3D generative model for unbounded nature scenes with a persistent world representation, consisting of a terrain map and skydome.
- We augment our generative pipeline to support camera extrapolation beyond the training camera distribution by extending the terrain features.
- Our model is learned entirely from single-view landscape photos with unknown camera poses.

2. Related Work

Image and view extrapolation. Pioneering work by Kaneva *et al.* [32] proposed the task of infinite image extrapolation by using a large image database to perform classical

2D image retrieval, stitching, and rendering. More recently, various learning-based 2D image inpainting [24, 41, 46, 67, 79, 94, 95, 97] and outpainting [2, 9, 43, 80, 89, 91] methods have been developed. These methods fill in missing image regions or expand the field of view by synthesizing realistic image content that is coherent with the partial input image. Beyond 2D, prior work has explored single-view 3D *view extrapolation*, often by applying 2D image synthesis techniques within a 3D representation [28, 30, 40, 65, 66, 74, 90]. However, these methods can only extrapolate content within a very limited range of viewpoints.

Video generation. Video generation aims to synthesize realistic videos from different types of input. Unconditional video generation produces long videos often from noise input [3, 17, 19, 47, 53, 77, 83], while conditional video generation generates sequences by conditioning on one or a few images [13, 15, 27, 36, 38, 84, 85, 85, 86, 88, 92, 96], or a text prompt [26, 75]. However, applying these ideas in 3D requires supervision from multi-view training data, and cannot achieve persistent 3D scene content at runtime, since there is no explicit 3D representation. Some recent work preserves global scene consistency via extra 3D geometry inputs such as point clouds [49] or voxel grids [22]. In contrast, our method synthesizes both the geometry and appearance of an entire world from scratch using a global feature representation to achieve consistent generated content.

Generative view synthesis. Novel view synthesis aims to produce new views of a scene from single [7, 31, 37, 57, 66, 73, 74, 81, 82, 90, 93] or multiple image observations [1, 10, 16, 39, 48, 50–52, 64, 68, 71, 87, 98] by constructing a local or global 3D scene representation. However, most prior methods can only interpolate or extrapolate a limited distance from the input views, and do not possess a generative ability.

On the other hand, a number of generative view synthesis methods have been recently proposed utilizing neural volumetric representations [5, 14, 20, 54–56, 62, 69, 78]. These methods can learn to generate 3D representations from 2D supervision, and have demonstrated impressive results on generating novel objects [59], faces [5, 12, 20, 58], or indoor environments [14, 63]. However, none of these methods can generate unbounded outdoor scenes due to lack of multi-view data for supervision, and due to the larger and more complex scene geometry and appearance that is difficult to model with prior representations. In contrast, our approach can generate globally consistent, large-scale nature scenes by training solely from unstructured 2D photo collections.

Our work is particularly inspired by recent perpetual view generation methods, including InfiniteNature [45] and InfiniteNature-Zero [42], which can generate unbounded fly-through videos of natural scenes, and are trained on nature videos or photo collections. However, these methods generate video sequences in an auto-regressive manner, and therefore cannot achieve globally consistent 3D scene content.

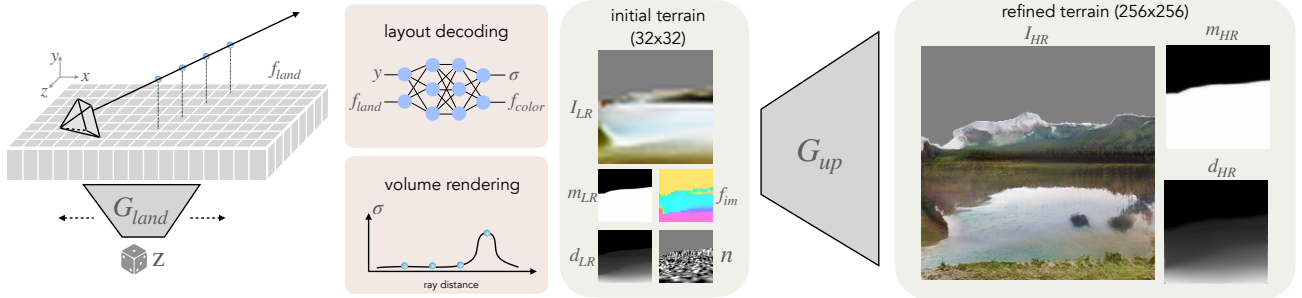


Figure 2. *Overview of scene layout decoding.* The layout generator G_{land} samples a random latent code to produce a 2D scene layout grid f_{land} representing the shape and appearance of a terrain map, and which can be spatially extended using a grid of latent codes (see § 3.2). To render an image from a given camera, sampled points along camera rays passing over the feature plane are decoded via an MLP into a color feature f_{color} and density σ , which are then volume rendered. This produces a low-resolution image, mask, depth, image features, and a projected noise pattern, which are provided to a refinement network G_{up} to produce final image, mask, and depth outputs.

Our approach instead adopts a global scene representation that can be trained to generate consistent-by-construction and realistic novel views spanning large-scale scenes. Concurrent works for scene synthesis InfiniCity [44] and SceneDreamer [8] leverage birds-eye-view representations, while SceneScape [18] builds a mesh representation from text.

3. Method

Our scene representation for unbounded landscapes consists of two components, a *scene layout grid* and a *skydome*. The scene layout grid models the landscape terrain, and is a 2D grid of features defined on a “ground plane.” These 2D features are intended to describe both the height and appearance content of the terrain, representing the full 3D scene — in fact, we decode these features to a 3D radiance field, which can then be rendered to an image (§3.1). To enable camera motion beyond the training volume, we spatially extend the 2D feature grid to arbitrary sizes (§3.2). Because it is computationally expensive to generate and volume render highly detailed 3D content at the scale we aim for, we use an image-space refinement network that adds additional texture detail to rendered images (§3.3).

The second scene component is a *skydome* (§3.4), which is a spherical (panoramic) image intended to model very remote content, such as the sun and sky, as well as distant mountains. The skydome is generated to harmonize with the terrain content described by the scene layout grid.

All the stages of our approach are trained with GAN losses (§3.5). In what follows, we use the 3D coordinate convention that the ground plane is the xz -plane, and the y -axis represents height above or below this plane. Generally, the camera used to view the scene will be positioned some height above the ground.

3.1. Scene layout generation and rendering

To represent a distribution over landscapes, we take a generative approach following the layout representation of

GSN [14]. First, a 2D scene layout grid is synthesized from a sampled random noise code \mathbf{z} passed to a StyleGAN2 [34] generator G_{land} . This creates a 2D feature grid f_{land} , which we bilinearly interpolate to obtain a 2D function over spatial coordinates x and z :

$$f_{\text{land}}(x, z) = \text{Interpolate}(G_{\text{land}}(\mathbf{z}), (x, z)) \quad (1)$$

To define a full 3D scene, we need a way to compute the content at any 3D location (x, y, z) . We define a multi-layer perceptron M that takes a scene grid feature, as well as the height y of the point at which we want to evaluate the scene content. The outputs of M are the 2D-to-3D lifted feature f_{color} and the density σ at point (x, y, z) :

$$f_{\text{color}}, \sigma = M(f_{\text{land}}(x, z), y). \quad (2)$$

In this way, the 2D scene layout grid determines a radiance field over all 3D points within the bounds of the grid [14, 70, 93]. That is, feature vectors in the grid encode not just appearance information, but also the height (or possibly multiple heights) of the terrain at their ground location.

