

Infinite Nature: Perpetual View Generation of Natural Scenes from a Single Image

Andrew Liu*
Ameesh Makadia

Richard Tucker*
Noah Snavely
Google Research

Varun Jampani
Angjoo Kanazawa

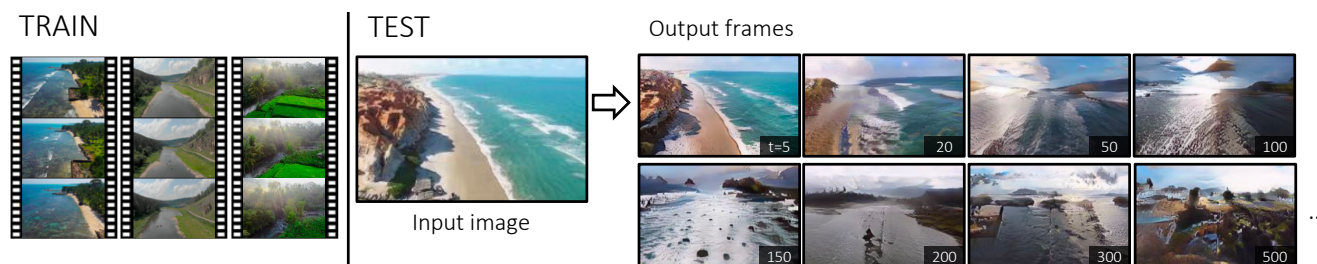


Figure 1. **Perpetual View Generation.** Using a collection of aerial videos of nature scenes for training (left), our method learns to take a single image and perpetually generate novel views for a camera trajectory covering a long distance (right). Our method can successfully generate hundreds of frames of an aerial video from a single input image (up to 500 shown here).

Abstract

We introduce the problem of perpetual view generation—long-range generation of novel views corresponding to an arbitrarily long camera trajectory given a single image. This is a challenging problem that goes far beyond the capabilities of current view synthesis methods, which quickly degenerate when presented with large camera motions. Methods for video generation also have limited ability to produce long sequences and are often agnostic to scene geometry. We take a hybrid approach that integrates both geometry and image synthesis in an iterative ‘render, refine and repeat’ framework, allowing for long-range generation that cover large distances after hundreds of frames. Our approach can be trained from a set of monocular video sequences. We propose a dataset of aerial footage of coastal scenes, and compare our method with recent view synthesis and conditional video generation baselines, showing that it can generate plausible scenes for much longer time horizons over large camera trajectories compared to existing methods. Project page at <https://infinite-nature.github.io>.

1. Introduction

Consider the input image of a coastline in Fig. 1. Imagine flying through this scene as a bird. Initially, we would see objects grow in our field of view as we approach them.

* indicates equal contribution

Beyond, we might find a wide ocean or new islands. At the shore, we might see cliffs or beaches, while inland there could be mountains or forests. As humans, we are adept at imagining a plausible world from a single picture, based on our own experience.

How can we emulate this ability on a computer? One approach would be to attempt to generate an entire 3D planet with high-resolution detail from a single image. However, this would be extremely expensive and far beyond the current state of the art. So, we pose the more tractable problem of *perpetual view generation*: given a single image of scene, the task is to synthesize a video corresponding to an arbitrary camera trajectory. Solving this problem can have applications in content creation, novel photo interactions, and methods that use learned world models like model-based reinforcement learning.

Perpetual view generation, though simple to state, is an extremely challenging task. As the viewpoint moves, we must extrapolate new content in unseen regions and also synthesize new detail in existing regions that are now closer to the camera. Two active areas of research, video synthesis and view synthesis, both fail to scale to this problem for different reasons.

Recent video synthesis methods apply developments in image synthesis [20] to the temporal domain, or rely on recurrent models [10]. But they can generate only limited numbers of novel frames (e.g., 25 [41] or 48 frames [9]). Additionally, such methods often neglect an important el-

ement of the video’s structure—they model neither scene geometry nor camera movement. In contrast, many view synthesis methods do take advantage of geometry to synthesize high-quality novel views. However, these approaches can only operate within a limited range of camera motion. As shown in Figure 6, once the camera moves outside this range, such methods fail catastrophically.

We propose a hybrid framework that takes advantage of both geometry and image synthesis techniques to address these challenges. We use disparity maps to represent a scene’s geometry, and decompose the perpetual view generation task into the framework of *render-refine-and-repeat*. First, we **render** the current frame from a new viewpoint, using disparity to ensure that scene content moves in a geometrically correct manner. Then, we **refine** the resulting image and geometry. This step adds detail and synthesizes new content in areas that require inpainting or outpainting. Because we refine both the image and disparity, the whole process can be **repeated** in an recurrent manner, allowing for perpetual generation with arbitrary trajectories.

To train our system, we curated a large dataset of drone footage of nature and coastal scenes from over 700 videos, spanning 2 million frames. We run a structure from motion pipeline to recover 3D camera trajectories, and refer to this as the Aerial Coastline Imagery Dataset (ACID). Our trained model can generate sequences of hundreds of frames while maintaining the aesthetic feel of an aerial coastal video, even though after just a few frames, the camera has moved beyond the limits of the scene depicted in the initial view.

Our experiments show that our novel render-refine-repeat framework, with propagation of geometry via disparity maps, is key to tackling this problem. Compared to recent view synthesis and video generation baselines, our approach can produce plausible frames for much longer time horizons. This work represents a significant step towards perpetual view generation, though it has limitations such as a lack of global consistency in the hallucinated world. We believe our method and dataset will lead to further advances in generative methods for large-scale scenes.

