

# Worldsheet: Wrapping the World in a 3D Sheet for View Synthesis from a Single Image

Ronghang Hu  
Facebook AI Research  
ronghanghu@fb.com

Deepak Pathak  
Carnegie Mellon University  
dpathak@cs.cmu.edu

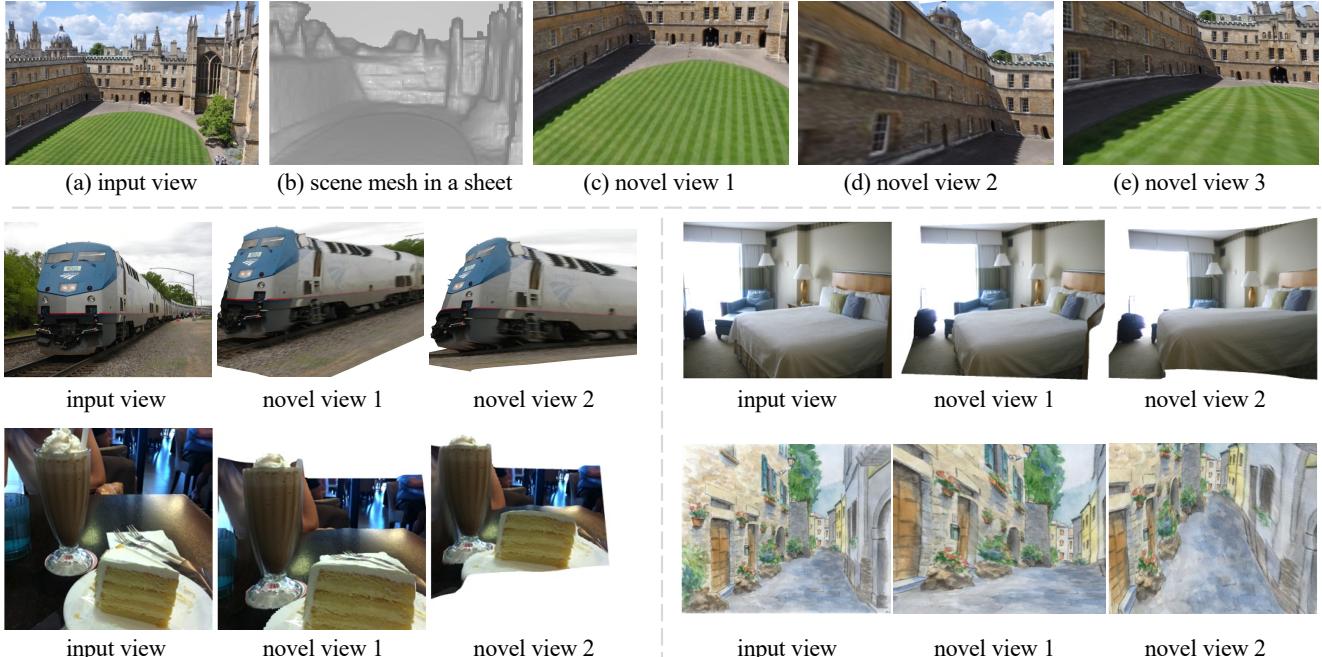


Figure 1: We synthesize novel views from large view-point changes given a single input RGB image (shown in a) by wrapping a mesh sheet (shown in b) onto the image, and rendering it from novel view-points (shown in c, d, e). Plausible novel views are generated for outdoor scenes, outdoor objects, indoor scenes, indoor objects and even paintings with high resolution ( $960 \times 960$ ) input. Please see continuously synthesized views at [worldsheet.github.io](https://worldsheet.github.io) (Image sources: [44, 11, 22, 25]).

## Abstract

We present *Worldsheet*, a method for novel view synthesis using just a single RGB image as input. This is a challenging problem as it requires an understanding of the 3D geometry of the scene as well as texture mapping to generate both visible and occluded regions from new viewpoints. Our main insight is that simply shrink-wrapping a planar mesh sheet onto the input image, consistent with the learned intermediate depth, captures underlying geometry sufficient enough to generate photorealistic unseen views with arbitrarily large view-point changes. To operationalize this, we propose a novel differentiable texture sampler that allows our wrapped mesh sheet to be textured; which is then transformed into a target image via differentiable

rendering. Our approach is category-agnostic, end-to-end trainable without using any 3D supervision and requires a single image at test time. *Worldsheet* consistently outperforms prior state-of-the-art methods on single-image view synthesis across several datasets. Furthermore, this simple idea captures novel views surprisingly well on a wide range of high resolution in-the-wild images in converting them into a navigable 3D pop-up. Video results and code at <https://worldsheet.github.io>

## 1. Introduction

A 2D image is the projection of an underlying 3D world, but we as humans have no trouble in understanding this structure and imagining how an image will look from other

views. Consider the train shown in Figure 1, the reason we can seamlessly predict its other views just from a single image is because of the abstractions we have learned from past experience of watching several trains, or similar shaped objects from different views. Enabling machines with such an ability to reason about 3D from a single image will bring trillions of still photos to life with several applications in virtual reality, animation, image editing, and robotics.

This goal of synthesizing novel views from 2D images has been pursued for decades, from early efforts relying completely on geometry [7, 54, 35] to more recent learning based approaches [52, 40, 1, 27, 23, 32]. Over the years, significant progress has been made in this direction. However, although generating impressive photorealistic output views, most of these previous approaches require multiple images or ground-truth depth at test time, which severely hinders their practicality as a general tool.

On the other hand, the holy grail of developing a general-purpose tool to convert any arbitrary single image into a corresponding 3D structure for generating novel views has kept the community hooked for years. To compensate for the lack of multiple views or 3D models at test time, methods on single-image 3D rely on statistical learning from data. This line of work can be traced back to classic works of Hoiem *et al.* [11], followed by Saxena *et al.* [34], to obtain ‘qualitative 3D’ from a single image by marrying geometry and learning to see images as a collection of planes.

An ideal approach to general-purpose view synthesis should not only rely on a single image at test time, but also learn from easy-to-collect supervision at training. In the deep learning era, there is growing interest in end-to-end methods with intermediate 3D representations supervised by multiple images and no explicit 3D information during training. However, they are mostly applied to objects [18, 42, 48, 13, 39] and not scenes, and are either category-specific, applied to synthetic scenes, or both. Recent works [5, 47, 43] alleviate this issue by training with multiple views of real-world scenes, relying on point cloud or multiplane images as intermediate representations. However, multiplane images are mostly good for small view-point changes as each plane is at a constant depth, and for point cloud, one needs to represent each point in a scene individually, making it inefficient to scale to high-res data or large view-point changes. In contrast, meshes provide a sparser scene representation, *e.g.*, a single mesh face can theoretically represent the whole flat surface of a wall, making it ideal for single-image view synthesis. However, mesh recovery from single images has been studied mostly for object images and in a category-specific manner [2, 13, 17].

In this paper, we present an end-to-end approach for novel view synthesis from *a single image* via an intermediate *mesh representation*. Unlike mesh reconstruction for objects of certain categories, generating meshes for a scene

is challenging as there is no notion of mean or canonical shape to start from, or silhouette from segmentation for supervision. We side-step this problem by wrapping by a deformable mesh sheet over the 3D world – much like wrapping a 2D tinfoil onto a 3D pan before baking! We name this shrink-wrapped mesh *Worldsheet*, a term borrowed from physics for the 2D manifold of high-dimensional strings. After generating this *Worldsheet* for a given view, novel views are obtained by moving the camera in 3D geometry space (Figure 2), which allows us to train from just two views of a scene *without any 3D or depth supervision*.

To train our model end-to-end, reconstructing mesh texture from inputs and rendering it in novel views both need to be differentiable. Latter is easily handled thanks to recent differentiable mesh renderers [14, 24, 31]. To address the former, we propose a differentiable texture sampler over projected 2D views, enabling gradient computation of the reconstructed texture map over the 3D mesh geometry.

*Worldsheet* allows generating novel views from a single image with large view-point changes. It consistently outperforms prior state-of-the-art by a significant margin on three benchmark datasets (Matterport [3], Replica [41], and RealEstate10K [52]) and is applicable to very high-resolution images in-the-wild as shown in Figure 1.

## 2. Related work

**Novel view synthesis from multiple images.** Traditionally, novel view synthesis used multiple input views at test time [4, 20, 9], and are often based on different representations. Among recent works, Waechter *et al.* [45] build scene meshes with diffuse appearance. StereoMag [52] proposes multiplane images (MPIs) from a stereo image pair as a layered scene representation. NPBG [1] captures the scene as a point cloud with neural descriptors. NeRF [27] proposes a neural radiance field representation for scene appearance. NSVF [23] adopts sparse voxel octrees as scene representations. FVS [32] uses recurrent mapping and blending based on depth map of the scene from multi-view stereo. Yoon *et al.* [49] combine depth from both single and multiple views to generate novel views of dynamic scenes. Access to multiple input views greatly simplifies the task, allowing the scene geometry to be recovered via multi-view stereo [35].