To render an image from a desired camera pose, we cast rays \mathbf{r} from the camera origin through 3D space, sample points (x, y, z) along them, and compute f_{color} and σ at each point. We then use volume rendering to composite f_{color} along each ray into projected 2D image features f_{im} , a disparity image d_{LR} , and a sky segmentation mask m_{LR} . We form an initial RGB image of the terrain, I_{LR} , via a learned linear projection P of these image features. This process is depicted in the left half of Fig. 2, and is defined as:

$$\begin{aligned} f_{\text{im}}(\mathbf{r}) &= \sum_{i=1}^N w_i f_{\text{color}, i}, & d_{\text{LR}}(\mathbf{r}) &= \sum_{i=1}^N w_i d_i, \\ m_{\text{LR}}(\mathbf{r}) &= \sum_{i=1}^N w_i, & I_{\text{LR}} &= P f_{\text{im}}, \end{aligned} \quad (3)$$

where $i \in \{1..N\}$ refers to the index of each sampled point along ray \mathbf{r} in order of increasing distance from the camera,

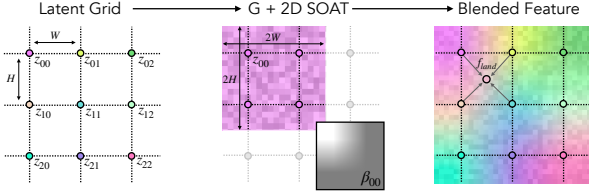


Figure 3. *Layout extension procedure.* To extend the layout at inference time, we sample noise codes \mathbf{z} in a grid arrangement. To smoothly transition between adjacent feature grids, we use the SOAT (StyleGAN of All Trades) procedure [11] in 2D. Operating on a 2×2 sub-grid, we apply each generator layer four times in fully convolutional manner over the entire sub-grid, each time conditioned on a different corner latent code \mathbf{z} , before multiplying by bilinear blending weights. This process is repeated for each layer of the generator and each sub-grid. Each 2×2 sub-grid produces a $2H \times 2W$ feature grid, and sub-grids are blended together in an overlapping fashion to obtain an extended feature grid f_{land} of arbitrary spatial size.

d_i is the inverse-depth (disparity) of point i , and weights w_i are determined from the volume rendering equations used in NeRF [51] (see supplemental).

We intend the mask m_{LR} to distinguish sky regions (which will be empty and filled later using the skydome) from non-sky regions, and achieve this by training using segmented real images in which color and disparity for sky pixels are replaced with zero. Since to achieve zero disparity all weights along a ray must be zero (which also results in a zero-valued color feature), this approach encourages the generator to omit sky content. However, while we find that the model indeed learns to generate transparent sky regions, land geometry can also become partially transparent. To counter this, we penalize visible decreases in opacity along viewing rays using finite differences of opacity α :

$$\mathcal{L}_{\text{transparent}}(\mathbf{r}) = \sum_{i=2}^N w_i \frac{\max(\alpha_{i-1} - \alpha_i, 0)}{\delta_i}. \quad (4)$$

3.2. Layout Extension

While G_{land} creates a fixed-size feature grid, our objective is to generate geometry of arbitrary size, enabling long-distance camera motion at inference time. Hence, we devise a way to *extend* the feature grid in the x and z dimensions. We illustrate this process in Fig. 3, where we first sample noise codes \mathbf{z} in a grid arrangement, where each \mathbf{z} generates a 2D layout feature grid of size $H \times W$. To obtain a smooth transition between these independently sampled layout features, we generalize the image interpolation approach from SOAT (StyleGAN of all Trades) [11] to two dimensions. We operate on 2×2 sub-grids and blend intermediate features

from each layer of the generator as follows:

$$\begin{aligned} f_{k,l+1} &= G_l(f_l, \mathbf{z}_k); \quad k = \{00, 01, 10, 11\} \\ f_{l+1} &= \sum_{k=\{00,01,10,11\}} \beta_k(x, z) f_{k,l+1}. \end{aligned} \quad (5)$$

For each of the four corner anchors k , we construct the modulated feature $f_{k,l+1}$ by applying G_l (the l -th layer of G_{land}) in a fully convolutional manner over the entire sub-grid. We then interpolate between the four feature grids using bilinear interpolation weights $\beta_k(x, z)$. By stitching these 2×2 sub-grids in an overlapping manner, we can obtain a scene layout feature grid of arbitrary size to use as f_{land} . Additional details are provided in the supplemental.

3.3. Image refinement

Due to the computational cost of volume rendering, training the layout generator at higher resolutions becomes impractical. We therefore use a refinement network G_{up} to up-sample the initial generated image I_{LR} to a higher-resolution result I_{HR} , while adding textural details (Fig. 2-right). We use a StyleGAN2 backbone for G_{up} , replacing the earlier feature layers with feature output f_{im} and the RGB residual layers with a concatenation of I_{LR} , d_{LR} , and m_{LR} . To encourage the refined terrain image I_{HR} to be consistent with the sky mask, the network also predicts a refined disparity map and sky mask for compositing with the skydome (see §3.4):

$$I_{\text{HR}}, d_{\text{HR}}, m_{\text{HR}} = G_{\text{up}}(f_{\text{im}}, I_{\text{LR}}, d_{\text{LR}}, m_{\text{LR}}). \quad (6)$$

We compute a reconstruction loss between the initial and refined disparity and mask outputs, and penalize G_{up} for producing gray sky pixels in I_{HR} outside the predicted mask m_{HR} . Please see the supplemental for more details.