2. Related Work

Image extrapolation. Our work is inspired by the seminal work of Kaneva *et al.* [19], which proposed a non-parametric approach for generating ‘infinite’ images through stitching 2D-transformed images, and by patch-based non-parametric approaches for image extension [29, 1]. We revisit the ‘infinite images’ concept in a learning framework that also reasons about the 3D geometry behind each image. Also related to our work are recent deep learning approaches to the problem of *outpainting*, i.e., inferring unseen content outside image boundaries [44, 46, 36], as well as *inpainting*, the task of filling in missing content within an image [15, 50]. These approaches use adversarial frameworks and

semantic information for in/outpainting. Our problem also incorporates aspects of super-resolution [14, 22]. Image-specific GAN methods also demonstrate a form of image extrapolation and super-resolution of textures and natural images [53, 34, 30, 33]. In contrast to the above methods, we reason about the 3D geometry behind each image and study image extrapolation in the context of temporal image sequence generation.

View synthesis. Many view synthesis methods operate by interpolating between multiple views of a scene [23, 3, 24, 12, 7], although recent work can generate new views from just a single input image, as in our work [5, 39, 25, 38, 31, 6]. However, in both settings, most methods only allow for a very limited range of output viewpoints. Even methods that explicitly allow for view extrapolation (not just interpolation) typically restrict the camera motion to small regions around a reference view [52, 35, 8].

One factor that limits camera motion is that many methods construct a static scene representation, such as a layered depth image [39, 32], multiplane image [52, 38], point cloud [25, 45], or radiance field [48, 37], and inpaint disoccluded regions. Such representations can allow for fast rendering, but the range of viable camera positions is limited by the finite bounds of the scene representation. Some methods augment this scene representation paradigm, enabling a limited increase in the range of output views. Niklaus *et al.* perform inpainting *after* rendering [25], while SynSin uses a post-rendering refinement network to produce realistic images from feature point-clouds [45]. We take inspiration from these methods by rendering and then refining our output. In contrast, however, our system does not construct a single 3D representation of a scene. Instead we proceed iteratively, generating each output view from the previous one, and producing a geometric scene representation in the form of a disparity map for each frame.

Some methods use video as training data. Monocular depth can be learned from 3D movie left-right camera pairs [27] or from video sequences analysed with structure-from-motion techniques [4]. Video can also be directly used for view synthesis [38, 45]. These methods use pairs of images, whereas our model is trained on sequences of several widely-spaced frames since we want to generate long-range video.

Video synthesis. Our work is related to methods that generate a video sequence from one or more images [42, 11, 43, 10, 40, 47]. Many such approaches have focused on predicting the future of dynamic objects with a static camera, often using simple videos of humans walking [2] or robot arms [11]. In contrast, we focus on mostly static scenes with a moving camera, using real aerial videos of nature. Some recent research addresses video synthesis from in-the-wild videos with moving cameras [9, 41], but without taking geometry explicitly into account, and with strict limits on the length of the generated video. In contrast, in our work

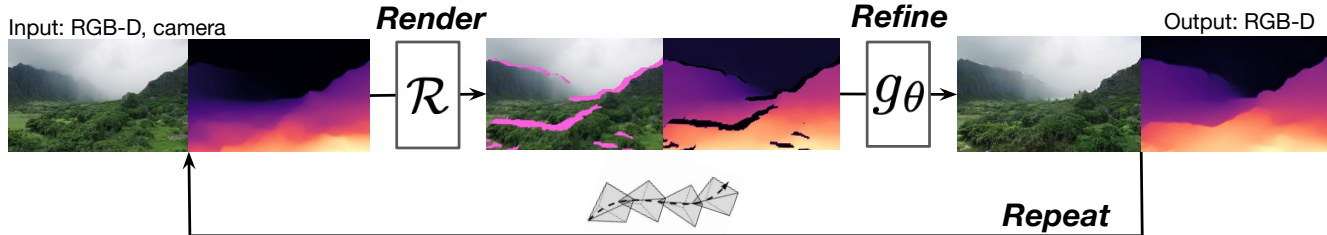


Figure 2. **Overview.** We first *render* an input image to a new camera view using the disparity. We then *refine* the image, synthesizing and super-resolving missing content. As we output both RGB and geometry, this process can be *repeated* for perpetual view generation.

the movement of pixels from camera motion is explicitly modeled using 3D geometry.

3. Perpetual View Generation

Given an RGB image I_0 and a camera trajectory (P_0, P_1, P_2, \dots) of arbitrary length, our task is to output a new image sequence (I_0, I_1, I_2, \dots) that forms a video depicting a flythrough of the scene captured by the initial view. The trajectory is a series of 3D camera poses $P_t = \begin{pmatrix} R^{3 \times 3} & \mathbf{t}^{3 \times 1} \\ 0 & 1 \end{pmatrix}$, where R and \mathbf{t} are 3D rotations and translations, respectively. In addition, each camera has an intrinsic matrix K . At training time camera data is obtained from video clips via structure-from-motion as in [52]. At test time, the camera trajectory may be pre-specified, generated by an auto-pilot algorithm, or controlled via a user interface.

3.1. Approach: Render, Refine, Repeat

Our framework applies established techniques (3D rendering, image-to-image translation, auto-regressive training) in a novel combination. We decompose *perpetual view generation* into the three steps, as illustrated in Figure 2:

1. **Render** a new view from an old view, by warping the image according to a disparity map using a differentiable renderer,
2. **Refine** the rendered view and geometry to fill in missing content and add detail where necessary,
3. **Repeat** this process, propagating *both image and disparity* to generate each new view from the one before.

Our approach has several desirable characteristics. *Representing geometry* with a disparity map allows much of the heavy lifting of moving pixels from one frame to the next to be handled by differentiable rendering, ensuring local temporal consistency. The synthesis task then becomes one of *image refinement*, which comprises: 1) inpainting disoccluded regions 2) outpainting of new image regions and 3) super-resolving image content. Because every step is *fully differentiable*, we can train our refinement network by back-propagating through several view synthesis iterations. Our *auto-regressive* framework means that novel views may be infinitely generated with explicit view control, even though

training data is finite in length.