**Novel view synthesis from a single image.** In early works, Debevec *et al.* [7] recover 3D scene models and Horry *et al.* [12] fit a regular mesh to generate novel views. Liebowitz *et al.* [21] and Criminsi *et al.* [6] generate meshes via projective geometry constraints. These methods came at the expense of manual editing until Hoiem *et al.* [11] generate automatic 3D pop-up from 2D images based on vertical and ground planes. More recently in [29, 16, 36], layered depth images are used for single image view synthesis based on a pretrained depth estimator. In [43], online videos

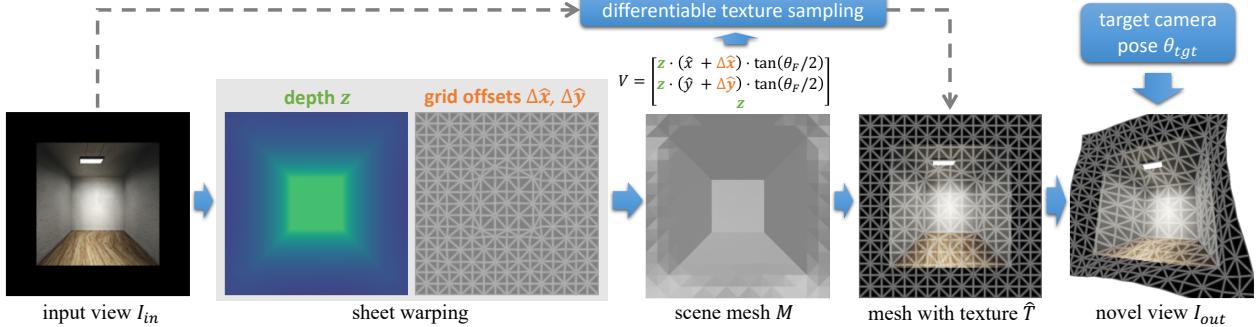


Figure 2: An overview of our *Worldsheet* approach. Given an input view  $I_{in}$ , we build a scene mesh by warping a  $W_m \times H_m$  grid sheet onto the scene geometry via grid offset ( $\Delta\hat{x}, \Delta\hat{y}$ ) and depth  $z$  (Sec. 3.1). Then, we sample the UV texture map  $\hat{T}$  of the scene mesh differentiably (Sec. 3.2) and render it from the target camera pose to output a novel view  $I_{out}$ . Our mesh warping is learned end-to-end using the losses on the novel view (Sec. 3.3). We further apply an inpainting network over  $I_{out}$  (not shown above; see Sec. 3.4) to inpaint invisible regions and refine image details, outputting a refined novel view  $I_{paint}$ . We train with two views without any 3D or depth supervision and require just a single RGB image at test time.

are used to train scale-invariant MPI representation for view synthesis. SynSin [47] synthesizes novel views from a single image with a feature point cloud. We follow the same experimental setup as SynSin but learns to construct scene meshes instead of point clouds as scene representations to generate novel views from large viewpoint changes.

**Differentiable mesh rendering.** Recent work on differentiable mesh renders [14, 24, 31] allow learning 3D structures through synthesis. NMR [14] and SoftRas [24] reconstruct the 3D object shape as a mesh by rendering it, comparing it with the input image, and back-propagating losses to refine the mesh geometry. CMR [13], CSM [17] and U-CMR [8] build category-specific object meshes from images by deforming from a mean or template category shape through silhouette (and keypoints in [13]) supervision.

Our method is aligned with the analysis-by-synthesis paradigm above. However, unlike most previous works that apply differentiable mesh rendering to objects, we learn the 3D geometry of *scenes* through the rendering losses on the novel view. Moreover, instead of predicting a texture flow as in [13, 8], we propose to analytically sample the mesh texture from the input view with a differentiable texture sampler. Unlike [30], our differentiable texture sampler considers multiple mesh faces in the z-buffer (soft rasterization instead of only the closest one), and assumes perspective (instead of orthographic) camera projection.

### 3. Worldsheet: Rendering the World in a Sheet

In this work, we propose *Worldsheet* to synthesize novel views from a single image, as shown in Figure 2. Given an input image  $I_{in}$ , we build a 3D scene mesh  $M$  by warping a lattice grid (*i.e.* a ‘sheet’) onto the scene geometry through depth and grid offset, and sample the mesh’s UV texture map in a differentiable manner. We then render the mesh from the target camera pose to output the novel view image  $I_{out}$ . We further apply an inpainting network over the

rendered mesh to imagine the invisible regions, outputting a 2D image  $I_{paint}$  as the final target view prediction. Our model is trained with only losses over the 2D output image of the target view, *without any 3D or depth supervision*.

#### 3.1. Scene mesh prediction by warping a sheet

From the input view image  $I_{in}$  of size  $W_{im} \times H_{im}$ , we build a scene mesh by warping a  $W_m \times H_m$  lattice grid (*i.e.* a sheet) onto the scene, as shown in Figure 2. We first extract a  $W_m \times H_m$  visual feature map  $\{q_{w,h}\}$  from  $I_{in}$  with a convolutional neural network. Each  $q_{w,h}$  is a feature vector at spatial location  $(w, h)$  on the  $W_m \times H_m$  network output. In our implementation, we use ResNet-50 [10] (pretrained on ImageNet) with dilation [50] to output features  $\{q_{w,h}\}$ .

From each  $q_{w,h}$  on the feature map, we predict the grid offset  $\Delta\hat{x}_{w,h}$  and  $\Delta\hat{y}_{w,h}$  to decide how much the vertex  $(w, h)$  on the grid should move away from its anchor positions within the image plane (we output  $\Delta\hat{x}_{w,h}$  and  $\Delta\hat{y}_{w,h}$  in NDC space [37] between -1 to 1). We also predict how far each vertex is from the camera, *i.e.* its depth  $z_{w,h}$ . These values are predicted using learned linear mappings as

$$\Delta\hat{x}_{w,h} = \tanh(W_1 q_{w,h} + b_1) / (W_m - 1) \quad (1)$$

$$\Delta\hat{y}_{w,h} = \tanh(W_2 q_{w,h} + b_2) / (H_m - 1) \quad (2)$$

$$z_{w,h} = g(W_3 q_{w,h} + b_3) \quad (3)$$

where division by  $(W_m - 1)$  and  $(H_m - 1)$  ensures that the vertices can only move within a certain range.  $g(\cdot)$  is a scalar nonlinear function to scale the network prediction into depth values. We use  $g(\psi) = \alpha_g / (\sigma(\psi) + \epsilon_g) + \beta_g$  in our implementation, where  $\sigma(\cdot)$  is the sigmoid function and  $\alpha_g, \beta_g$  and  $\epsilon_g$  are fixed hyper-parameters.

**Building the 3D scene mesh.** We first build the mesh vertices  $\{V_{w,h}\}$  from the grid offset and depth as

$$V_{w,h} = \begin{bmatrix} z_{w,h} \cdot (\hat{x}_{w,h} + \Delta\hat{x}_{w,h}) \cdot \tan(\theta_F/2) \\ z_{w,h} \cdot (\hat{y}_{w,h} + \Delta\hat{y}_{w,h}) \cdot \tan(\theta_F/2) \\ z_{w,h} \end{bmatrix} \quad (4)$$

for  $w = 1, \dots, W_m$  and  $h = 1, \dots, H_m$ . Here  $\theta_F$  is the camera field-of-view, and  $\hat{x}_{w,h}$  and  $\hat{y}_{w,h}$  are anchor positions on the grid equally spaced from  $-1$  to  $1$ .

Then, we connect the mesh vertices  $\{V_{w,h}\}$  along the edges on the grid to form mesh faces  $\{F\}$  as shown in Figure 2 and obtain a 3D mesh  $M = (\{V_{w,h}\}, \{F\})$ . A vertex in the mesh is connected to its 4 or 8 neighbours on the grid.

**Regularization.** The predicted scene mesh is regularized with two losses. First, we apply an L2 regularization term  $L_g = \sum_{w,h} (\Delta \hat{x}_{w,h}^2 + \Delta \hat{y}_{w,h}^2)$  to the grid offset to avoid unnecessary large grid deformations. To encourage the mesh surface to be smooth unless it needs to bend to fit the scene geometry, we also apply a Laplacian term  $L_m = \sum_{w,h} \left\| \sum_{(\bar{w},\bar{h}) \in N(w,h)} (V_{\bar{w},\bar{h}} - V_{w,h}) \right\|_1$  on the mesh vertices, where  $N(w, h)$  are the adjacent vertices to  $(w, h)$ .

