

---

# Zero-Shot Scene Reconstruction from Single Images with Deep Prior Assembly

---

**Anonymous Author(s)**

Affiliation

Address

email

## Abstract

Large language and visual models have been leading a revolution in visual computing. By greatly scaling up sizes of data and model parameters, the large models learn deep priors which lead to remarkable performance in various tasks. In this work, we present deep prior assembly, a novel framework that assembles diverse deep priors from large models for scene reconstruction from single images in a zero-shot manner. We show that this challenging task can be done without extra knowledge but just simply generalizing one deep prior in one sub-task. To this end, we introduce novel methods related to poses, scales, and occlusion parsing which are keys to enable deep priors to work together in a robust way. Deep prior assembly does not acquire any 3D or 2D data-driven training in the task and demonstrates superior performance in generalization priors on open-world scenes. We conduct evaluations on large-scale dataset, and report analysis, numerical and visual comparisons with the latest methods to show our superiority.

## 1 Introduction

Reconstructing scenes from images is a vital task in 3D computer vision and computer graphics. It bridges the gap between the 2D images that can be easily captured by phone cameras and the 3D geometries of scenes for various real-world applications, e.g., autonomous driving, augmented/virtual reality and robotics. Reconstructing scenes from multi-view images [53, 56] is well-explored to recover 3D geometries with multi-view consistency and camera poses. However, reconstructing a scene from a single RGB image is still challenging, which is extremely difficult due to the quite limited information. Recent works [11, 35] try to solve this task as a reconstruction problem which leverages neural networks with encoder-decoders to draw supervisions from pairs of images and 3D ground truth geometries and layouts. Nevertheless, they are limited in the amount of image-scene pairs and also struggle to generalize to out-of-distribution images in open world.

Large language and visual models have been being extensively studies in the past few years, which revolutionizes neural language processing [51, 7], 2D/3D representation learning [16, 58] and content generation [43, 25], etc. By greatly scaling up sizes of training samples and model parameters, large models show brilliant capabilities and remarkable performance. However, they are limited in a specific task with a specific modality, which limits their capability in high level perception tasks. Driven by the observation, we raise an interesting question: can we assemble series of deep priors from large models, which are experts with different modalities in different tasks, to solve an extremely challenging task that none of them can accomplish alone?

In this work, we propose *deep prior assembly*, a novel framework which assembles diverse deep priors from large models for scene reconstruction from single images in a zero-shot manner. We rethink this task from a new perspective, and decompose it into a set of sub-tasks instead of seeking to a data-driven solution. We narrow down the responsibility of each deep prior on a sub-task that it

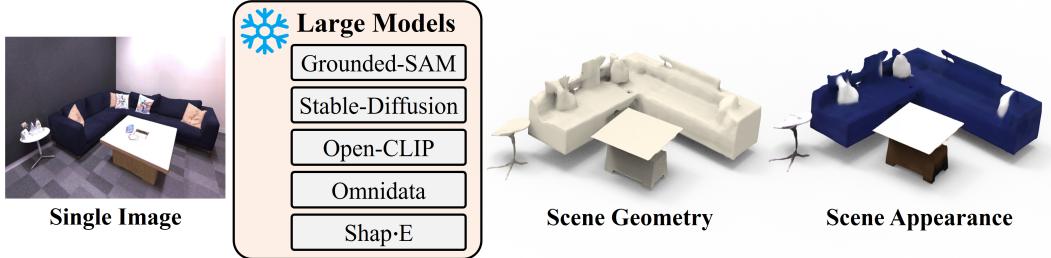


Figure 1: **An illustration of our work.** We assemble diverse deep priors from large models with frozen parameters for scene reconstruction from single images in a zero-shot manner.

is good at, and introduce novel methods related to poses, scales, and occlusion parsing to enable deep priors to work together in a robust way.