Formally, for an image I_t with pose P_t we have an associated disparity (i.e., inverse depth) map $D_t \in \mathbb{R}^{H \times W}$, and we compute the next frame I_{t+1} and its disparity D_{t+1} as

$$\hat{I}_{t+1}, \hat{D}_{t+1}, \hat{M}_{t+1} = \mathcal{R}(I_t, D_t, P_t, P_{t+1}), \quad (1)$$

$$I_{t+1}, D_{t+1} = g_\theta(\hat{I}_{t+1}, \hat{D}_{t+1}, \hat{M}_{t+1}). \quad (2)$$

Here, \hat{I}_{t+1} and \hat{D}_{t+1} are the result of rendering the image I_t and disparity D_t from the new camera P_{t+1} , using a differentiable renderer \mathcal{R} [13]. This function also returns a mask \hat{M}_{t+1} indicating which regions of the image are missing and need to be filled in. The refinement network g_θ then inpaints, outpaints and super-resolves these inputs to produce the next frame I_{t+1} and its disparity D_{t+1} . The process is repeated iteratively for T steps during training, and at test time for an arbitrarily long camera trajectory. Next we discuss each step in detail.

Geometry and Rendering. Our render step \mathcal{R} uses a differentiable mesh renderer [13]. First, we convert each pixel coordinate (u, v) in I_t and its corresponding disparity d in D_t into a 3D point in the camera coordinate system: $(x, y, z) = K^{-1}(u, v, 1)/d$. We then convert the image into a 3D triangular mesh where each pixel is treated as a vertex connected to its neighbors, ready for rendering.

To avoid stretched triangle artifacts at depth discontinuities and aid our refinement network by identifying regions to be inpainted, we compute a per-pixel binary mask $M_t \in \mathbb{R}^{H \times W}$ by thresholding the gradient of the disparity image $\nabla \hat{D}_t$, computed with a Sobel filter:

$$M_t = \begin{cases} 0 & \text{where } \|\nabla \hat{D}_t\| > \alpha, \\ 1 & \text{otherwise.} \end{cases} \quad (3)$$

We use the 3D mesh to render both image and mask from the new view P_{t+1} , and multiply the rendered image element-wise by the rendered mask to give \hat{I}_{t+1} . The renderer also outputs a depth map as seen from the new camera, which we invert and multiply by the rendered mask to obtain \hat{D}_{t+1} . This use of the mask ensures that any regions in \hat{I}_{t+1} and \hat{D}_{t+1} that were occluded in I_t are masked out and set to zero

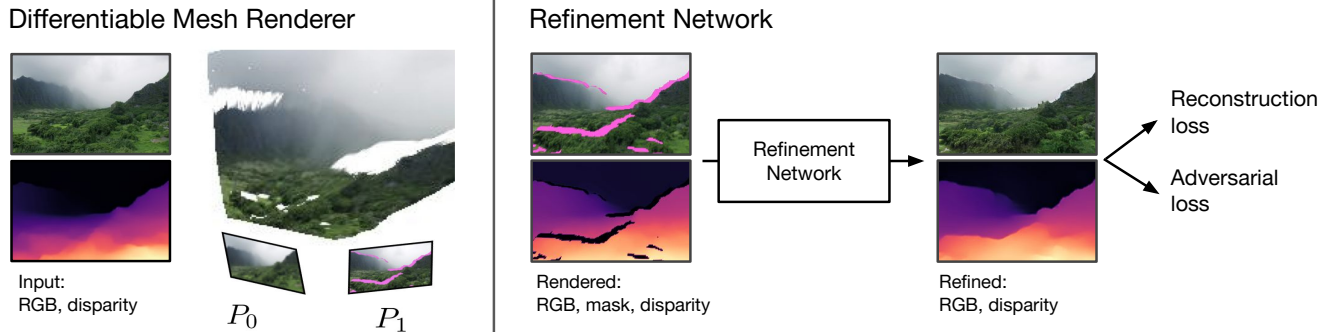


Figure 3. **Illustration of the rendering and refinement steps.** Left: Our differentiable rendering stage takes a paired RGB image and disparity map from viewpoint P_0 and creates a textured mesh representation, which we render from a new viewpoint P_1 , warping the textures, adjusting disparities, and returning a binary mask representing regions to fill in. Right: The refinement stage takes the output of the renderer and uses a deep network to fill in holes and add details. The output is a new RGB image and disparity map that can be supervised with reconstruction and adversarial losses.

(along with regions that were outside the field of view of the previous camera). These areas are ones that the refinement step will have to inpaint (or outpaint). See Figures 2 and 3 for examples of missing regions shown in pink.

Refinement and Synthesis. Given the rendered image \hat{I}_{t+1} , its disparity \hat{D}_{t+1} and its mask \hat{M}_{t+1} , our next task is to refine this image, which includes blurry regions and missing pixels. In contrast to prior inpainting work [49, 36], the refinement network also has to perform super-resolution and thus we cannot use a compositing operation in refining the rendered image. Instead we view the refinement step as a generative image-to-image translation task, and adopt the state-of-the-art SPADE network architecture [26] for our g_θ , which directly outputs I_{t+1}, D_{t+1} . We encode I_0 to provide the additional GAN noise input required by this architecture. See the supplementary materials for more details.

Rinse and Repeat. The previous steps allow us to generate a single novel view. A crucial aspect of our approach is that we refine not only RGB but also disparity, so that scene geometry is propagated between frames. With this setup, we can use the refined image and disparity as the next input to train in an auto-regressive manner, with losses backpropagated over multiple steps. Other view synthesis methods, although not designed in this manner, may also be trained and evaluated in a recurrent setting, although naively repeating these methods without propagating the geometry as we do requires the geometry to be re-inferred from scratch in every step. As we show in Section 6, training and evaluating these baselines with a *repeat* step is still insufficient for perpetual view generation.