For fine texture details, StyleGAN2 also uses layer-wise spatial noise in intermediate generator layers (in addition to the global latent \mathbf{z}). Using a fixed 2D noise pattern results in texture ‘sticking’ as we move the camera [33], but resampling it every frame reduces spatial coherence and removing it entirely results in convolutional gridding artifacts. To avoid these issues and improve spatial consistency, we replace the 2D image-space noise with projected 3D world-space noise, where the noise input to G_{up} is the projection of samples from a grid of noise, n . This noise pattern is drawn from a standard Gaussian distribution defined on the ground plane at the same resolution of the layout features, which is then lifted into 3D and volume rendered along each ray \mathbf{r} :

$$n(\mathbf{r}) = \sum_{i=1}^N w_i n(x, z). \quad (7)$$

3.4. Skydome

We model remote content (sky and distant mountains) separately with a skydome generator G_{sky} (Fig. 4). This generator follows the StyleGAN3 architecture [33], with a mapping

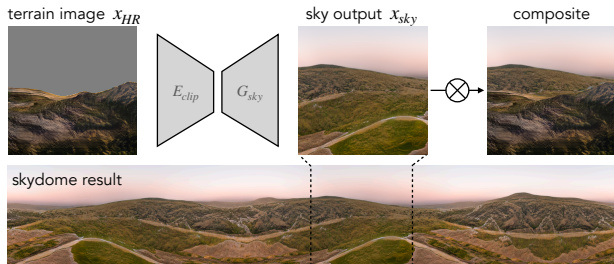


Figure 4. *Skydome generator*. Conditioned on the terrain image, the skydome generator G_{sky} synthesizes distant content (e.g., sky pixels and remote mountains) that is consistent with the generated terrain using encoder E_{clip} . G_{sky} is conditioned on cylindrical coordinates which can be unwrapped to produce a panoramic skydome image.

network and synthesis network conditioned on cylindrical coordinates [4]. We adapt it by conditioning on the terrain output: we encode terrain images I_{HR} using the pretrained CLIP image encoder E_{clip} [60], and concatenate this to the style-code output of the mapping network as input into G_{sky} :

$$I_{sky} = G_{sky}(\text{concat}(E_{clip}(I_{HR}), \text{mapping}(\mathbf{z}))). \quad (8)$$

Conditioning on the foreground terrain image encourages the skydome generator to generate a sky that is consistent with the terrain content. This model trains on single-view landscape images but can produce a full panorama at inference-time by passing in coordinates that correspond to a 360° cylinder. The skydome is rendered to an individual camera viewpoint using camera ray directions, giving the skydome image I_{dome} which is then composited with the terrain image using the sky mask:

$$I_{full} = I_{HR} \odot m_{HR} + I_{dome} \odot (1 - m_{HR}). \quad (9)$$

3.5. Training

We train the layout generator (rendering at 32×32), refinement network (upsampling to 256×256), and skydome generator separately. To train the refinement network, we operate on outputs of the layout generator, freezing the weights of that model. For the skydome generator, we train using real landscape images, and apply it only to the outputs of the refinement network at inference time. We follow the StyleGAN2 objective [34], with additional losses for each training stage, architecture, and hyperparameters provided in the supplemental.

Dataset and camera poses. We train on LHQ [76], a dataset of 90K unposed, single-view images of natural landscapes. A number of LHQ images contain geometry that is not amenable to “flying”, such as a landscape pictured through a window, or a closeup of trees. Therefore, we perform a filtering process on LHQ prior to training (see supplemental). We also obtain auxiliary outputs – disparity and sky segmentation – using the pretrained DPT [61] model. Disparity

and sky segmentation are used to construct the real image distribution in the GAN training phases.

After filtering, we use 56,982 images for training, and augment with horizontal flipping. During training we also need to sample camera poses. Prior 3D generators [5, 6, 14, 20, 58, 69] either use ground-truth poses from a simulator, or assume an object-centric camera distribution in which the camera looks at a fixed origin from some radius. Because our dataset lacks ground truth poses, we first sample a bank of training poses uniformly across the layout feature grid with random small height offsets, and rotate such that the near half of the camera view frustum falls entirely within the layout grid. Since the aerial layout should not be specific to any given camera pose, we generate f_{land} without any camera pose information, and then adopt the sampling scheme from GSN [14] which samples a camera pose from the initial training pose bank proportional to the inverse terrain density at each camera position, to avoid placing the camera within occluding geometry.

4. Experiments

Given its persistent scene representation and the extensibility of its layout grid, our model enables arbitrary motion through a synthesized landscape, including long camera trajectories. We show sample outputs from our model under a variety of camera movements (§ 4.1); present qualitative and quantitative comparisons with alternate scene representations, including auto-regressive prediction models and unconditional generators defined for bounded or object-centric scenes (§ 4.2); and investigate variations of our model to evaluate design decisions (§ 4.3).

4.1. Persistent, unbounded scene synthesis

Figure 5 shows example landscapes generated by our model with various camera motions. As the camera moves (by rotating and/or translating) the generated imagery changes in a way that is consistent with the underlying geometry, e.g. hills move across the image or become closer. Extending the generated aerial feature grid allows us to place the camera *outside* the distribution of training camera poses, while maintaining both geometric and stylistic consistency. As illustrated in Figure 1 and our project page, the persistent and extendable layout features enables synthetic ‘flights’ over large distances that can also return to a consistent starting point.

4.2. Comparing scene representations

We compare our model with three state-of-the-art methods. InfiniteNature-Zero is an auto-regressive method that, given an initial frame, generates successive frames sequentially by warping each image to the next based on depth [42]. It allows for unbounded camera trajectories, but has no persistent world model. GSN [14] and EG3D [5] are uncon-