Specifically, we first detect and segment the instances in the input image with Grounded-SAM [27, 30], which is a variation of Segment-Anything Model [27]. For the segmented instances that are often corrupted due to occlusions or of low resolution, we leverage Stable-Diffusion [43] to enhance and inpaint images containing the segmented instances. However, the Stable-Diffusion often produces some predictions which drift away from the input instances and do not conform to the original appearance. To solve this issue, we introduce to use Open-CLIP [39, 24] to filter out the bad samples and select the ones matching the input instance most. We then generate the 3D models for each instance with Shap-E [25] using the amended instance images as input. Additionally, we estimate the depth of the origin image with Omnidata [14] as the 3D scene geometry prior. To recover a layout consistent to the image, we propose an approach to optimize the location, orientation and size for each 3D instance to fit it to the estimated segmentation masks and the depths.

Deep prior assembly merely generalizes deep priors and does not require additional data-driven training for extra prior knowledge. Our evaluations on various open-world scenes show our capability of reconstructing diverse objects and recovering plausible layout merely from a single view angle. Our main contributions can be summarized as follows.

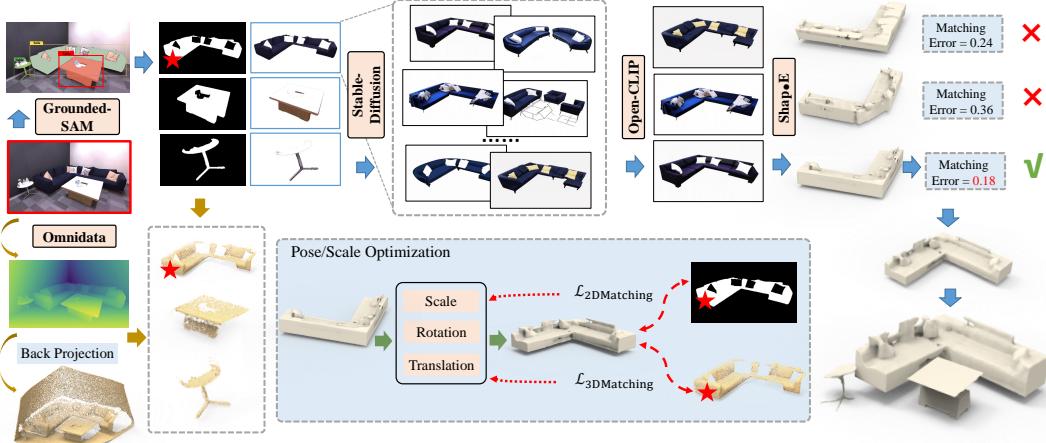
- We propose deep prior assembly, a flexible framework that assembles diverse deep priors from large models together for reconstructing scenes from single images in a zero-shot manner.
- We introduce a novel approach on optimizing the location, orientation and size of instances by matching them with both 2D and 3D supervision.
- We evaluate deep prior assembly for generating diverse open-world 3D scenes, and show our superiority over the state-of-the-art supervised methods.

## 2 Related Work

### 2.1 Large Models in Different Modalities

Recently, it has been drawing significant attention to scaling up deep models for much more powerful representations and higher performances with different modalities (e.g. NLP, 2D vision). Starting from NLP, recent works in scaling up pre-trained language models [7, 31, 40] have largely revolutionized natural language processing. Some research in computer vision [39, 13, 4, 22, 16] translates the progress from language to 2D vision via model and data scaling. CLIP models [39, 20] learn large-scale cross-modality representations with contrastive learning to align text and image concepts. More recently, Uni3D [58] investigates powerful and large-scale 3D representations.