Geometric Grounding to Prevent Drift. A key challenge in generating long sequences is dealing with the accumulation of errors [28]. In a system where the current prediction affects future outputs, small errors in each iteration can compound, eventually generating predictions outside the distribution seen during training and causing unexpected behav-

iors. Repeating the generation loop in the training process and feeding the network with its own output ameliorates drift and improves visual quality as shown in our ablation study (Section 6.2). However, we notice that the disparity in particular can still drift at test time, especially over time horizons far longer than seen during training. Therefore we add an explicit geometric re-grounding of the disparity maps.

Specifically, we take advantage of the fact that the rendering process provides the correct range of disparity from a new viewpoint \hat{D}_{t+1} for visible regions of the previous frame. The refinement network may modify these values as it refines the holes and blurry regions, which can lead to drift as the overall disparity becomes gradually larger or smaller than expected. However, we can geometrically correct this by rescaling the refined disparity map to the correct range by computing a scale factor γ via solving

$$\min_{\gamma} \|M \odot (\log(\gamma D_{t+1}) - \log(\hat{D}_{t+1}))\| \quad (4)$$

By scaling the refined disparity by γ , our approach ensures that the disparity map stays at a consistent scale, which significantly reduces drift at test time as shown in Section 6.3.

4. Aerial Coastline Imagery Dataset

Learning to generate long sequences requires real image sequences for training. Many existing datasets for view synthesis do not use sequences, but only a set of views from slightly different camera positions. Those that do have sequences are limited in length: RealEstate10K, for example, has primarily indoor scenes with limited camera movement [52]. To obtain long sequences with a moving camera and few dynamic objects, we turn to aerial footage of beautiful nature scenes available on the Internet. Nature scenes are a good starting point for our challenging problem, as GANs have shown promising results on nature textures [30, 33]. We collected 765 videos using keywords such as ‘coastal’ and ‘aerial footage’, and processed these videos with SLAM

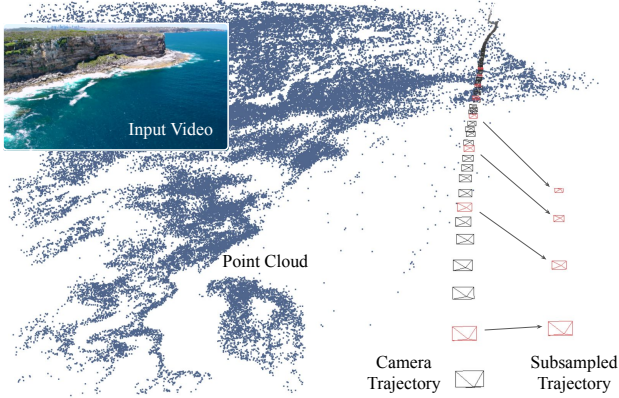


Figure 4. **Processing video for ACID.** We run structure from motion on coastline drone footage collected from YouTube to create the Aerial Coastline Imagery Dataset (ACID). See Section 4.

and structure from motion following the approach of Zhou *et al.* [52], yielding over 13,000 sequences with a total of 2.1 million frames. We have released the list of videos and SfM camera trajectories. See Fig. 4 for an illustrative example of our SfM pipeline running on a coastline video.

To obtain disparity maps for every frame, we use the off-the-shelf MiDaS single-view depth prediction method [27]. We find that MiDaS is quite robust and produces sufficiently accurate disparity maps for our method. Because MiDaS disparity is only predicted up to scale and shift, it must first be rescaled to match our data. To achieve this, we use the sparse point-cloud computed for each scene during structure from motion. For each frame we consider only the points that were tracked in that frame, and use least-squares to compute the scale and shift that minimize the disparity error on these points. We apply this scale and shift to the MiDaS output to obtain disparity maps (D_i) that are scale-consistent with the SfM camera trajectories (P_i) for each sequence.

Due to the difference in camera motions between videos, we strategically sub-sample frames to ensure consistent camera speed in training sequences. See more details in the supplementary materials.

5. Experimental Setup

Losses. We train our approach on a collection of image sequences $\{I_t\}_{t=0}^T$ with corresponding camera poses $\{P_t\}_{t=0}^T$ and disparity maps for each frame $\{D_t\}_{t=0}^T$. Following the literature on conditional generative models, we use an L1 reconstruction loss on RGB and disparity, a VGG perceptual loss on RGB [18] and a hinge-based adversarial loss with a discriminator (and feature matching loss) [26] for the T frames that we synthesize during training. We also use a KL-divergence loss [21] on our initial image encoder $\mathcal{L}_{\text{KLD}} = \mathcal{D}_{\text{KL}}(q(z|x) || \mathcal{N}(0, 1))$. Our complete loss function

is

$$\mathcal{L} = \mathcal{L}_{\text{reconst}} + \mathcal{L}_{\text{perceptual}} + \mathcal{L}_{\text{adv}} + \mathcal{L}_{\text{feat matching}} + \mathcal{L}_{\text{KLD}} \quad (5)$$

The loss is computed over all iterations and over all samples in the mini-batch.

Metrics. Evaluating the quality of the generated images in a way that correlates with human judgement is a challenge. We use the Fréchet inception distance (FID), a common metric used in evaluating generative models of images. FID computes the difference between the mean and covariance of the embedding of real and fake images through a pre-trained Inception network [17] to measure the realism of the generated images as well as their diversity. We precompute real statistics using 20k real image samples from our dataset. To measure changes in generated quality over time, we report FID over a sliding window: we write $\text{FID-}w$ at t to indicate a FID value computed over all image outputs within a window of width w centered at time t , i.e. $\{I_i\}$ for $t - w/2 < i \leq t + w/2$. For short-range trajectories where ground truth images are available, we also report mean squared error (MSE) and LPIPS [51], a perceptual similarity metric that correlates better with human perceptual judgments than traditional metrics such as PSNR and SSIM.

Implementation Details. We train our model with $T = 5$ steps of render-refine-repeat at an image resolution of 160×256 (as most aerial videos have a 16:9 aspect ratio). The choice of T is limited by both memory and available training sequence lengths. The refinement network architecture is the same as that of SPADE generator in [26], and we also employ the same multi-scale discriminator. We implement our models in TensorFlow, and train with a batch size of 4 over 10 GPUs for 7M iterations, which takes about 8 days. We then identify the model checkpoint with the best FID score over a validation set.