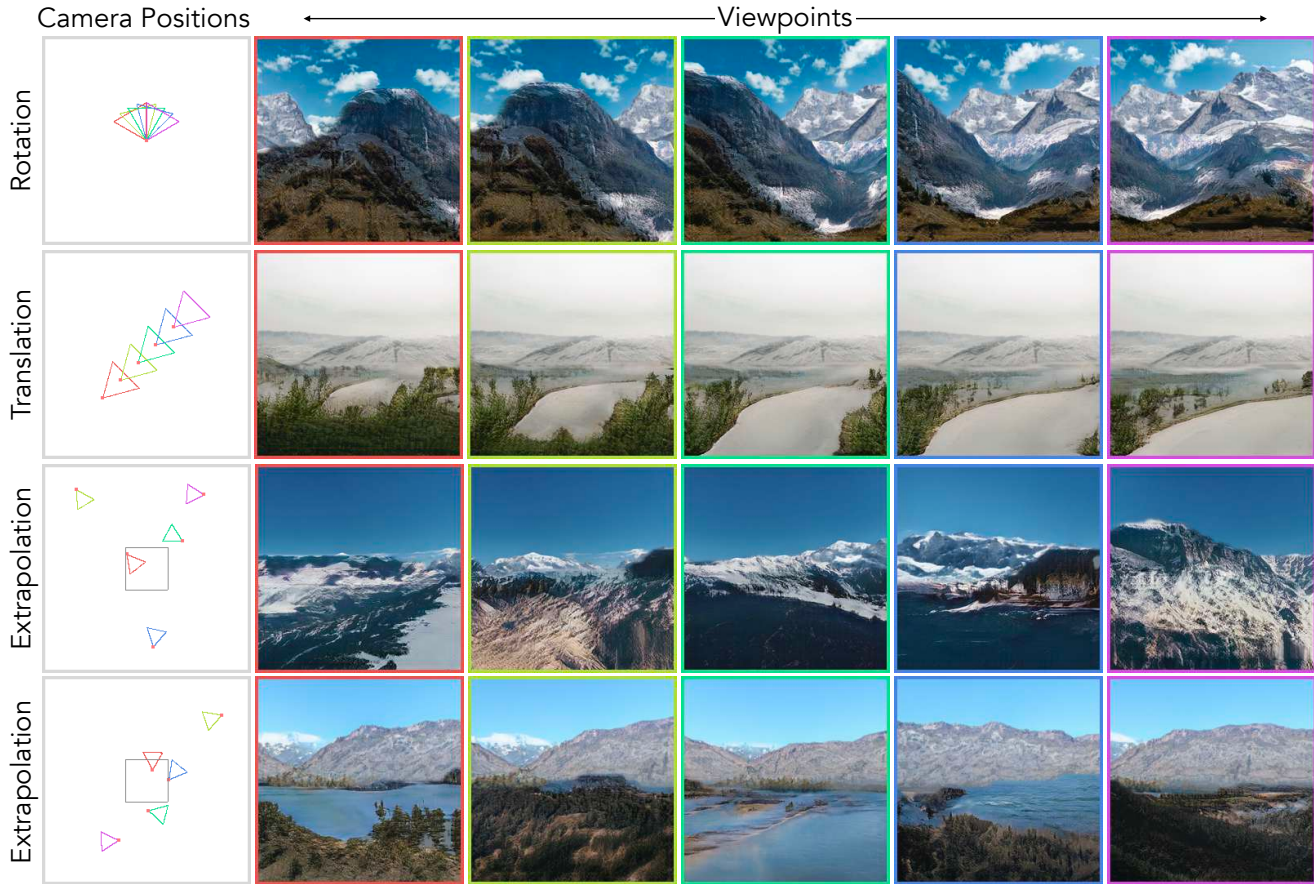


Figure 5. *Visualization of nearby and extrapolated camera motion.* Each row shows a set of sampled viewpoints, shown in an overhead view in the first column, and the corresponding rendered images in the other columns. Our model enables 3D-consistent view synthesis, visible under rotating or translating camera trajectories. We can also extrapolate the layout features at inference time, enabling camera motions outside of the training camera distribution (shown as a black square in the last two rows) with a consistent scene style.

Model	Persistent	Unbounded	FID		Consistency	
			C_{forward}	1-step	cycle	
Inf Nat Zero [42]	✗	✓	28.15	1.84	3.94	
Ours (128px)	✓	✓	26.09	2.12	0.00	

Table 1. *Comparison with InfiniteNature-Zero.* Using camera motions from InfiniteNature-Zero, we evaluate image quality as FID on 5K images after moving 100 steps forward (C_{forward}), one-step consistency as the L1 error when backwards warping one camera step, and cycle consistency as the L1 error between the original frame and the result after a pair of forward/backward steps. InfiniteNature-Zero is more consistent for a single step, but it has non-zero cyclic consistency error, and image quality degrades after repeated model applications. L1 values are multiplied by 100 throughout.

ditional generative models: GSN uses a layout feature grid which is also the basis of our model, but focuses on bounded indoor scenes with ground-truth camera pose trajectories, while EG3D uses a tri-plane representation and primarily focuses on objects and portraits. These methods have persistent world models (feature grid and tri-plane representation) but do not allow for unbounded trajectories.

Model	Persistent	Unbounded	FID			Consistency
			C_{train}	C_{forward}	C_{random}	
GSN [14]	✓	✗	29.95	50.22	45.48	12.80
EG3D [5]	✓	✗	9.85	30.17	32.08	3.01
Ours	✓	✓	21.42	26.67	23.39	3.56

Table 2. *Quantitative comparison to unconditional GANs.* We evaluate image quality as FID on 5K images on (a) training camera poses C_{train} , (b) forward motion C_{forward} (See Table 1), (c) random camera poses C_{random} . One-step consistency error is measured as the L1 error when backwards warping the result after one camera step to the initial frame, multiplied by 100. Once outside the training pose distribution our model generates better images than other methods, with consistency close to that of EG3D.

Quantitative comparisons. We evaluate image quality using FID [25], and multi-view consistency using photometric error. To compare with InfiniteNature-Zero (Table 1), we initialize with an image and depth map from our model, move the camera forwards using a forward motion trajectory from InfiniteNature-Zero, and evaluate image quality at a distance of 100 forward steps. Our model attains better FID, showing that it does not suffer from image degradation due to succes-

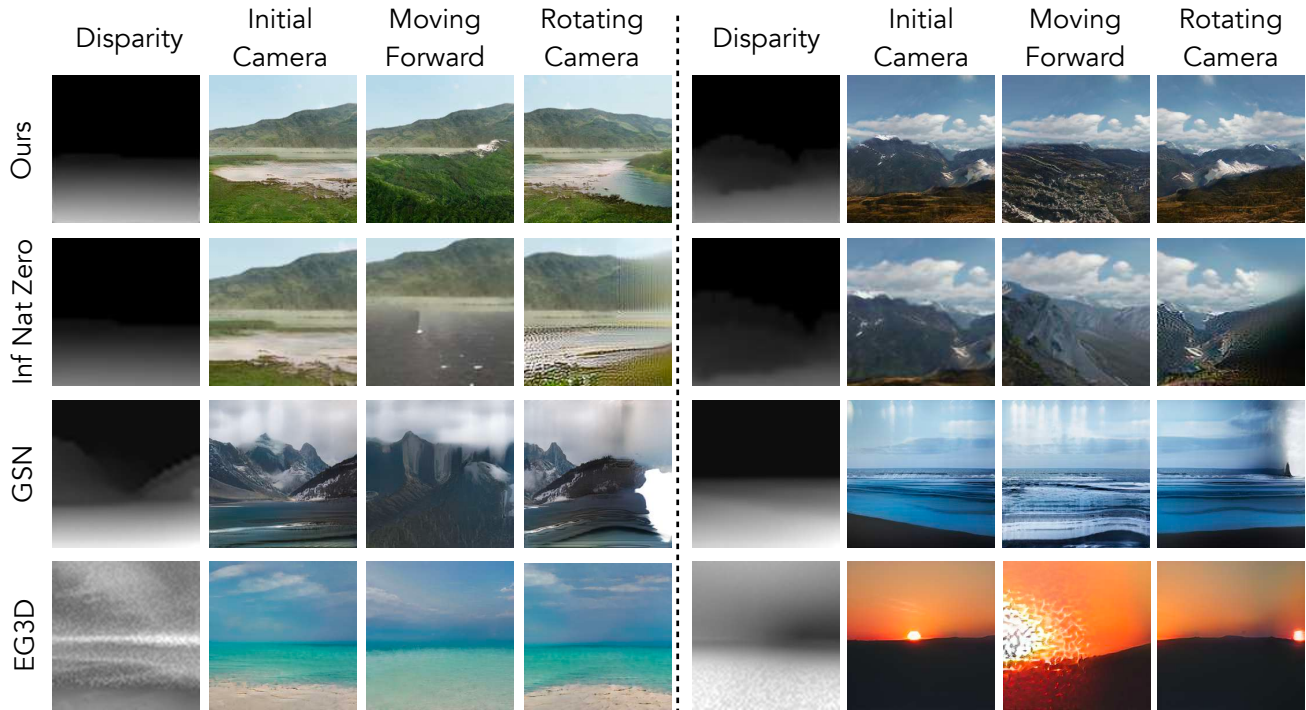


Figure 6. *Comparison to auto-regressive and bounded-volume 3D generative models.* Each row shows results for a given method on two generated scenes under different camera motion, along with a disparity map. Compared to InfiniteNature-Zero, our model enables long-range view synthesis by rendering a global scene description from different viewpoints, rather than auto-regressively predicting successive frames. 3D generative models like EG3D and GSN do not support view extrapolation on unbounded scenes. See our webpage for animated results.