### 3.2. Differentiable texture sampler

To render the mesh in another camera pose for novel view synthesis, we would like to reconstruct its texture from the scene appearance in the input view. While a few approaches [13, 8] build mesh texture with a learned texture flow on objects, it is hard to apply the same to scenes, which are unlike birds or cats and do not have canonical shapes. In this work, we take an alternative route and propose a differentiable texture sampler. We *analytically* sample the mesh texture  $\hat{T}$  as a UV texture map [37] from the input view  $I_{in}$ , where gradients  $\partial \hat{T} / \partial V$  and  $\partial \hat{T} / \partial I_{in}$  over the vertex coordinates and the input image can be computed.

To implement this texture sampler, we first build a z-buffer containing top  $K$  mesh faces on each image pixel  $(i, j)$  as in PyTorch3D [31], and splat the RGB pixel intensities from the image  $I_{in}$  onto the UV texture map  $\hat{T}$ . Specifically, we compute the blending probability  $p_{i,j}^k$  for the  $k$ -th covering mesh face (sorted in ascending z-order in the z-buffer) on pixel  $(i, j)$  based on the softmax RGB blend in [31] as  $p_{i,j}^k = D_{i,j}^k \cdot \exp(-Z_{i,j}^k) / \sum_{k'} D_{i,j}^{k'} \cdot \exp(-Z_{i,j}^{k'})$ , where  $D_{i,j}^k$  is the soft-rasterization probability map [24] of the  $k$ -th face on the image (nearly 0 outside the face, nearly 1 inside the face, and blurry on the edges; see Figure 4 in Liu *et al.* [24]), and  $Z_{i,j}^k$  is the depth value of the 3D point on the  $k$ -th face covering pixel  $(i, j)$  in the z-buffer.

We then decompose the input image  $I_{in}$  into  $K$  images  $I_{in}^k$ , where  $I_{in}^k(i, j) = I_{in} \cdot p_{i,j}^k$ , and splat the RGB pixel intensities from each  $I_{in}^k$  to a texture map layer  $\hat{T}^k$  as

$$\hat{T}^k = \text{splat}(I_{in}^k, f^k) \quad (5)$$

where the flow  $(u, v) = f^k(i, j)$  maps image coordinates  $(i, j)$  to UV coordinates  $(u, v)$  on the  $k$ -th mesh face in the z-buffer. Here splat is a differentiable splatting operation (based on [28]; see supplemental for details) from the image space  $I_{in}^k$  to the texture space  $\hat{T}^k$ . Finally, we sum all the  $K$  texture maps as the final UV texture map  $\hat{T} = \sum_k \hat{T}^k$ .



Figure 3: Parts of the target view (the gray area in b) are often invisible from the input and must be imagined based on prior knowledge. We make plausible predictions over the invisible regions with an inpaint network (c). However, this task is inherently uncertain (*e.g.* one cannot be sure about the rightmost cabinet in d).

In summary, the image pixels are splatted onto the texture space through each face, and blended together to obtain the final texture map. The entire process is differentiable with respect to both  $I_{in}$  and the mesh vertex coordinates  $\{V\}$ , as one can analytically compute  $\partial \hat{T}^k / \partial I_{in}^k$ ,  $\partial \hat{T}^k / \partial f^k$ ,  $\partial f^k / \partial V$ ,  $\partial D / \partial V$ , and  $\partial Z / \partial V$ .

### 3.3. Learning scene geometry by view synthesis

To synthesize a novel view, we project the mesh vertex coordinates  $\{V\}$  from the input camera pose  $\theta_{in}$  to  $\{V^{tgt}\}$  in the camera coordinate space of the target view-point  $\theta_{tgt}$ . Then, we render the mesh  $M^{tgt} = (\{V^{tgt}\}, \{F\})$  in the target camera pose along with its texture map  $\hat{T}$  to output a 2D image  $I_{out}$  of size  $W_{im} \times H_{im}$  as the target view:

$$I_{out} = \text{render\_mesh}(\{V^{tgt}\}, \{F\}, \hat{T}). \quad (6)$$

We use the differentiable mesh renderer in [31] so that we can compute the gradients  $\partial I_{out} / \partial V^{tgt}$  and  $\partial I_{out} / \partial \hat{T}$ . Through mesh rendering, we also obtain a foreground mask  $F_{out}$  with the same size as  $I_{out}$ , indicating which pixels in the rendered image  $I_{out}$  are covered by the mesh and which pixels are from background color, as shown in Figure 3 (b).

Our model is supervised with paired input and target views of a scene (along with their camera poses). We use a pixel L1 loss  $L_{out}^{rgb} = \|I_{out} - I_{tgt}\|_1 / (W_{im} \cdot H_{im})$  and a perceptual loss [46, 47]  $L_{out}^{pc} = P(I_{out}, I_{tgt})$ , where  $I_{tgt}$  is the ground-truth target view image. The model then needs just a single image at test time.

### 3.4. Inpainting and image refinement

The target view image consists of two parts: things that can be directly seen from the input view  $I_{in}$ , and things that need to be imagined based on our prior knowledge of the visual world, as illustrated in Figure 3. As our mesh warping and rendering procedure in Sec. 3.1, 3.2 and 3.3 builds a pixel-to-pixel correspondence between the input and the target view, it only renders pixels that are *visible* from the input view. To obtain a plausible imagination of the *invisible* image regions, we apply an inpainting network  $G$  on the

#	Method	Matterport [3]						Replica [41]		
		PSNR $\uparrow$			SSIM $\uparrow$			Perc Sim $\downarrow$		
		Both	InVis	Vis	Both	InVis	Vis	Both	InVis	Vis
1	Im2Im [53]	15.87	16.20	15.97	0.53	0.60	0.48	2.99	0.58	2.05
2	Tatarchenko <i>et al.</i> [42]	14.79	14.83	15.05	0.57	0.62	0.53	3.73	0.74	2.50
3	Vox [38] w/ UNet	18.52	17.85	19.05	0.57	0.57	0.57	2.98	0.77	1.96
4	Vox [38] w/ ResNet	20.62	19.64	21.22	0.70	0.69	0.68	1.97	0.47	1.19
5	SynSin [47]	20.91	19.80	21.62	0.71	0.71	0.70	1.68	0.43	0.99
6	ours w/o inpainting	—	—	25.42	—	—	0.80	—	—	0.68
7	ours	<b>24.67</b>	<b>22.90</b>	<b>26.00</b>	<b>0.82</b>	<b>0.77</b>	<b>0.82</b>	<b>1.05</b>	<b>0.35</b>	<b>0.54</b>
									<b>23.51</b>	<b>0.85</b>
										<b>1.32</b>

Table 1: Novel View Synthesis: Performance of our and previous approaches on the Matterport dataset and the Replica dataset. All models are trained on Matterport and evaluated on both datasets. See Sec. 4.1 for details.

rendered mesh  $I_{out}$  to fill the missing regions and output a new image  $I_{paint} = G(I_{out})$  as the final target view.

We build our inpainting network based on the generator in pix2pixHD [46], which translates a 4-channel input (the rendered mesh  $I_{out}$  and its foreground mask  $F_{out}$ ) into a 3-channel output image  $I_{paint}$ . Our inpainting network outputs an entire image – it not only fills the invisible regions but also refines the image details in the visible regions. We apply the same RGB pixel L1 loss  $L_{paint}^{rgb}$  and perceptual loss  $L_{paint}^{pc}$  as in Sec. 3.3 on the inpainting output  $I_{paint}$ .

### 3.5. Training

During training, we train our model using the Adam optimizer [15] with a weighted combination each loss as  $L = \lambda_1 L_{out}^{rgb} + \lambda_2 L_{out}^{pc} + \lambda_3 L_{paint}^{rgb} + \lambda_4 L_{paint}^{pc} + \lambda_5 L_g + \lambda_6 L_m$  with  $\lambda_1 = \lambda_3 = 8$ ,  $\lambda_2 = \lambda_4 = 2$ ,  $\lambda_5 = 0.2$ , and  $\lambda_6 = 10^{-4}$ . Our model is trained for a total of 50000 iterations with batch size 64 and  $10^{-4}$  learning rate.

We use a grid mesh with size  $W_m \times H_m = 33 \times 33$  (and also  $65 \times 65$  in Sec. 4.2). Following SynSin [47], we use  $W_{im} \times H_{im} = 256 \times 256$  as the input and output image size. Our mesh implementation is based on PyTorch3D [31].