6. Evaluation

We compare our approach with three recent state-of-the-art single-image view synthesis methods—the 3D Photography method (henceforward ‘3D Photos’) [32], SynSin [45], and single-view MPIs [38]—as well as the SVG-LP video synthesis method [10]. We retrain each method on our ACID training data, with the exception of 3D Photos which is trained on in-the-wild imagery and, like our method, takes MiDaS disparity as an input. SynSin and single-view MPI were trained at a resolution of 256×256 . SVG-LP takes two input frames for context, and operates at a lower resolution of 128×128 .

The view synthesis baseline methods were not designed for long camera trajectories; every new frame they generate comes from the initial frame I_0 even though after enough camera movement there may be very little overlap between

Method	Over frames 1–10		Over frames 1–50
	LPIPS ↓	MSE ↓	FID ↓
<i>Baseline methods</i>			
SVG-LP [10]	0.60	0.020	135.9
SynSin [45]	0.32	0.018	98.1
MPI [38]	0.35	0.019	65.0
3D Photos [32]	0.30	0.020	123.6
<i>Applied iteratively at test time</i>			
SynSin–Iterated	0.40	0.021	143.6
MPI–Iterated	0.47	0.020	201.2
<i>Trained with repeat ($T = 5$)</i>			
SynSin–Repeat	0.44	0.036	153.3
MPI–Repeat	0.55	0.020	203.0
Ours	0.32	0.020	50.6

Table 1. **Quantitative evaluation.** We compute LPIPS and MSE against ten frames of ground truth, and FID-50 over 50 frames generated from an input test image. See Section 6.1.

the two. Therefore we also compare against two variants of each of these methods. First, variants with *iterated evaluation* (Synsin–Iterated, MPI–Iterated): these methods use the same trained models as their baseline counterparts, but we apply them iteratively at test time to generate each new frame from the *previous* frame rather than the initial one. Second, variants *trained with repeat* (Synsin–Repeat, MPI–Repeat): these methods are trained autoregressively, with losses backpropagated across $T = 5$ steps, as in our full model. (We omit these variations for the 3D Photos method, which was unfortunately too slow to allow us to apply it iteratively, and which we are not able to retrain.)

6.1. Short-to-medium range view synthesis

To evaluate short-to-medium-range synthesis, we select ACID test sequences with an input frame and 10 subsequent ground truth frames (subsampling as described in the supplementary), with the camera moving forwards at an angle of up to 45° . Although our method is trained on all types of camera motions, this forward motion is appropriate for comparison with view synthesis methods which are not designed to handle extreme camera movements.

We then extrapolate the camera motion from the last two frames of each sequence to extend the trajectory for an additional 40 frames. To avoid the camera colliding with the scene, we check the final camera position against the disparity map of the last ground-truth frame, and discard sequences in which it is outside the image or at a depth large enough to be occluded by the scene.

This yields a set of 279 sequences with camera trajectories of 50 steps and ground truth images for the first 10 steps. For short-range evaluation, we compare to ground truth on the first 10 steps. For medium-range evaluation, we compute FID scores over all 50 frames.

We apply each method to these sequences to generate novel views corresponding to the camera poses in each se-

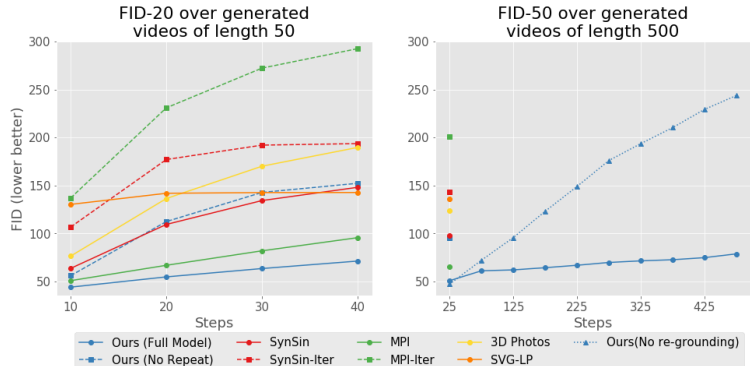


Figure 5. **FID over time.** Left: FID-20 over time for 50 frames generated by each method. Right: FID-50 over 500 frames generated by our method via autopilot. For comparison, we plot FID-50 for the baselines on the first 50 steps. Despite generating sequences an order of magnitude longer, our FID-50 is still lower than that of the baselines. See Sec. 6.1, 6.3.

quence (SVG-LP is the exception in that it does not take account of camera pose). See results in Table 1. While our goal is perpetual view generation, we find that our approach is competitive with recent view synthesis approaches for short-range synthesis on LPIPS and MSE metrics. For mid-range evaluation, we report FID-50 over 50 generated frames. Our approach has a dramatically lower FID-50 score than other methods, reflecting the more naturalistic look of its output. To quantify the degradation of each method over time, we report a sliding window FID-20 computed from $t = 10$ to 40. As shown in Fig. 5 (left), the image quality (measured by FID-20) of the baseline methods deteriorates more rapidly with increasing t compared to our approach.

Qualitative comparisons of these methods are shown in Fig. 6 and our supplementary video, which illustrates how the quality of each method’s output changes over time. Notable here are SVG-LP’s blurriness and inability to predict any camera motion at all; the increasingly stretched textures of 3D Photos’ output; and the way the MPI-based method’s individual layers become noticeable. SynSin does the best job of generating plausible texture, but still produces holes after a while and does not add new detail.

The –Iterated and –Repeat variants are consistently worse than the original SynSin and MPI methods, which suggests that simply applying an existing method iteratively, or re-training it autoregressively, is insufficient to deal with large camera movement. These variants show more drifting artifacts than their original versions, likely because (unlike our method), they do not propagate geometry from step to step. The MPI methods additionally become very blurry on repeated application, as they have no ability to add detail, lacking our refinement step.