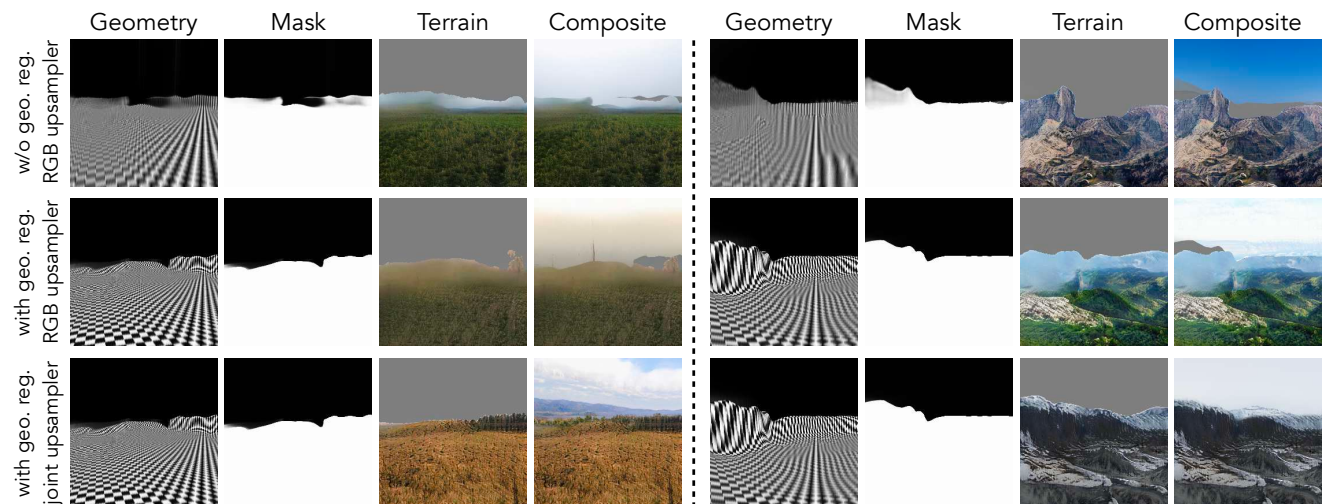


Figure 7. *Qualitative comparison of model variations.* Each row shows a model variant, visualizing generated geometry (as a rendered scene filled with a checkerboard pattern), sky mask, rendered terrain, and final image composite. (Top) Without geometry regularization, the model produces semi-transparent terrain. (Middle) Adding geometry regularization (Eqn. 4) makes the terrain more solid, but there are inconsistencies between the terrain and mask prediction. (Bottom) Our full model uses geometry regularization and also adds a upsampler that operates on inverse-depth and sky mask inputs in addition to RGB (Eqn. 6) to discourage boundary effects between the terrain and sky.

sive applications of an auto-regressive model. To compute one-step consistency error, we generate a new frame at a position equivalent to one forward step of InfiniteNature-Zero, warp it back to the original camera position using depth, and compute L1 error with the original frame in the overlapping

region. Because InfiniteNature-Zero uses explicit warping as part of its model, it can achieve better one-step consistency, whereas our 2D upsampling operation is more susceptible to geometric inconsistency. We measure cyclic consistency error as the L1 error between the initial frame to the result

after a step forward and back. Because InfiniteNature-Zero lacks a persistent global representation, it has non-zero cyclic consistency error, whereas our model is fully consistent with zero cyclic consistency error.

To compare with the unconditional generative models GSN and EG3D, we compute FID on sets of output images corresponding to different distributions of camera positions: camera poses used in training which are intended to overlap with the layout, camera poses 100 steps forward from these mimicking InfiniteNature-Zero trajectories, and a uniform distribution of randomly oriented cameras over the layout grid. As seen in Table 2, GSN is the least successful method when applied to this domain. EG3D generates high-quality images at training camera poses, but tends to represent the scene as floating nearby clouds with planar mountains at the edges of the volume (incorrect geometry). Our method generalizes better to new camera positions. GSN has the highest one-step consistency error, while the consistency error of our model is close to that of EG3D (which relies less on 2D upsampling). In the supplemental, we experiment with an alternative architecture that builds on extendable triplane units with lower consistency error and faster rendering speed.

Qualitative comparisons. In Fig. 6 we show example outputs of each model over forward-moving and rotating trajectories. Due to its auto-regressive nature, the quality of InfiniteNature-Zero’s output degrades somewhat as the camera trajectory becomes longer. A more serious limitation is that, trained only on forward movement, it is unable to synthesize plausible views under camera rotation. GSN and EG3D also struggle with long camera trajectories, producing unrealistic outputs as the cameras approach the spatial limits of the training camera distribution. In the case of GSN applied to our setting, the results contain flickering and grid-like artifacts, which our projected noise (§ 3.3) mitigates.

4.3. Model Variations

To investigate individual components of our model, we separately evaluate variations of the layout generator and refinement network.

Layout generator. The resolution of the scene layout grid and the number of samples per ray affect the quality of the volume-rendered output I_{LR} . As shown in Table 3, higher resolution and more samples lead to the best image quality (FID computed on 32×32 pixel images for speed, compared to segmented real images with gray sky pixels). To maximize the capacity of layout generation and rendering within computational limits we opt for a 256×256 feature grid with 128 samples per ray.

Refinement network. Next, we investigate the refinement stage, which upsamples and refines the layout generator output. In our full model, the refinement network operates not only on RGB images but also on inverse-depth and sky mask

Model	Samples Per Ray	Layout Resolution	FID (I_{LR})
Low	64	32	33.66
Medium	128	32	32.02
High	128	64	22.62
Full	128	256	16.06

Table 3. *Variations on layout generation.* Higher feature grid resolution and more samples per ray yield the best results, bounded by computational limits. For speed, FID is computed on 5K samples rendered at 32×32 .

Refinement Output	Projected Noise	FID (I_{HR})		Consistency
		C_{train}	C_{random}	
I_{HR}	✗	26.30	27.08	5.08
I_{HR}, d_{HR}, m_{HR}	✗	23.75	27.25	5.81
I_{HR}, d_{HR}, m_{HR}	✓	21.42	23.39	3.91

Table 4. *Variations on the refinement network.* We find refining not only the low-resolution image but also the depth and sky-mask improves image quality, but can lead to jittery results. The addition of projected noise into the upsampler results in smoother frames with lower consistency error.