Except for the large foundation models which focus on producing large-scale representations for language, 2D images or 3D point clouds, some research explore large models for specific tasks (e.g. text-to-image generation, image segmentation, 3D object generation) and have shown remarkable performance. Stable Diffusion trains a large model of latent diffusion models and achieves commercially available 2D content generation effects. Segment Anything Model (SAM) [27] and its following works (e.g. Grounded-SAM, Semantic-SAM) revolutionize the field of image segmentation by training models with large-scale annotated data. Omnidata trains the large depth estimation model



**Figure 2: The overview of deep prior assembly.** Given a single image of a 3D scene, we detect the instances and segment them with Grounded-SAM. After normalizing the size and center for the instances, we attempt to amend the quality of the instance images by enhancing and inpainting them. Here, we take a sofa in the image for example. Leveraging the Stable-Diffusion model, we generate a set of candidate images with the image-to-image generation and a text prompt of the instance category predicted by Grounded-SAM. We then filter out the bad generation samples with Open-CLIP by evaluating the cosine similarity between the generated instances and original one. After that, we generate multiple 3D model proposals for this instance with Shap-E from the Top- $K$  generated instance images. Additionally, we estimate the depth of the origin input image with Omnidata as a 3D geometry prior. To estimate the layout, we propose an approach to optimize the location, orientation and size for each 3D proposal by matching it with the estimated segmentation masks and the depths (the  $\star$  for the example sofa). Finally, we choose the 3D model proposal with minimal matching error as the final prediction of this instance, and the final scene is generated by combining the generated 3D models for all detected instances.

77 with various data sources for bringing robust 3D awareness to pure RGB images. In the 3D domain,  
 78 the recent works Point-E and Shap-E collect millions of 3D objects to train large 3D models for  
 79 generating 3D geometries from rendering-style images or texts. In this work, we aim at leveraging  
 80 the powerful capabilities of the large models in different modalities and different domains to solve a  
 81 challenging task, i.e. scene generation from single images, by assembling deep priors together.

## 82 2.2 Scene Reconstruction from Images

83 Recovering the underlying 3D surfaces of scenes from images is a long-standing task in 3D computer  
 84 vision. Most of the previous works focus on the multi-view reconstruction with the input dense images  
 85 captured around the scene. Classic multi-view stereo methods [2, 6, 5] mainly represent the scene  
 86 by estimating depths for dense images by feature matching. Inspired by NeRF [34] which performs  
 87 volume rendering for scene representation, a series of works [53, 36, 55, 52, 56, 28, 1] introduce  
 88 the neural implicit surface reconstruction by learning signed distance fields [37] or occupancy fields  
 89 [33] for scenes from multi-view images. Manhattan-SDF [21] obtains depth maps from COLMAP  
 90 [46] and adopts Manhattan world priors to handle the texture-less regions of the scene. NeuRIS [52]  
 91 proposes to use normal priors for indoor scene reconstruction, and MonoSDF [56] further introduces  
 92 monocular depth cues for improving scene geometries.

## 93 3 Method

94 **Overview.** The overview of deep prior assembly is shown in Fig. 2. We will start from an  
 95 introduction of the task decomposition in Sec. 3.1 and then present the pipeline for solving each of  
 96 the decomposed sub-tasks using a deep prior from a specific large model in Sec. 3.2. Finally, we  
 97 introduce an optimization-based approach for layout estimation in Sec. 3.3.

98    **3.1 Task Decomposition**

99    Revealing 3D scene geometries from a single image is an extremely challenging task due to limited  
100   context and supervisions. Instead of using a data-driven strategy to learn priors [35, 11], we reformulate  
101   this task from a new perspective. We decompose it into a set of sub-tasks, each of which can  
102   be done with one deep prior without a need of learning extra knowledge. More specifically, we can  
103   progressively resolve the task by:

- 104    1) First, performing detection and segmentation on the input image to acquire the segmentation  
105        images, masks and category labels for all detected instances.
- 106    2) Amending instance images through enhancing and inpainting to improve the image qualities.
- 107    3) Generating a set of 3D model proposals for each instance from 2D segmented images.
- 108    4) Estimating the layout by predicting the scale, location, rotation for each 3D proposal to put  
109        them to the correct position of the 3D scene.
- 110    5) Producing a scene reconstruction by applying the estimated layout and shape poses with  
111        reconstructed instances.

112    **3.2 Assembling Large Models**

113   Inspired by the remarkable performances of recent large models, we propose to assign an expert large  
114   model in each sub-task, which maximizes their abilities for modeling a scene in a zero-shot manner.