## 4. Experiments

We evaluate our model on three datasets: Matterport [3], Replica [41], and RealEstate10K [52], following the prior work’s [47] experimental setup and details. In addition, we analyze the advantages and limitations of wrapping a mesh sheet by testing on in-the-wild high resolution images.

### 4.1. Evaluation on Matterport and Replica

We first train and evaluate our approach on the Matterport [3] dataset, which contains 3D scans of homes. We load the Matterport dataset in the Habitat simulator [33], following the same training, validation, and test splits as in SynSin [47]. During training, we supervise our model with paired 2D images of the input and the target views. We empirically find that it works slightly better to first train the scene mesh predictor (Sec. 3.1) and then freeze the scene

mesh to further train the inpainting network (Sec. 3.4), rather than training both components jointly from scratch.

**Metrics.** Following SynSin [47], we evaluate the quality of our predicted novel view images  $I_{paint}$  using three metrics: Peak Signal-to-Noise Ratio (**PSNR**; higher is better), Structural Similarity (**SSIM**; higher is better), and Perceptual Similarity distance (**Perc Sim**; lower is better). The Perc Sim metric is based on the convolutional feature distance between the prediction and the ground-truth, which is shown highly correlated with human judgement [51, 47]. Since only a part of the target view image can be seen from the input image as illustrated in Figure 3, we separately evaluate these metrics on visible regions (**Vis**, which can be seen from the input view), invisible regions (**InVis**, which cannot be seen but must be imagined), and the entire image (**Both**). Note that the visible region masks are obtained from the ground-truth scene geometry and camera frustum (available from the Habitat simulator) instead of predicted by our mesh, and are the same as in SynSin’s evaluation.

**Baselines.** we compare our method to several previous approaches: **Im2Im** [53] is an image-to-image translation method which predicts a appearance flow to warp an input view to the target view based on an input camera transformation. **Tatarchenko *et al.*** [42] is similar to Im2Im, but directly predicts the target view image instead of an appearance flow. **Vox w/ UNet** and **Vox w/ ResNet** are two variants of the deep voxel representation [38] with different encoder-decoder architectures based on UNet or ResNet, implemented in [47]. **SynSin** [47] projects a dense feature point cloud (extracted from every image pixel) to the target camera pose and applies a refinement network on the point cloud projection to output the target view image. See Sec. 2 for more discussions on SynSin.

We also evaluate the prediction of our model before inpainting (*i.e.* directly using the mesh rendering output  $I_{out}$  as the target view) to analyze how well our method performs with texture sampling and mesh rendering alone.

**Results.** The results are shown in Table 1. Even without inpainting, the mesh rendering output  $I_{out}$  from our method already outperforms previous approaches by a large mar-

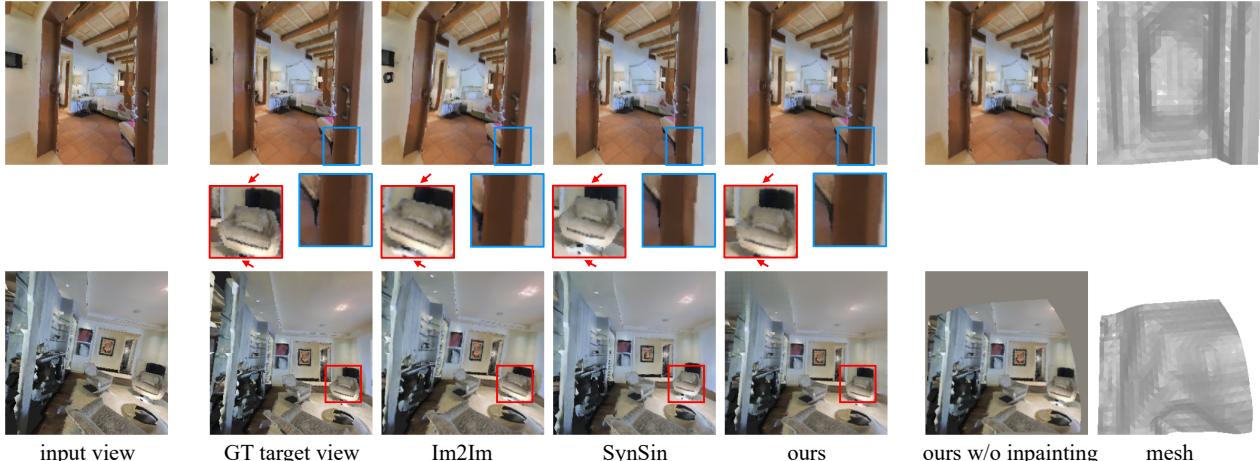


Figure 4: Novel views from our and previous methods on the Matterport dataset (scene mesh shown in the last column). The first row have same view-point change as in [47] while the second row has  $2\times$  larger camera angle change.

#	Method	Matterport [3] ( $2\times$ cam. change)						Replica [41] ( $2\times$ cam. change)					
		PSNR $\uparrow$			SSIM $\uparrow$			Perc Sim $\downarrow$					
		Both	InVis	Vis	Both	InVis	Vis	Both	InVis	Vis			
1	Im2Im [53]	14.93	15.16	15.28	0.51	0.56	0.46	3.26	0.93	1.91	15.91	0.63	2.63
2	Tatarchenko <i>et al.</i> [42]	14.71	14.77	15.08	0.56	0.61	0.52	3.74	1.04	2.14	14.19	0.68	3.37
3	SynSin [47]	19.15	17.76	20.69	0.67	0.66	0.66	2.06	0.78	0.96	19.63	0.77	1.94
4	ours w/o inpainting	—	—	24.20	—	—	0.76	—	—	0.69	—	—	—
5	ours	<b>22.62</b>	<b>20.89</b>	<b>24.76</b>	<b>0.77</b>	<b>0.72</b>	<b>0.77</b>	<b>1.41</b>	<b>0.63</b>	<b>0.56</b>	<b>21.12</b>	<b>0.81</b>	<b>1.70</b>

Table 2: Novel View Synthesis: Generalization performance of our and previous approaches to larger view-point changes on the Matterport dataset and the Replica dataset. All models are trained on the Matterport dataset and evaluated on both datasets with  $2\times$  camera angle changes than in the training data. See Sec. 4.1 for details.

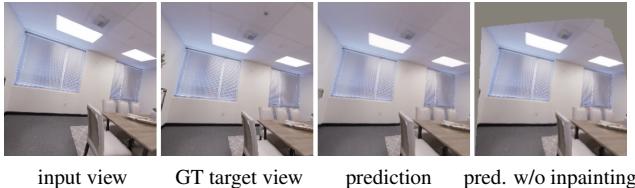


Figure 5: Generalization of our model (trained on Matterport) to the Replica dataset without retraining.

gin under all the three metrics on the visible regions. With the help of an inpainting network, our final output  $I_{paint}$  has significantly higher performance than previous work on both invisible and visible regions, achieving a new state-of-the-art performance on this dataset. Figure 4 shows view synthesis examples from our method and previous work on the Matterport dataset, where our method can paint things such as doorframe or sofa at more precise locations.

**Generalization to the Replica dataset.** Following [47], we also evaluate how well our model generalizes to another scene dataset, Replica [41], which contains high-quality laser scans of both homes and offices. We take our model trained on the Matterport dataset and directly evaluate on the Replica dataset without re-training. The results are

shown in Table 1, where all methods are trained and evaluated under the same setting. It can be seen that our method achieves noticeably better generalization to this dataset and outperforms previous approaches by a large margin. Figure 5 shows view synthesis examples on the Replica dataset.

**Generalization to larger view-point changes.** We further analyze how well our and previous approaches generalize to larger camera pose changes beyond their training data. In this analysis, we sample new input-target view pairs on the test scenes with  $2\times$  larger camera angle changes than in the training data, and directly evaluate all approaches on these new view-points without retraining. The results are shown in Table 2, where our method largely outperforms other approaches under all metrics. Figure 4 (second row) shows an example under  $2\times$  larger camera angle change.

## 4.2. Evaluation on RealEstate10K

The RealEstate10K dataset [52] is built upon video frames extracted from YouTube videos on houses. The input view and the target view are different frames sampled from the video within a time range, and the camera poses are estimated using Structure-from-Motion over the video.