(Eqn. 6), and uses projected noise for spatial consistency of texture detail (Eqn. 7). As shown in Table 4, both help to improve our model’s FID and consistency error.

As shown in Fig. 7 (second row), upsampling only the RGB image I_{LR} can lead to output that is inconsistent with the generated sky mask, leading to temporally unstable gaps in the final composited image. This figure also shows the effect of our geometric regularization (Eqn. 4) in reducing unwanted transparency, especially in distant terrain.

5. Discussion and conclusion

Limitations. A few drawbacks of our model include costly volume rendering limiting the resolution of I_{LR} , imperfect 3D consistency due to image-space refinement, and imperfect or repeating geometry decoded from the scene layout features. We elaborate in the supplemental.

Conclusion. We present an unconditional world generator for unbounded synthesis of persistent 3D nature scenes. We build persistent world representation by modeling scene content with a spatially extendable layout feature grid which can be decoded via volume rendering to form a terrain image. This rendered terrain is combined with a separate sky-dome, representing infinitely far content, to synthesize novel viewpoints supporting nearby and distant camera motions. Altogether, our model enables 3D consistent image generation and view synthesis of unbounded scenes learned from single-view, unposed landscape photos.

Acknowledgements. Thanks to Andrew Liu and Richard Bowen for the fruitful discussions and helpful comments. This project was part of an internship at Google.

References

- [1] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-NeRF 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022. [2](#)
- [2] Richard Strong Bowen, Huiwen Chang, Charles Herrmann, Piotr Teterwak, Ce Liu, and Ramin Zabih. OCONet: Image extrapolation by object completion. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 2307–2317, 2021. [2](#)
- [3] Tim Brooks, Janne Hellsten, Miika Aittala, Ting-Chun Wang, Timo Aila, Jaakko Lehtinen, Ming-Yu Liu, Alexei A Efros, and Tero Karras. Generating long videos of dynamic scenes. *arXiv preprint arXiv:2206.03429*, 2022. [1](#), [2](#)
- [4] Lucy Chai, Michael Gharbi, Eli Shechtman, Phillip Isola, and Richard Zhang. Any-resolution training for high-resolution image synthesis. In *European Conference on Computer Vision*, 2022. [5](#)
- [5] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, 2022. [1](#), [2](#), [5](#), [6](#)
- [6] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-GAN: Periodic implicit generative adversarial networks for 3D-aware image synthesis. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 5799–5809, 2021. [5](#)
- [7] Xu Chen, Jie Song, and Otmar Hilliges. Monocular neural image based rendering with continuous view control. In *Proc. Int. Conf. on Computer Vision (ICCV)*, pages 4090–4100, 2019. [2](#)
- [8] Zhaoxi Chen, Guangcong Wang, and Ziwei Liu. Scenedreamer: Unbounded 3d scene generation from 2d image collections. *arXiv preprint arXiv:2302.01330*, 2023. [3](#)
- [9] Yen-Chi Cheng, Chieh Hubert Lin, Hsin-Ying Lee, Jian Ren, Sergey Tulyakov, and Ming-Hsuan Yang. In&Out: Diverse image outpainting via GAN inversion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11431–11440, 2022. [2](#)
- [10] Inchang Choi, Orazio Gallo, Alejandro Troccoli, Min H Kim, and Jan Kautz. Extreme view synthesis. In *Proc. Int. Conf. on Computer Vision (ICCV)*, pages 7781–7790, 2019. [2](#)
- [11] Min Jin Chong, Hsin-Ying Lee, and David Forsyth. StyleGAN of All Trades: Image Manipulation with Only Pretrained StyleGAN. *arXiv preprint arXiv:2111.01619*, 2021. [4](#)
- [12] Yu Deng, Jiaolong Yang, Jianfeng Xiang, and Xin Tong. Gram: Generative radiance manifolds for 3D-aware image generation. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 10673–10683, June 2022. [2](#)
- [13] Emily Denton and Rob Fergus. Stochastic video generation with a learned prior. In *Proc. Int. Conf. on Machine Learning (ICML)*, pages 1174–1183. PMLR, 2018. [2](#)
- [14] Terrance DeVries, Miguel Angel Bautista, Nitish Srivastava, Graham W Taylor, and Joshua M Susskind. Unconstrained scene generation with locally conditioned radiance fields. In *Proc. Int. Conf. on Computer Vision (ICCV)*, pages 14304–14313, 2021. [2](#), [3](#), [5](#), [6](#)
- [15] Chelsea Finn, Ian Goodfellow, and Sergey Levine. Unsupervised learning for physical interaction through video prediction. In *Neural Information Processing Systems*, 2016. [2](#)
- [16] John Flynn, Michael Broxton, Paul Debevec, Matthew DuVall, Graham Fyffe, Ryan Overbeck, Noah Snavely, and Richard Tucker. DeepView: View synthesis with learned gradient descent. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 2367–2376, 2019. [2](#)
- [17] Gereon Fox, Ayush Tewari, Mohamed Elgharib, and Christian Theobalt. StyleVideoGAN: A temporal generative model using a pretrained StyleGAN. In *Proc. British Machine Vision Conf. (BMVC)*, 2021. [2](#)
- [18] Rafail Fridman, Amit Abecasis, Yoni Kasten, and Tali Dekel. Scenescape: Text-driven consistent scene generation. *arXiv preprint arXiv:2302.01133*, 2023. [3](#)
- [19] Songwei Ge, Thomas Hayes, Harry Yang, Xi Yin, Guan Pang, David Jacobs, Jia-Bin Huang, and Devi Parikh. Long video generation with time-agnostic VQGAN and time-sensitive transformer. *arXiv preprint arXiv:2204.03638*, 2022. [2](#)
- [20] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. StyleNeRF: A style-based 3D-aware generator for high-resolution image synthesis. *arXiv preprint arXiv:2110.08985*, 2021. [2](#), [5](#)
- [21] David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. *Advances in neural information processing systems*, 31, 2018. [1](#)
- [22] Zekun Hao, Arun Mallya, Serge Belongie, and Ming-Yu Liu. Gancraft: Unsupervised 3D neural rendering of Minecraft worlds. In *Proc. Int. Conf. on Computer Vision (ICCV)*, pages 14072–14082, 2021. [2](#)
- [23] Erik Härkönen, Aaron Hertzmann, Jaakko Lehtinen, and Sylvain Paris. Ganspace: Discovering interpretable gan controls. In *Neural Information Processing Systems*, 2020. [1](#)
- [24] James Hays and Alexei A Efros. Scene completion using millions of photographs. In *ACM Trans. Graphics (SIGGRAPH North America)*, 2007. [2](#)
- [25] 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. *Neural Information Processing Systems*, 30, 2017. [6](#)
- [26] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. [2](#)
- [27] Jun-Ting Hsieh, Bingbin Liu, De-An Huang, Li F Fei-Fei, and Juan Carlos Niebles. Learning to decompose and disentangle representations for video prediction. In *Neural Information Processing Systems*, volume 31, 2018. [2](#)
- [28] Ronghang Hu, Nikhila Ravi, Alexander C. Berg, and Deepak Pathak. Worldsheet: Wrapping the world in a 3D sheet for view synthesis from a single image. In *Proc. Int. Conf. on Computer Vision (ICCV)*, 2021. [2](#)