In summary, our thoughtful combination of render-refine-repeat shows better results than these existing methods and variations. Figure 7 shows additional qualitative results from generating 15 and 30 frames using on a variety of inputs.

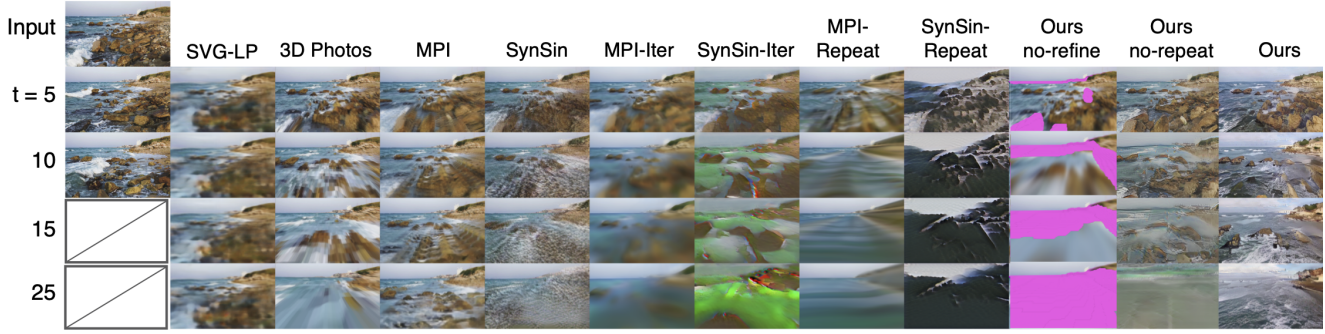


Figure 6. **Qualitative comparison over time.** We show a generated sequence for each method at different time steps. Note that we only have ground truth images for 10 frames; the subsequent frames are generated using an extrapolated trajectory. Pink region in Ours no-refine indicate missing content uncovered by the moving camera.

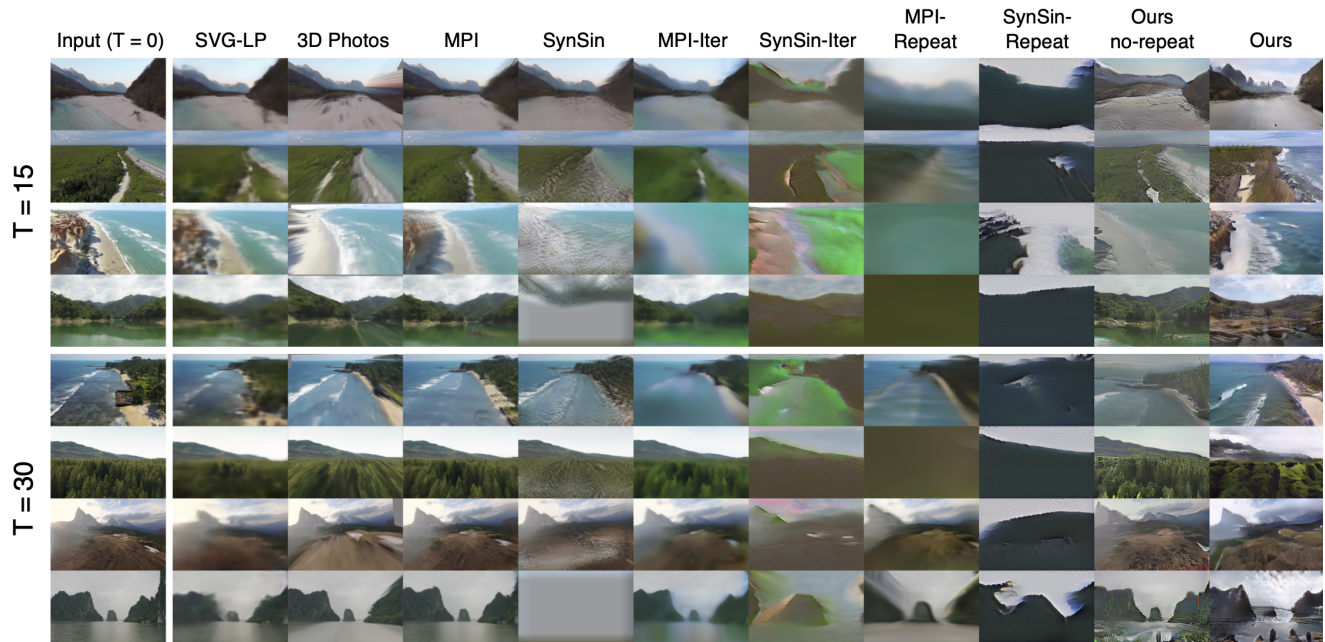


Figure 7. **Qualitative comparison.** We show the diversity and quality of many generated results for each method on the $t=15$ and 30 frame generation. Competing approaches result in missing or unrealistic frames. Our approach is able to generate plausible views of the scene.

6.2. Ablations

We investigate the benefit of training over multiple iterations of our *render-refine-repeat* loop by also training our model with $T = 1$ (‘No repeat’). As shown in in Table 2, the performance on short-range generation, as measured in LPIPS and MSE, is similar to our full model, but when we look at FID, we observe that this method generates lower quality images and that they get substantially worse with increasing t (see Fig. 5, left). This shows the importance of using a recurrent training setup to our method.

We next consider the *refine* step. Omitting this step completely results in a larger and larger portion of the image being completely missing as t increases: examples are shown as ‘Ours (no refine)’ in Fig. 6, where for clarity the miss-

ing pixels are highlighted in pink. In our full model, these regions are inpainted or outpainted by the refinement network at each step. Note also that even non-masked areas of the image are much blurrier when the refinement step is omitted, showing the benefit of the refinement network in super-resolving image content.

Table 2 also shows results on two further variations of our refinement step. First, replacing our refinement network with a simpler U-Net architecture yields substantially worse results (‘U-Net refinement’). Second, disabling geometric grounding (Section 3.1) also leads to slightly lower quality on this short-to-medium range view synthesis task (‘No re-grounding’).