On this dataset, we follow the experimental setup in

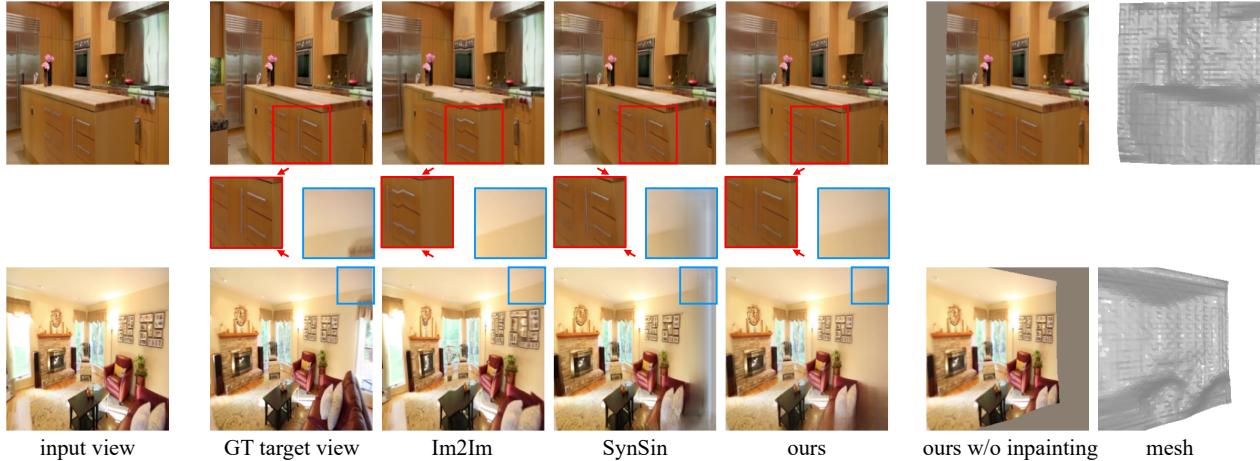


Figure 6: Novel views from our and previous methods on the RealEstate10K dataset (scene mesh shown in the last column).

RealEstate10K [52]				
#	Method	PSNR $\uparrow$	SSIM $\uparrow$	Perc Sim $\downarrow$
1	Im2Im [53]	17.05	0.56	2.19
2	Tatarchenko <i>et al.</i> [42]	11.35	0.33	3.95
3	Vox [38] w/ UNet	17.31	0.53	2.30
4	Vox [38] w/ ResNet	21.88	0.71	1.30
5	3DView (similar to [26])	21.88	0.66	1.52
6	SynSin [47]	22.31	0.74	1.18
7	SynSin $\dagger$ [47]	22.83	0.75	1.13
8	StereoMag $\ddagger$ [52]	25.34	<b>0.82</b>	1.19
9	ours ( $33 \times 33$ mesh)	26.24	<b>0.82</b>	0.83
10	ours ( $65 \times 65$ mesh)	<b>26.74</b>	<b>0.82</b>	<b>0.80</b>

Table 3: Results on the RealEstate10K dataset. We follow the evaluation protocol in SynSin and take the best metrics from predictions over two input views for the same target view. See Sec. 4.2 for details. ( $\dagger$ : SynSin trained 5 $\times$  longer for 250k iterations;  $\ddagger$ : unlike all other approaches, StereoMag takes two views as input to predict the target view.)

SynSin [47] and use the same training, validation, and test data. In addition to using a  $33 \times 33$  mesh, we also train our model with a higher resolution  $W_m \times H_m = 65 \times 65$  mesh, which is initialized from a trained  $33 \times 33$ -mesh model with a new transposed convolution layer to upsample on the feature map  $\{q_{w,h}\}$  in Sec. 3.1 to  $65 \times 65$  spatial dimensions.

We compare our method to several previous approaches. In addition to the baselines in Sec. 4.1, we also compared to two additional approaches. **3DView** is a system similar to the Facebook 3D Photo [26] based on layered depth images, implemented in [47]. **StereoMag** [52] builds multi-plane images (MPIs) from a stereo pair, using images from *two* different views as input at test time. Except for StereoMag, all other methods use a single view at test time.

**Results.** We follow the evaluation protocol of [47] on RealEstate10K. However, to compare with StereoMag [52] that uses two input views on this dataset, the best metrics of two views were reported for single-view methods in [47].

That is, make two separate predictions based on two different input views respectively for each target view at test time and the select the best metrics among the two predictions as the scores for this target view. We believe this metric does not fully capture the true performance of single-image view synthesis as it evaluates on two views at test time, but we report it for apples-to-apples comparison to baselines. Note that apart from Table 3, all other results (Table 1, 2, 4) are average single-view results. The results in Table 3 show that our method achieves the highest performance under the evaluation protocol described above, outperforming previous approaches by a noticeable margin. Besides, a higher resolution  $65 \times 65$  mesh gives a further performance boost.

**Detailed analyses.** As discussed above, we further evaluate by our usual protocol of taking the average metrics over all predictions (which is used on Matterport and Replica in Sec. 4.1), to measure how well the model does on average from a single input view. Apart from metrics over the entire image, we would also like to analyze how well each model does on rendering regions seen in the input view vs. invisible regions (where things must be imagined). However, we cannot compute the exact visibility map as there are no ground-truth geometry annotations on the RealEstate10K dataset. To get an approximation, we evaluate on the central  $\frac{W_{im}}{2} \times \frac{H_{im}}{2} = 128 \times 128$  crop of the target image (**Center**, which is nearly always visible) and the rest of the image (**Peripheral**, containing most of the invisible regions).

These detailed comparisons are shown in Table 4. It can be seen that our method achieves the highest performance on both the center regions (which are mostly visible) and the peripheral regions, outperforms the previous single-view based approaches by a large margin. Figure 6 shows predicted novel views on this dataset. In addition, we visualize the pixel-wise squared error map on the prediction from our method and SynSin [47] in Figure 7, where our method paints objects at more precise locations compared to SynSin, resulting in higher PSNR and better quality.

RealEstate10K [52] (averaged metrics over all predictions)										
#	Method	PSNR ↑			SSIM ↑			Perc Sim ↓		
		Both	Peripheral	Center	Both	Peripheral	Center	Both	Peripheral	Center
1	Im2Im [53]	15.56	15.65	15.80	0.51	0.53	0.44	2.59	1.93	0.65
2	Tatarchenko <i>et al.</i> [42]	11.11	11.32	10.68	0.32	0.35	0.24	3.96	2.96	1.13
3	SynSin [47]	20.11	20.02	21.15	0.67	0.67	0.67	1.52	1.18	0.31
4	SynSin <sup>†</sup> [47]	20.47	20.35	21.65	0.68	0.68	0.68	1.49	1.16	0.29
5	ours (33 × 33 mesh)	23.11	22.90	24.83	<b>0.75</b>	0.74	<b>0.76</b>	1.17	0.94	<b>0.21</b>
6	ours (65 × 65 mesh)	<b>23.41</b>	<b>23.20</b>	<b>25.17</b>	<b>0.75</b>	<b>0.75</b>	<b>0.76</b>	<b>1.14</b>	<b>0.92</b>	<b>0.21</b>

Table 4: Performance of our and previous approaches on RealEstate10K, with metrics averaged over all predictions. We evaluate on the entire  $256 \times 256$  image (**both**), the **center**  $128 \times 128$  crop (nearly always visible), and the remaining **peripheral** regions (containing most of the invisible regions). See Sec. 4.2 for details. (<sup>†</sup>: SynSin trained 5× longer for 250k iterations.)

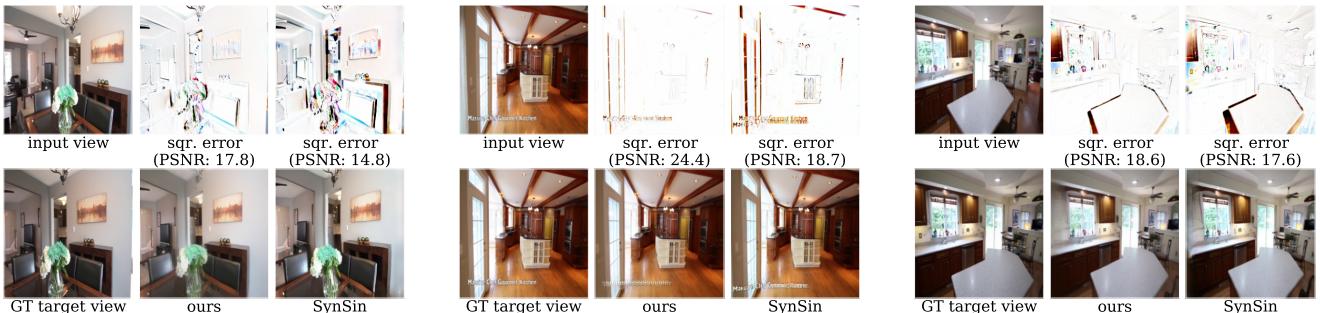


Figure 7: Squared error maps in target view from our method and SynSin [47] on the RealEstate10K dataset (darker is higher error). Our method paints things in the novel view at more precise locations, resulting in lower error and higher PSNR.