- [29] Ali Jahanian, Lucy Chai, and Phillip Isola. On the "steerability" of generative adversarial networks. In *Proc. Int. Conf. on Learning Representations (ICLR)*, 2020. 1
- [30] Varun Jampani, Huiwen Chang, Kyle Sargent, Abhishek Kar, Richard Tucker, Michael Krainin, Dominik Kaeser, William T Freeman, David Salesin, Brian Curless, et al. SLIDE: Single image 3D photography with soft layering and depth-aware inpainting. In *Proc. Int. Conf. on Computer Vision (ICCV)*, pages 12518–12527, 2021. 2
- [31] Wonbong Jang and Lourdes Agapito. CodeNeRF: Disentangled neural radiance fields for object categories. In *Proc. Int. Conf. on Computer Vision (ICCV)*, pages 12949–12958, 2021. 2
- [32] Biliana Kaneva, Josef Sivic, Antonio Torralba, Shai Avidan, and William T. Freeman. Infinite images: Creating and exploring a large photorealistic virtual space. In *Proceedings of the IEEE*, 2010. 2
- [33] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. In *Neural Information Processing Systems*, 2021. 4
- [34] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 8110–8119, 2020. 3, 5
- [35] Seung Wook Kim, Jonah Philion, Antonio Torralba, and Sanja Fidler. DriveGAN: Towards a controllable high-quality neural simulation. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, 2021. 1
- [36] Jing Yu Koh, Honglak Lee, Yinfei Yang, Jason Baldridge, and Peter Anderson. Pathdreamer: A world model for indoor navigation. In *Proc. Int. Conf. on Computer Vision (ICCV)*, pages 14738–14748, 2021. 2
- [37] Johannes Kopf, Kevin Matzen, Suhub Alsisan, Ocean Quigley, Francis Ge, Yangming Chong, Josh Patterson, Jan-Michael Frahm, Shu Wu, Matthew Yu, Peizhao Zhang, Zijian He, Peter Vajda, Ayush Saraf, and Michael Cohen. One shot 3D photography. *ACM Transactions on Graphics (Proceedings of ACM SIGGRAPH)*, 39(4), 2020. 2
- [38] Wonkwang Lee, Whie Jung, Han Zhang, Ting Chen, Jing Yu Koh, Thomas Huang, Hyungsuk Yoon, Honglak Lee, and Seunghoon Hong. Revisiting hierarchical approach for persistent long-term video prediction. *arXiv preprint arXiv:2104.06697*, 2021. 2
- [39] Marc Levoy and Pat Hanrahan. Light field rendering. In *ACM Trans. Graphics (SIGGRAPH North America)*, 1996. 2
- [40] Jiaxin Li, Zijian Feng, Qi She, Henghui Ding, Changhu Wang, and Gim Hee Lee. MINE: Towards continuous depth MPI with NeRF for novel view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 12578–12588, October 2021. 2
- [41] Wenbo Li, Zhe Lin, Kun Zhou, Lu Qi, Yi Wang, and Jiaya Jia. MAT: Mask-aware transformer for large hole image inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10758–10768, 2022. 2
- [42] Zhengqi Li, Qianqian Wang, Noah Snavely, and Angjoo Kanazawa. InfiniteNature-Zero: Learning perpetual view generation of natural scenes from single images. In *European Conference on Computer Vision*, pages 515–534. Springer, 2022. 1, 2, 5, 6
- [43] Chieh Hubert Lin, Hsin-Ying Lee, Yen-Chi Cheng, Sergey Tulyakov, and Ming-Hsuan Yang. InfinityGAN: Towards infinite-pixel image synthesis. *arXiv preprint arXiv:2104.03963*, 2021. 2
- [44] Chieh Hubert Lin, Hsin-Ying Lee, Willi Menapace, Menglei Chai, Aliaksandr Siarohin, Ming-Hsuan Yang, and Sergey Tulyakov. InfiniCity: Infinite-scale city synthesis. *arXiv preprint arXiv:2301.09637*, 2023. 3
- [45] 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. In *Proc. Int. Conf. on Computer Vision (ICCV)*, pages 14458–14467, 2021. 1, 2
- [46] Hongyu Liu, Ziyu Wan, Wei Huang, Yibing Song, Xintong Han, and Jing Liao. PD-GAN: Probabilistic diverse GAN for image inpainting. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 9371–9381, 2021. 2
- [47] Yuchen Liu, Zhixin Shu, Yijun Li, Zhe Lin, Federico Perazzi, and Sun-Yuan Kung. Content-aware GAN compression. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 12156–12166, 2021. 2
- [48] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM Trans. Graph.*, 38(4):65:1–65:14, July 2019. 2
- [49] Arun Mallya, Ting-Chun Wang, Karan Sapra, and Ming-Yu Liu. World-consistent video-to-video synthesis. In *Proc. European Conf. on Computer Vision (ECCV)*, pages 359–378. Springer, 2020. 2
- [50] 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. In *ACM Trans. Graphics (SIGGRAPH North America)*, 2019. 2
- [51] 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. In *Proc. European Conf. on Computer Vision (ECCV)*, pages 405–421. Springer, 2020. 2, 4
- [52] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *arXiv preprint arXiv:2201.05989*, 2022. 2
- [53] Andres Munoz, Mohammadreza Zolfaghari, Max Argus, and Thomas Brox. Temporal shift GAN for large scale video generation. In *Proc. Winter Conf. on Computer Vision (WACV)*, pages 3179–3188, 2021. 2
- [54] Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian Richardt, and Yong-Liang Yang. HoloGAN: Unsupervised learning of 3D representations from natural images. In *The IEEE International Conference on Computer Vision (ICCV)*, Nov 2019. 2