115   **Detect and Segment 2D instances.** To reveal a scene  $S$  from a single image  $I$ , we first detect the  
116   instances in  $I$  and separate multiple objects into single instances. In this way, we can reconstruct a  
117   scene at shape level, which simplifies the task.

118   **Mask R-CNN vs. SAM vs. Grounded SAM.** Detecting and segmenting images have been widely  
119   explored in the past few years. Mask R-CNN [23] is widely adopted as a popular and robust backbone.  
120   However, the performance of Mask R-CNN does not generalize well since it is only trained under a  
121   relative small dataset. Recently, the large SAM [27] have shown promising segmentation accuracy by  
122   scaling up parameters and using more training samples, nevertheless, it only predicts fine-grained  
123   masks but with few semantic concepts. Thus, we adopt Grounded-SAM, which is an improved  
124   version of SAM by introducing Grounding-DINO [30] as an open-set detector and using SAM to  
125   jointly predict detection boxes, segmentation masks and category labels for each instance, formulated  
126   as:

$$\{m_i, c_i, d_i\}_{i=1}^N = \text{GroundedSAM}(I), \text{ where } d_i > \sigma \quad (1)$$

127    $\{m_i, c_i, d_i\}$  indicate the predicted mask, category and the detection confidence score for the  $i$ -th  
128   instance, and  $N$  is the number of instances in this scene. We only keep the predictions with a high  
129   confidence score larger than a threshold  $\sigma$ . The low confident instances often contain large occlusions  
130   or wrong category predictions.

131   **Enhance and Inpaint 2D instances.** With the predicted masks  $\{m_i\}_{i=1}^N$ , we achieve the segmented  
132   instance images  $\{t_i\}_{i=1}^N$  by masking the original input image  $I$ . We then normalize each instance  
133   image by centering it at the origin and normalizing its scale in  $\{t_i\}_{i=1}^N$  to 0.6 of the max width  
134   or height of  $I$ . As shown in Fig. 2, the segmented instance images often suffer from occlusions  
135   or low-resolution of small instances. The low-quality images have a large negative impact on the  
136   followup 3D generation. Therefore, we propose to first improve the quality of instance images by  
137   enhancing and inpainting them with the large model Stable-Diffusion [43].

138   Specifically, we adopt the image-to-image generation [32] from Stable-Diffusion. We take the  
139   instance image  $t_i$  as the initialization, and add noises on it and then subsequently denoise the noise  
140   corrupted image to increase the realism through the guidance of the text prompt description from the  
141   predicted categories  $c_i$ . We find the prompt template ‘High quality, authentic style {category}’ works  
142   fine for most of the indoor instances. For other situations, we directly leverage the category as the  
143   prompt. We observe that Stable-Diffusion may produce some unreliable predictions which are ‘too  
144   creative’ and can not faithfully improve the image quality but turn it into another image, as shown in  
145   Fig. 3. To solve this issue, we generate multiple enhanced images  $\{e_i^j\}_{j=1}^M$  for each instance image  $t_i$   
146   and filter out the bad generation samples with the approaches described next.

147 **Filter Out Bad Generation Samples.**  
 148 To filter out the bad generation samples  
 149 produced by Stable-Diffusion and se-  
 150 lect the top  $K$  enhanced images for the  
 151 following 2D-to-3D generation, we pro-  
 152 pose to leverage the CLIP models as a  
 153 judge to determine which generated im-  
 154 ages  $\{e_i^j\}_{j=1}^M$  from Stable-Diffusion can  
 155 conform to the original appearance of  
 156 the instance  $t_i$ . Specifically, we adopt  
 157 the large open-sourced CLIP [39] model  
 158 Open-CLIP [24] as the implementation.