### 4.3. Analysis: testing the limits of wrapping sheets

So far, we showed that the idea of wrapping a mesh sheet onto an image achieves strong performance across all benchmarks. But one might wonder how good is a planar sheet prior for novel view synthesis in any arbitrary images. To test the limits of wrapping a mesh *sheet*, we test it over a large variety of images including outdoor scenes, outdoor objects, indoor scenes, indoor objects, and even artistic paintings. We analyze our underlying mesh data structure by pretraining depth [19] to fill in the  $z$  values, and examine how well it generates novel views in-the-wild. As shown in Figure 1 (top row), although missing a few details (such as tree branches), a scene mesh sheet captures the geometric structures sufficient enough to render high-resolution ( $960 \times 960$ ) photorealistic novel views even from very large viewpoint changes. Please see videos at [worldsheet.github.io](https://worldsheet.github.io) for animation of continuously generated views. This result confirms that our mesh sheet data structure and warping procedure, despite being simple, are flexible enough to handle the variety of the visual world.

**Limitations and Future Work.** Although shrink-wrapping a mesh sheet onto images works well on a wide range of scenes, one limitation is that we currently do not cut the mesh on depth discontinuities, assuming that the foreground objects (e.g. trees) are connected to the background (e.g. sky), which can create artifacts near object boundaries. Figure 8 shows failure cases due to depth discontinuity.

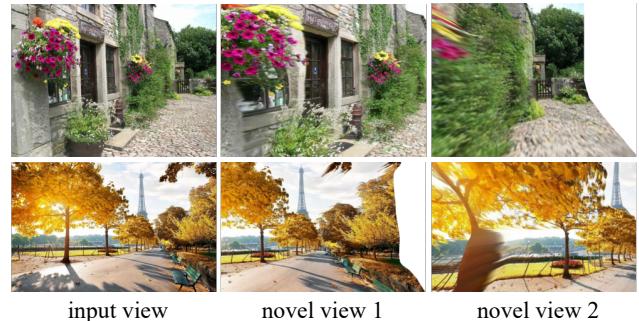


Figure 8: Some failures of our approach due to depth discontinuity, where the foreground objects (trees or flower) are blurry at their boundaries. See Sec. 4.3 for a discussion.

## 5. Conclusion

In this work, we propose *Worldsheet*, which synthesizes novel views from a single image by shrink-wrapping the scene in a grid mesh. Our approach jointly learns the scene geometry and generates novel views through differentiable texture sampling and mesh rendering, supervised with only 2D images of the input and the target views. The approach is category-agnostic and end-to-end trainable resulting in state-of-the-art performance on single-image view synthesis across three datasets by a large margin. Furthermore, we show that our approach to wrap a sheet is flexible enough to handle the variety of the visual world.

## Acknowledgements

We are grateful to Alyosha Efros, Alex Berg, Angjoo Kanazawa, Shubham Goel, Devi Parikh, Ross Girshick, Georgia Gkioxari, Justin Johnson, Nikhila Ravi, Brian Okorn and other colleagues at FAIR and CMU for fruitful discussions. This work was supported in part by DARPA Machine Common Sense grant.

## References

- [1] Kara-Ali Aliev, Artem Sevastopolsky, Maria Kolos, Dmitry Ulyanov, and Victor Lempitsky. Neural point-based graphics. *arXiv preprint arXiv:1906.08240*, 2019. [2](#)
- [2] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, 1999. [2](#)
- [3] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niebner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. In *2017 International Conference on 3D Vision (3DV)*. IEEE, 2017, License for Matterport dataset available at [http://kaldir.vc.in.tum.de/matterport/MP\\_TOS.pdf](http://kaldir.vc.in.tum.de/matterport/MP_TOS.pdf). [2, 5, 6, 12](#)
- [4] Shenchang Eric Chen. Quicktime vr: An image-based approach to virtual environment navigation. In *Proceedings of the 22nd annual conference on Computer graphics and interactive techniques*, 1995. [2](#)
- [5] Xu Chen, Jie Song, and Otmar Hilliges. Monocular neural image based rendering with continuous view control. In *ICCV*, 2019. [2](#)
- [6] Antonio Criminisi, Ian Reid, and Andrew Zisserman. Single view metrology. *International Journal of Computer Vision*, 2000. [2](#)
- [7] Paul E Debevec, Camillo J Taylor, and Jitendra Malik. Modeling and rendering architecture from photographs: A hybrid geometry-and image-based approach. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, 1996. [2](#)
- [8] Shubham Goel, Angjoo Kanazawa, and Jitendra Malik. Shape and viewpoint without keypoints. In *ECCV*, 2020. [3, 4](#)
- [9] Steven J Gortler, Radek Grzeszczuk, Richard Szeliski, and Michael F Cohen. The lumigraph. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, 1996. [2](#)
- [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. [3](#)
- [11] Derek Hoiem, Alexei A Efros, and Martial Hebert. Automatic photo pop-up. In *ACM SIGGRAPH 2005 Papers*. 2005. [1, 2](#)
- [12] Youichi Horry, Ken-Ichi Anjyo, and Kiyoshi Arai. Tour into the picture: using a spidery mesh interface to make animation from a single image. In *Proceedings of the 24th annual conference on Computer graphics and interactive techniques*, 1997. [2](#)
- [13] Angjoo Kanazawa, Shubham Tulsiani, Alexei A Efros, and Jitendra Malik. Learning category-specific mesh reconstruction from image collections. In *ECCV*, pages 371–386, 2018. [2, 3, 4](#)
- [14] Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. Neural 3d mesh renderer. In *CVPR*, 2018. [2, 3](#)
- [15] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. [5](#)
- [16] Johannes Kopf, Kevin Matzen, Suhib Alsian, Ocean Quigley, Francis Ge, Yangming Chong, Josh Patterson, Jan-Michael Frahm, Shu Wu, Matthew Yu, et al. One shot 3d photography. *ACM Transactions on Graphics (TOG)*, 2020. [2](#)
- [17] Nilesh Kulkarni, Abhinav Gupta, and Shubham Tulsiani. Canonical surface mapping via geometric cycle consistency. In *ICCV*, 2019. [2, 3](#)
- [18] Tejas D Kulkarni, William F Whitney, Pushmeet Kohli, and Josh Tenenbaum. Deep convolutional inverse graphics network. In *NIPS*, 2015. [2](#)
- [19] Katrin Lasinger, René Ranftl, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *arXiv preprint arXiv:1907.01341*, 2019. [8](#)
- [20] Marc Levoy and Pat Hanrahan. Light field rendering. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, 1996. [2](#)
- [21] David Liebowitz, Antonio Criminisi, and Andrew Zisserman. Creating architectural models from images. In *Computer Graphics Forum*. Wiley Online Library, 1999. [2](#)
- [22] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. [1](#)
- [23] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *NeurIPS*, 2020. [2](#)
- [24] Shichen Liu, Tianye Li, Weikai Chen, and Hao Li. Soft rasterizer: A differentiable renderer for image-based 3d reasoning. In *ICCV*, 2019. [2, 3, 4](#)
- [25] Madara Lukjanovica. Italy alleyway sketch. <https://www.pinterest.com/pin/546342998522093232>. [1](#)
- [26] Kevin Matzen, Matthew Yu, Jonathan Lehman, Peizhao Zhang, Jan-Michael Frahm, Peter Vajda, Johannes Kopf, and Matt Uyttendaele. Powered by ai: Turning any 2d photo into 3d using convolutional neural nets. <https://ai.facebook.com/blog/powerd-by-ai-turning-any-2d-photo-into-3d-using-convolutional-neural-nets/>. Feb 2020. [7](#)
- [27] 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. *ECCV*, 2020. [2](#)
- [28] Simon Niklaus and Feng Liu. Softmax splatting for video frame interpolation. In *CVPR*, 2020. [4, 11](#)