- [55] Michael Niemeyer and Andreas Geiger. CAMPARI: Camera-aware decomposed generative neural radiance fields. In *2021 International Conference on 3D Vision (3DV)*, pages 951–961. IEEE, 2021. 2
- [56] Michael Niemeyer and Andreas Geiger. GIRAFFE: Representing scenes as compositional generative neural feature fields. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 11453–11464, 2021. 2
- [57] Simon Niklaus, Long Mai, Jimei Yang, and Feng Liu. 3D Ken Burns effect from a single image. *ACM Trans. Graphics*, 38(6):1–15, 2019. 2
- [58] Roy Or-El, Xuan Luo, Mengyi Shan, Eli Shechtman, Jeong Joon Park, and Ira Kemelmacher-Shlizerman. StyleSDF: High-Resolution 3D-Consistent Image and Geometry Generation. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 13503–13513, June 2022. 2, 5
- [59] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3D using 2D diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 2
- [60] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *Proc. Int. Conf. on Machine Learning (ICML)*, pages 8748–8763. PMLR, 2021. 5
- [61] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *Proc. Int. Conf. on Computer Vision (ICCV)*, pages 12179–12188, 2021. 5
- [62] Daniel Rebain, Mark Matthews, Kwang Moo Yi, Dmitry Lagun, and Andrea Tagliasacchi. LOLNeRF: Learn from one look. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1558–1567, 2022. 2
- [63] Xuanchi Ren and Xiaolong Wang. Look outside the room: Synthesizing a consistent long-term 3D scene video from a single image. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 3563–3573, 2022. 2
- [64] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *Proc. European Conf. on Computer Vision (ECCV)*, 2020. 2
- [65] Chris Rockwell, David F Fouhey, and Justin Johnson. Pixel-Synth: Generating a 3D-consistent experience from a single image. In *Proc. Int. Conf. on Computer Vision (ICCV)*, pages 14104–14113, 2021. 2
- [66] Robin Rombach, Patrick Esser, and Björn Ommer. Geometry-free view synthesis: Transformers and no 3D priors. In *Proc. Int. Conf. on Computer Vision (ICCV)*, pages 14356–14366, 2021. 2
- [67] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–10, 2022. 2
- [68] Sara Fridovich-Keil and Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [69] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. GRAF: Generative radiance fields for 3D-aware image synthesis. *Neural Information Processing Systems*, 33:20154–20166, 2020. 2, 5
- [70] Prafull Sharma, Ayush Tewari, Yilun Du, Sergey Zakharov, Rares Ambrus, Adrien Gaidon, William T Freeman, Fredo Durand, Joshua B Tenenbaum, and Vincent Sitzmann. Seeing 3D objects in a single image via self-supervised static-dynamic disentanglement. *arXiv preprint arXiv:2207.11232*, 2022. 3
- [71] Yuan Shen, Wei-Chiu Ma, and Shenlong Wang. SGAM: Building a virtual 3D world through simultaneous generation and mapping. In *Neural Information Processing Systems*, 2022. 2
- [72] Yujun Shen and Bolei Zhou. Closed-form factorization of latent semantics in gans. In *arXiv preprint arXiv:2007.06600*, 2020. 1
- [73] Lixin Shi, Haitham Hassanieh, Abe Davis, Dina Katabi, and Fredo Durand. Light field reconstruction using sparsity in the continuous fourier domain. In *ACM Trans. Graphics (SIGGRAPH North America)*, 2014. 2
- [74] Meng-Li Shih, Shih-Yang Su, Johannes Kopf, and Jia-Bin Huang. 3D photography using context-aware layered depth inpainting. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [75] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022. 2
- [76] Ivan Skorokhodov, Grigorii Sotnikov, and Mohamed Elhoseiny. Aligning latent and image spaces to connect the unconnectable. In *Proc. Int. Conf. on Computer Vision (ICCV)*, pages 14144–14153, 2021. 5
- [77] Ivan Skorokhodov, Sergey Tulyakov, and Mohamed Elhoseiny. StyleGAN-V: A continuous video generator with the price, image quality and perks of StyleGAN2. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3626–3636, 2022. 2
- [78] Ivan Skorokhodov, Sergey Tulyakov, Yiqun Wang, and Peter Wonka. EpiGRAF: Rethinking training of 3D GANs. In *Neural Information Processing Systems*, 2022. 2
- [79] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky. Resolution-robust large mask inpainting with fourier convolutions. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2149–2159, 2022. 2
- [80] Piotr Teterwak, Aaron Sarna, Dilip Krishnan, Aaron Maschinot, David Belanger, Ce Liu, and William T Freeman. Boundless: Generative adversarial networks for image extension. In *Proc. Int. Conf. on Computer Vision (ICCV)*, pages 10521–10530, 2019. 2
- [81] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2
- [82] Shubham Tulsiani, Richard Tucker, and Noah Snavely. Layer-structured 3D scene inference via view synthesis. In *Proc. European Conf. on Computer Vision (ECCV)*, pages 302–317, 2018. 2

- [83] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. Mocogan: Decomposing motion and content for video generation. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 1526–1535, 2018. 2
- [84] Ruben Villegas, Jimei Yang, Seunghoon Hong, Xunyu Lin, and Honglak Lee. Decomposing motion and content for natural video sequence prediction. *arXiv preprint arXiv:1706.08033*, 2017. 2
- [85] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. In *Neural Information Processing Systems*, 2016. 2
- [86] Carl Vondrick and Antonio Torralba. Generating the future with adversarial transformers. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 1020–1028, 2017. 2
- [87] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. IBRNet: Learning multi-view image-based rendering. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 4690–4699, 2021. 2
- [88] Yunbo Wang, Mingsheng Long, Jianmin Wang, Zhifeng Gao, and Philip S Yu. PredRNN: Recurrent neural networks for predictive learning using spatiotemporal LSTMs. In *Neural Information Processing Systems*, pages 879–888, 2017. 2
- [89] Yi Wang, Xin Tao, Xiaoyong Shen, and Jiaya Jia. Wide-context semantic image extrapolation. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 1399–1408, 2019. 2
- [90] Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. SynSin: End-to-end view synthesis from a single image. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 7467–7477, 2020. 2
- [91] Zongxin Yang, Jian Dong, Ping Liu, Yi Yang, and Shuicheng Yan. Very long natural scenery image prediction by outpainting. In *Proc. Int. Conf. on Computer Vision (ICCV)*, pages 10561–10570, 2019. 2
- [92] Yufei Ye, Maneesh Singh, Abhinav Gupta, and Shubham Tulsiani. Compositional video prediction. In *Proc. Int. Conf. on Computer Vision (ICCV)*, 2019. 2
- [93] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. PixelNeRF: Neural radiance fields from one or few images. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 4578–4587, 2021. 2, 3
- [94] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Generative image inpainting with contextual attention. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, pages 5505–5514, 2018. 2
- [95] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Free-form image inpainting with gated convolution. In *Proc. Int. Conf. on Computer Vision (ICCV)*, pages 4471–4480, 2019. 2
- [96] Sihyun Yu, Jihoon Tack, Sangwoo Mo, Hyunsu Kim, Junho Kim, Jung-Woo Ha, and Jinwoo Shin. Generating videos with dynamics-aware implicit generative adversarial networks. In *The Tenth International Conference on Learning Representations*, 2022. 2
- [97] Shengyu Zhao, Jonathan Cui, Yilun Sheng, Yue Dong, Xiao Liang, Eric I Chang, and Yan Xu. Large scale image completion via co-modulated generative adversarial networks. In *Proc. Int. Conf. on Learning Representations (ICLR)*, 2021. 2
- [98] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. In *ACM Trans. Graphics (SIGGRAPH North America)*, 2018. 2