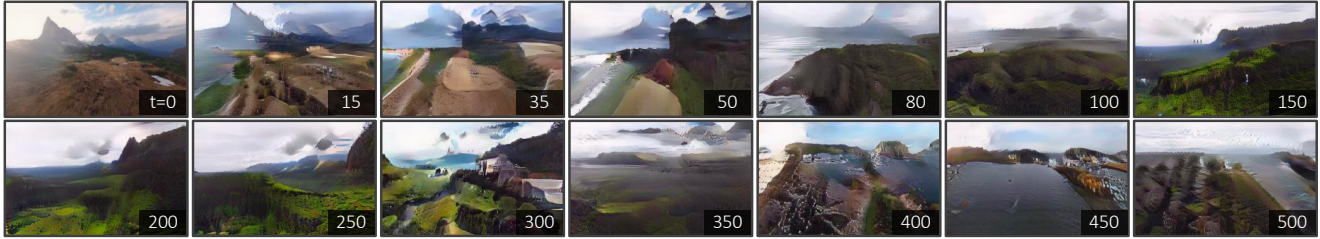


Figure 8. **Long trajectory generation.** From a single image, our approach can generate 500 frames of video without suffering visually. Please see the supplementary video for the full effect.

Ablations	LPIPS ↓	MSE ↓	FID-50 ↓
Full Model	0.32	0.020	50.6
No repeat ($T = 1$)	0.30	0.022	95.4
U-net Refinement	0.54	0.052	183.0
No re-grounding	0.34	0.022	64.3

Table 2. **Ablations.** We ablate aspects of our model to understand their contribution to the overall performance. See Section 6.2.

6.3. Perpetual view generation

We also evaluate the ability of our model to perform perpetual view generation by synthesizing videos of 500 frames, using an *auto-pilot* algorithm to create an online camera trajectory that avoids flying directly into the ground, sky or obstacles such as mountains. This algorithm works iteratively in tandem with image generation to control the camera based on heuristics which measure the proportion of sky and of foreground obstacles in the scene. See the supplementary materials for details.

We note that this task is exceptionally challenging and completely outside the capabilities of current generative and view synthesis methods. To further frame the difficulty, our refinement network has only seen videos of length 5 during training, yet we generate 500 frames for each of our test sequences. As shown in Fig. 5 (right), our FID-50 score over generated frames is remarkably robust: even after 500 frames, the FID is lower than that of all the baseline methods over 50 frames. Fig. 5 also shows the benefit of our proposed geometric grounding: when it is omitted, the image quality gradually deteriorates, indicating that resolving drift is an important contribution.

Fig. 8 shows a qualitative example of long sequence generation. In spite of the intrinsic difficulty of generating frames over large distances, our approach retains something of the aesthetic look of coastline, generating new islands, rocks, beaches, and waves as it flies through the world. The auto-pilot algorithm can receive additional inputs (such as a user-specified trajectory or random elements), allowing us to generate diverse videos from a single image. Please

see the supplementary video for more examples and the full effect of these generated fly-through videos.

6.4. User-controlled video generation

Because our rendering step takes camera poses as an input, we can render frames for arbitrary camera trajectories at test time, including trajectories controlled by a user in the loop. We have built a HTML interface that allows the user to steer our auto-pilot algorithm as it flies through this imaginary world. This demo runs over the internet and is capable of generating a few frames per second. Please see the supplementary video for a demonstration.

7. Discussion

We introduce the new problem of perpetual view generation and present a novel framework that combines both geometric and generative techniques as a first step in tackling it. Our system can generate video sequences spanning hundreds of frames, which to our knowledge has not been shown for prior video or view synthesis methods. The results indicate that our hybrid approach is a promising step. Nevertheless, many challenges remain.

First, our render-refine-repeat loop is by design memory-less, an intentional choice which allows us to train on finite length videos yet generate arbitrarily long output using a finite memory and compute budget. As a consequence it aims for local consistency between nearby frames, but does not directly tackle questions of long-term consistency or a global representation. How to incorporate long-term memory in such a system is an exciting question for future work. Second, our refinement network, like other GANs, can produce images that seem realistic but not recognizable [16]. Further advancements in image and video synthesis generation methods that incorporate geometry would be an interesting future direction. Last, we do not model dynamic scenes: combining our geometry-aware approach with methods that can reason about object dynamics could be another fruitful direction.

Acknowledgements. We would like to thank Dominik Kaeser for directing and helping prepare our videos and Huiwen Chang for making the MiDaS models easily accessible.