- [29] Simon Niklaus, Long Mai, Jimei Yang, and Feng Liu. 3d ken burns effect from a single image. *ACM Transactions on Graphics (TOG)*, 2019. 2
- [30] Georgios Pavlakos, Nikos Kolotouros, and Kostas Daniilidis. Texturepose: Supervising human mesh estimation with texture consistency. In *ICCV*, 2019. 3
- [31] Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. Accelerating 3d deep learning with pytorch3d. *arXiv preprint arXiv:2007.08501*, 2020. 2, 3, 4, 5
- [32] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *ECCV*, 2020. 2
- [33] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A platform for embodied ai research. In *ICCV*, 2019. 5
- [34] Ashutosh Saxena, Min Sun, and Andrew Y Ng. Make3d: Learning 3d scene structure from a single still image. *IEEE transactions on pattern analysis and machine intelligence*, 2008. 2
- [35] Steven M Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. A comparison and evaluation of multi-view stereo reconstruction algorithms. In *CVPR*, 2006. 2
- [36] Meng-Li Shih, Shih-Yang Su, Johannes Kopf, and Jia-Bin Huang. 3d photography using context-aware layered depth inpainting. In *CVPR*, 2020. 2
- [37] Dave Shreiner, Bill The Khronos OpenGL ARB Working Group, et al. *OpenGL programming guide: the official guide to learning OpenGL, versions 3.0 and 3.1*. Pearson Education, 2009. 3, 4
- [38] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhofer. Deepvoxels: Learning persistent 3d feature embeddings. In *CVPR*, 2019. 5, 7, 12
- [39] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *NeurIPS*, 2019. 2
- [40] 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 *CVPR*, 2019. 2
- [41] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, et al. The replica dataset: A digital replica of indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019. 2, 5, 6, 12
- [42] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. Multi-view 3d models from single images with a convolutional network. In *ECCV*. Springer, 2016. 2, 5, 6, 7, 8, 12
- [43] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *CVPR*, 2020. 2
- [44] University of Oxford. Oxford New College. <https://www.biology.ox.ac.uk/applications>. 1
- [45] Michael Waechter, Nils Moehrle, and Michael Goesele. Let there be color! large-scale texturing of 3d reconstructions. In *ECCV*. Springer, 2014. 2
- [46] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *CVPR*, 2018. 4, 5
- [47] Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. Synsin: End-to-end view synthesis from a single image. In *CVPR*, pages 7467–7477, 2020. 2, 3, 4, 5, 6, 7, 8, 11, 12, 14
- [48] Daniel E Worrall, Stephan J Garbin, Daniyar Turmukhametov, and Gabriel J Brostow. Interpretable transformations with encoder-decoder networks. In *ICCV*, 2017. 2
- [49] Jae Shin Yoon, Kihwan Kim, Orazio Gallo, Hyun Soo Park, and Jan Kautz. Novel view synthesis of dynamic scenes with globally coherent depths from a monocular camera. In *CVPR*, 2020. 2
- [50] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *ICLR*, 2016. 3
- [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 Transactions on Graphics (TOG)*, 2018. 2, 5, 6, 7, 8
- [53] Tinghui Zhou, Shubham Tulsiani, Weilun Sun, Jitendra Malik, and Alexei A Efros. View synthesis by appearance flow. In *ECCV*. Springer, 2016. 5, 6, 7, 8, 12
- [54] C Lawrence Zitnick, Sing Bing Kang, Matthew Uyttendaele, Simon Winder, and Richard Szeliski. High-quality video view interpolation using a layered representation. *ACM transactions on graphics (TOG)*, 2004. 2

# Worldsheet: Wrapping the World in a 3D Sheet for View Synthesis from a Single Image

(Supplementary Material)

## A. Continuous and large viewpoint changes

Our approach allows synthesizing continuous novel views by smoothly moving to a new camera pose largely different from the input. We kindly request the readers to check out the videos at <https://worldsheet.github.io> for a better understanding of the performance of our method in an absolute sense and in comparison to other approaches.

In these videos, we compare our synthesized novel views (from a single image) to SynSin [47] on the RealEstate10K dataset with simulated large view-point changes (the first frame contains the input view and the rest of the frames are synthesized). From the videos, it can be seen that our model can generate novel views with much larger camera translation and rotation than in the training data, while SynSin often suffers from severe artifacts in these cases, likely because its refinement network does not generalize well to a sparser point cloud (resulting from large camera zoom-in or rotation).

We also show in these videos continuously synthesized novel views on high resolution ( $960 \times 960$ ) images over a wide range of scenes (the first frame is the input view), as described in our analysis in Sec. 4.3 in the main paper.

## B. Details on differentiable texture sampler

Our differentiable texture sampler (Sec. 3.2 in the main paper) splat image pixels onto the texture space through each face. This splatting procedure involves three main steps: forward-mapping, normalization, and hole filling.

Suppose the flow  $(u, v) = f(i, j)$  maps image coordinates  $(i, j)$  to UV coordinates  $(u, v)$  on the texture map, which can be obtained through mesh rasterization (here we drop the face index  $k$  for simplicity). To implement  $\hat{T} = \text{splat}(I_{in}, f)$  (Eqn. 5 in the main paper), we first forward-map the image pixels to the texture map, using bilinear assignment when the mapped texture coordinates  $(u, v)$  do not fall on integers, as forward\_map below.

```
function T = forward_map(I, f)
    T = all_zeros(W_uv, H_uv)
    for i in 1:W_im
        for j in 1:H_im
            u, v = f(i, j)
            uf, uc = floor(u), ceil(u)
            vf, vc = floor(v), ceil(v)
            T[uf,vf] += I[i,j] * (uc - u) * (vc - v)
```

```
T[uc,vf] += I[i,j] * (u - uf) * (vc - v)
T[uf,vc] += I[i,j] * (uc - u) * (v - vf)
T[uc,vc] += I[i,j] * (u - uf) * (v - vf)
end
return T
```

Gradients can be taken over  $f$  through the bilinear weights. Our implementation is based on the sum splatting in [28].

However, forward mapping alone will lead to incorrect pixel intensity (*e.g.* imagine down-scaling an image to half its width and height by forward-mapping – each pixel in the low-resolution image will receive assignment from 4 pixels and become  $4\times$  brighter). Hence, a second normalization step is applied:

$$\hat{T}_{sum} = \text{forward\_map}(I_{in}, f) \quad (\text{B.1})$$

$$\hat{W}_{sum} = \text{forward\_map}(I_{one}, f) \quad (\text{B.2})$$

$$\hat{T}_{norm} = \hat{T}_{sum} / \max(\hat{W}_{sum}, 10^{-4}) \quad (\text{B.3})$$

where  $I_{one}$  is an  $W_{im} \times H_{im}$  image with all ones as its pixel intensity.  $\hat{T}_{norm}$  contains the normalized splatting result. A threshold  $10^{-4}$  is applied to avoid division by zero (which could happen due to holes described below).

**Filling holes with Gaussian filtering.** It is well known that the bilinear forward mapping above often leads to holes in the output (*e.g.* imagine up-scaling an image to a much larger size – there will be gaps in the output image as some pixels will not receive assignments). To minimize hole occurrence, in our UV texture map we assigned the UV coordinates of each mesh vertex with an equally-spaced  $W_m \times H_m$  lattice grid, and use the same image size as the texture map size ( $W_{uv} \times H_{uv} = W_{im} \times H_{im}$ ). This ensures that most texels on the texture map receive assignments in forward mapping, so that holes rarely occur in  $\hat{T}_{norm}$ . However, to address corner cases, we further apply a Gaussian filter to fill the holes on  $\hat{T}_{norm}$  (where  $\hat{W}_{sum}$  is zero as no assignment is received from forward-mapping):

$$\hat{M} = I[\hat{W}_{sum} > 0] \quad (\text{B.4})$$

$$\hat{T}_g = (F_g * \hat{T}_{norm}) / (F_g * \hat{M}) \quad (\text{B.5})$$

$$\hat{T} = \hat{M} \cdot \hat{T}_{norm} + ((1 - \hat{M}) \cdot \hat{T}_g) \quad (\text{B.6})$$

where  $F_g$  is a discrete 2D Gaussian kernel for image filtering (we use kernel size 7 and standard deviation 2 for

#	Method	Matterport [3]						Replica [41]					
		PSNR $\uparrow$			SSIM $\uparrow$			Perc Sim $\downarrow$					
		Both	InVis	Vis	Both	InVis	Vis	Both	InVis	Vis			
1	Im2Im [53]	15.87	16.20	15.97	0.53	0.60	0.48	2.99	0.58	2.05	17.42	0.66	2.29
2	Tatarchenko <i>et al.</i> [42]	14.79	14.83	15.05	0.57	0.62	0.53	3.73	0.74	2.50	14.36	0.68	3.36
3	Vox [38] w/ UNet	18.52	17.85	19.05	0.57	0.57	0.57	2.98	0.77	1.96	18.69	0.71	2.68
4	Vox [38] w/ ResNet	20.62	19.64	21.22	0.70	0.69	0.68	1.97	0.47	1.19	19.77	0.75	2.24
5	SynSin [47]	20.91	19.80	21.62	0.71	0.71	0.70	1.68	0.43	0.99	21.94	0.81	1.55
6	SynSin (w/ depth sup.) [47]	21.59	20.32	22.46	0.72	0.71	0.71	1.60	0.43	0.92	22.54	0.80	1.55
7	ours	<b>24.67</b>	<b>22.90</b>	<b>26.00</b>	<b>0.82</b>	<b>0.77</b>	<b>0.82</b>	<b>1.05</b>	<b>0.35</b>	<b>0.54</b>	<b>23.51</b>	<b>0.85</b>	<b>1.32</b>
8	ours (w/ depth sup.)	24.75	22.85	26.18	0.82	0.77	0.82	1.06	0.36	0.54	22.78	0.84	1.51

Table C.1: Analyses of our model and SynSin using explicit depth supervision during training (line 8 and 6) on the Matterport dataset. Our model without depth (the default setting; line 7) performs nearly equally well as its counter part with depth supervision (line 8) on Matterport and generalizes better to Replica, outperforming both variants of SynSin (line 5 and 6).

$F_g$  in our implementation). Here  $\hat{M}$  is a binary mask indicating which pixels have received assignments in forward mapping (*i.e.* 1 means valid and 0 means holes),  $\hat{T}_g$  is the Gaussian-blurred version of  $\hat{T}_{norm}$  (where the division ensures the correct pixel intensity; otherwise it will be darker due to holes in  $\hat{T}_{norm}$ ) and is used to fill only the holes in  $\hat{T}_{norm}$ . We use  $\hat{T}$  in Eqn. B.6 as the final splatting output.

### C. Ablation study: using depth supervision

In our experiments in the paper, we show that our model can be trained using only two views of a scene *without* 3D or depth supervision. In this section, we further analyze our approach by training it *with* depth supervision on the Matterport dataset, where the ground-truth depth can be obtained from the Habitat simulator.

In this analysis, we modify the differentiable mesh renderer to render RGB-D images from our mesh, and apply an L1 loss between the ground-truth and the rendered depth as additional supervision. We also compare to the performance of SynSin with depth supervision (reported in [47]). The results are shown in Table C.1. It can be seen that our model without depth supervision (the default setting; line 7) works almost equally as well as its counterpart using depth supervision (line 8) on the Matterport dataset and generalizes better to the Replica dataset. In addition, it outperforms SynSin under both supervision settings (line 5 and 6).

### D. Hyper-parameters in our model

In our nonlinear function  $g(\cdot)$  to scale the network prediction into depth values (in Eqn. 3 in the main paper), we use different output scales based on the depth range in each dataset. On Matterport and Replica, we use

$$g(\psi) = 1/(0.75 \cdot \sigma(\psi) + 0.01) - 1. \quad (\text{D.1})$$

On RealEstate10K (which has larger depth range), we double the output depth scale and use

$$g(\psi) = 2/(0.75 \cdot \sigma(\psi) + 0.01) - 2. \quad (\text{D.2})$$

However, we find that the performance of our model is quite insensitive to the hyper-parameters in  $g(\cdot)$ .

In our differentiable texture sampler and the mesh renderer, we mostly follow the hyper-parameters in PyTorch3D. We use  $K = 10$  faces per pixel and 1e-8 blur radius in mesh rasterization, 1e-4 sigma and 1e-4 gamma in softmax RGB blending, and background color filled with the mean RGB intensity on each dataset. On Matterport and Replica, the input views have 90-degree field-of-view. On RealEstate10K, we multiply the actual camera intrinsic matrix of each frame into its camera extrinsic  $R$  and  $T$  matrices, so that we can still use the same intrinsics and 90-degree field-of-view in the renderer. On high resolution images in the wild (Sec. 4.3 in the main paper), we assume 45-degree field-of-view.

We choose our mesh size  $W_m$  and  $H_m$  based on the image size. In our experiments on Matterport and Replica (Sec. 4.1 in the main paper), we use  $256 \times 256$  input image resolution following SynSin [47], and use pixel stride 8 on the lattice grid sheet (from which our mesh is built), resulting in  $W_m = H_m = 1 + 256/8 = 33$ . On the RealEstate10K dataset (Sec. 4.2 in the main paper), we additionally experiment with pixel stride 4 on the grid sheet, giving  $W_m = H_m = 1 + 256/4 = 65$ . In our analysis on high resolution images in the wild (resized to have the image long side equal to 960 and padded to  $960 \times 960$  square size for ease of rendering in PyTorch3D; Sec. 4.3 in the main paper), we use  $W_m \times H_m = 129 \times 129$  mesh on the actual image regions (not including the padding regions).

### E. More visualized examples

Figure E.1 shows additional visualization and error map comparisons between our approach and SynSin on the RealEstate10K dataset, where our method paints things in the novel view at more precise locations with lower error and higher PSNR.

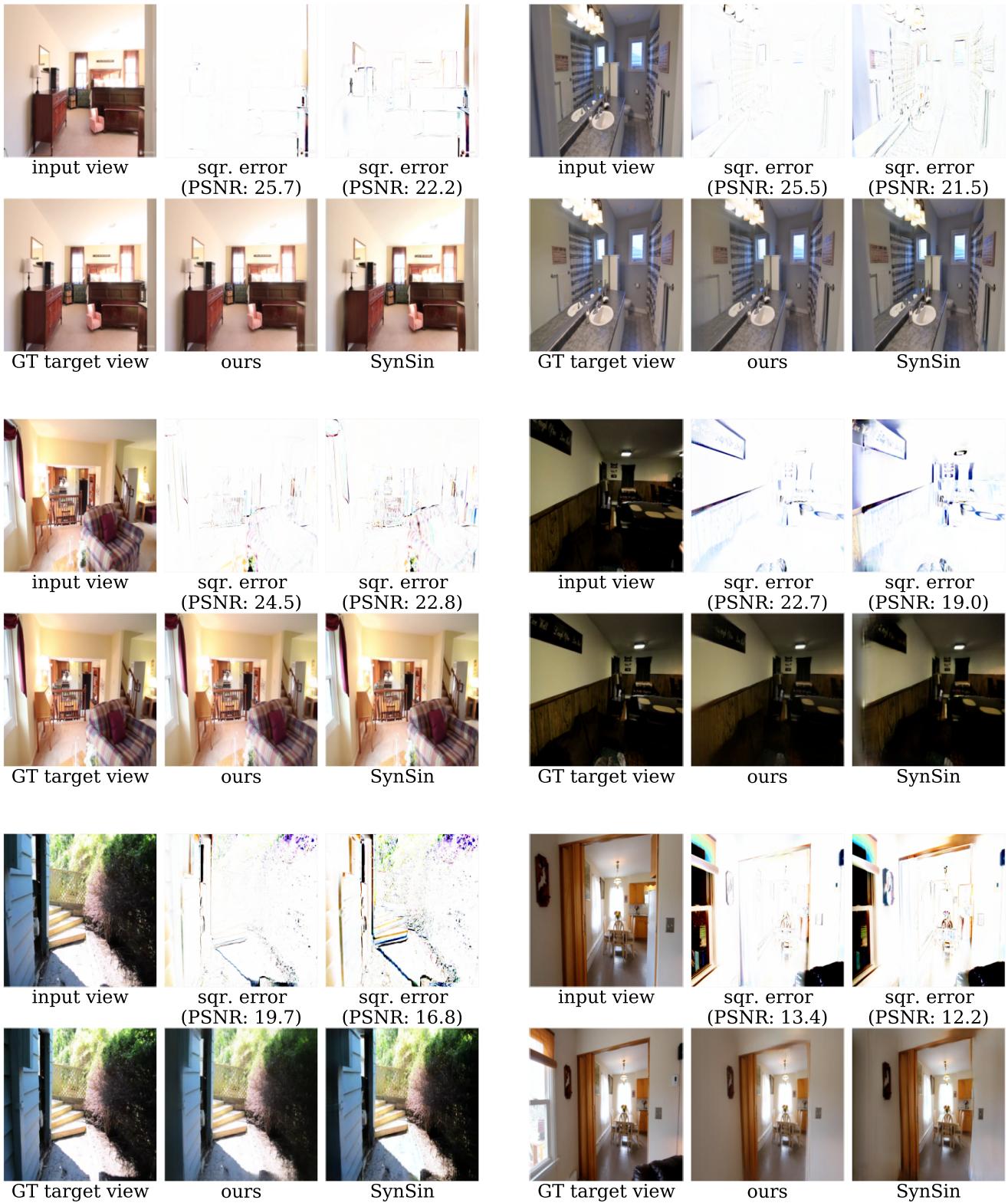




Figure E.1: Additional synthesized novel views (with squared error maps of the target view) from our method and SynSin [47] on the RealEstate10K dataset (darker is higher error). Our method paints things in the novel view at more precise locations, resulting in lower error and higher PSNR.