References

- [1] Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. PatchMatch: A randomized correspondence algorithm for structural image editing. *ACM Transactions on Graphics (Proc. SIGGRAPH)*, 28(3), Aug. 2009. 2
- [2] Moshe Blank, Lena Gorelick, Eli Shechtman, Michal Irani, and Ronen Basri. Actions as space-time shapes. In *ICCV*, pages 1395–1402. IEEE, 2005. 2
- [3] Gaurav Chaurasia, Sylvain Duchêne, Olga Sorkine-Hornung, and George Drettakis. Depth synthesis and local warps for plausible image-based navigation. *Trans. on Graphics*, 32:30:1–30:12, 2013. 2
- [4] Weifeng Chen, Shengyi Qian, and Jia Deng. Learning single-image depth from videos using quality assessment networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2
- [5] Xu Chen, Jie Song, and Otmar Hilliges. Monocular neural image based rendering with continuous view control. In *ICCV*, 2019. 2
- [6] Xu Chen, Jie Song, and Otmar Hilliges. Monocular neural image based rendering with continuous view control. In *ICCV*, pages 4090–4100, 2019. 2
- [7] Inchang Choi, Orazio Gallo, Alejandro Troccoli, Min H Kim, and Jan Kautz. Extreme view synthesis. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7781–7790, 2019. 2
- [8] Inchang Choi, Orazio Gallo, Alejandro Troccoli, Min H Kim, and Jan Kautz. Extreme view synthesis. In *ICCV*, pages 7781–7790, 2019. 2
- [9] Aidan Clark, Jeff Donahue, and Karen Simonyan. Efficient video generation on complex datasets. *arXiv preprint arXiv:1907.06571*, 2019. 1, 2
- [10] Emily Denton and Rob Fergus. Stochastic video generation with a learned prior. *arXiv preprint arXiv:1802.07687*, 2018. 1, 2, 5, 6
- [11] Chelsea Finn, Ian Goodfellow, and Sergey Levine. Unsupervised learning for physical interaction through video prediction. In *NeurIPS*, pages 64–72, 2016. 2
- [12] 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 *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2
- [13] Kyle Genova, Forrester Cole, Aaron Maschinot, Aaron Sarna, Daniel Vlasic, and William T. Freeman. Unsupervised training for 3d morphable model regression. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 3
- [14] Daniel Glasner, Shai Bagon, and Michal Irani. Super-resolution from a single image. In *ICCV*, pages 349–356, 2009. 2
- [15] James Hays and Alexei A Efros. Scene completion using millions of photographs. *ACM Transactions on Graphics (TOG)*, 26(3):4–es, 2007. 2
- [16] Aaron Hertzmann. Visual indeterminacy in generative neural art. *arXiv preprint arXiv:1910.04639*, 2019. 8
- [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. In *NeurIPS*, pages 6626–6637, 2017. 5
- [18] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016. 5
- [19] 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
- [20] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019. 1
- [21] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 5
- [22] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew P Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *CVPR*, 2017. 2
- [23] Marc Levoy and Pat Hanrahan. Light field rendering. In *Proceedings of SIGGRAPH 96*, Annual Conference Series, 1996. 2
- [24] 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)*, 2019. 2
- [25] Simon Niklaus, Long Mai, Jimei Yang, and Feng Liu. 3D Ken Burns effect from a single image. *ACM Transactions on Graphics (TOG)*, 2019. 2
- [26] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 4, 5
- [27] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 2, 5
- [28] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 627–635, 2011. 4
- [29] Arno Schödl, Richard Szeliski, David H Salesin, and Irfan Essa. Video textures. In *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, pages 489–498, 2000. 2
- [30] Tamar Rott Shaham, Tali Dekel, and Tomer Michaeli. Singan: Learning a generative model from a single natural image. In *ICCV*, pages 4570–4580, 2019. 2, 4
- [31] Lixin Shi, Haitham Hassanieh, Abe Davis, Dina Katabi, and Fredo Durand. Light field reconstruction using sparsity in the

- continuous fourier domain. *Trans. on Graphics*, 34(1):12:1–12:13, Dec. 2014. 2
- [32] Meng-Li Shih, Shih-Yang Su, Johannes Kopf, and Jia-Bin Huang. 3d photography using context-aware layered depth inpainting. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 5, 6
- [33] Assaf Shocher, Shai Bagon, Phillip Isola, and Michal Irani. Ingan: Capturing and remapping the “DNA” of a natural image. *arXiv preprint arXiv:1812.00231*, 2018. 2, 4
- [34] Assaf Shocher, Nadav Cohen, and Michal Irani. “zero-shot” super-resolution using deep internal learning. In *CVPR*, pages 3118–3126, 2018. 2
- [35] 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 *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2
- [36] Piotr Teterwak, Aaron Sarna, Dilip Krishnan, Aaron Maschinot, David Belanger, Ce Liu, and William T Freeman. Boundless: Generative adversarial networks for image extension. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 10521–10530, 2019. 2, 4
- [37] Alex Trevithick and Bo Yang. Grf: Learning a general radiance field for 3d scene representation and rendering. In *arXiv:2010.04595*, 2020. 2
- [38] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2, 5, 6
- [39] Shubham Tulsiani, Richard Tucker, and Noah Snavely. Layer-structured 3D scene inference via view synthesis. In *The European Conference on Computer Vision (ECCV)*, September 2018. 2
- [40] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. Mocogan: Decomposing motion and content for video generation. In *CVPR*, pages 1526–1535, 2018. 2
- [41] Ruben Villegas, Arkanath Pathak, Harini Kannan, Dumitru Erhan, Quoc V Le, and Honglak Lee. High fidelity video prediction with large stochastic recurrent neural networks. In *Advances in Neural Information Processing Systems*, pages 81–91, 2019. 1, 2
- [42] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. In *NeurIPS*, pages 613–621, 2016. 2
- [43] Carl Vondrick and Antonio Torralba. Generating the future with adversarial transformers. In *CVPR*, pages 1020–1028, 2017. 2
- [44] Yi Wang, Xin Tao, Xiaoyong Shen, and Jiaya Jia. Wide-context semantic image extrapolation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1399–1408, 2019. 2
- [45] Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. SynSin: End-to-end view synthesis from a single image. In *CVPR*, 2020. 2, 5, 6
- [46] Zongxin Yang, Jian Dong, Ping Liu, Yi Yang, and Shuicheng Yan. Very long natural scenery image prediction by outpainting. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 10561–10570, 2019. 2
- [47] Yufei Ye, Maneesh Singh, Abhinav Gupta, and Shubham Tulsiani. Compositional video prediction. In *ICCV*, 2019. 2
- [48] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. 2
- [49] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Generative image inpainting with contextual attention. *arXiv preprint arXiv:1801.07892*, 2018. 4
- [50] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Free-form image inpainting with gated convolution. In *ICCV*, 2019. 2
- [51] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 5
- [52] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. *ACM Trans. Graph.*, 37(4):65:1–65:12, 2018. 2, 3, 4, 5
- [53] Yang Zhou, Zhen Zhu, Xiang Bai, Dani Lischinski, Daniel Cohen-Or, and Hui Huang. Non-stationary texture synthesis by adversarial expansion. *arXiv preprint arXiv:1805.04487*, 2018. 2