159 We use cosine similarities  $\{z_i^j\}_{j=1}^M$   
 160 between the generated instance image  
 161  $\{e_i^j\}_{j=1}^M$  and the original one  $t_i$  as a met-  
 162 ric for the selection, formulated as:

$$z_i^j = \frac{f_\theta(e_i^j) \cdot f_\theta(t_i)}{\|f_\theta(e_i^j)\| \|f_\theta(t_i)\|}, \quad (2)$$

163 where  $f_\theta$  is the frozen image encoder  
 164 from Open-CLIP. We use the Top- $K$   
 165 generated instance images with the largest similarities to the original  $t_i$  as the amended images.  
 166 As shown in Fig. 3, Open-CLIP successfully filters out the bad generations.

167 **Generate 3D models from 2D instances.** The object-level 3D generation from single images [33, 54]  
 168 is a well-explored task in 3D computer vision, however, most previous works show limited generation  
 169 performance on open-world images. Shap-E [25] is a large model for 3D generation by training a  
 170 3D hyper diffusion model with millions of non-public 3D objects. Therefore, we leverage Shap-E to  
 171 provide the deep prior to convert 2D instance images into 3D reconstructions. By this way, we obtain  
 172  $K$  3D reconstruction proposals  $\{s_i^k\}_{k=1}^K$  for each 2D instance  $t_i$  by employing Shap-E on the top  $K$   
 173 amended images.

### 174 3.3 Recovering Scene Layout

175 The final step is to select the most accurate 3D model proposal  $\hat{s}_i$  from the  $K$  candidates  $\{s_i^k\}_{k=1}^K$   
 176 and put it to the right position in a 3D scene to recover the scene layout in the input image. To achieve  
 177 this, we propose a novel approach to optimize the location, orientation and size for each 3D proposal  
 178 in  $\{s_i^k\}_{k=1}^K$  by matching it with the estimated segmentation mask and the estimated depth. We further  
 179 introduce a RANSAC-like solution for robust position optimizing. We select the reconstruction with  
 180 the minimum matching error as the reconstruction  $\hat{s}_i$  of  $t_i$ .

181 **Depth Estimation.** For more accurate  
 182 layout estimations, we first estimate the  
 183 depth map of the input image  $I$  as a 3D  
 184 geometry prior. We leverage the large  
 185 model Omnidata [14, 26] as the depth  
 186 estimator which is trained under a col-  
 187 lected huge depth dataset [14] contain-  
 188 ing 14-million indoor, outdoor, scene-  
 189 level and object-level samples.

190 An issue here is that the depth  $D$  esti-  
 191 mated for the input image  $I$  with Om-  
 192 nidata is not scale-aware, which can not  
 193 be directly used as supervisions. We  
 194 solve this problem by estimating the  
 195 scale  $h$  and shift  $q$  of the predicted depth  
 196 using only one pair of predicted and real  
 197 depth of a randomly selected scene for

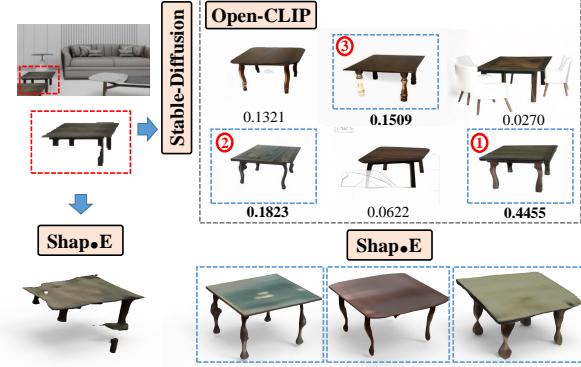


Figure 3: **Examples on the effect of our pipeline.** For the corrupted 2D instant segmented from the scene image, we leverage Stable-Diffusion to produce 6 amended generations. We then adopt Open-CLIP to filter out bad samples by judging the similarities and producing confidence scores for the generations, and keep the Top-3 generated images. The shape generations with Shap-E from the amended images are significantly more complete and accurate than the one produced by the original corrupted image.

As shown in Fig. 3, Open-CLIP successfully filters out the bad generations.

**Generate 3D models from 2D instances.** The object-level 3D generation from single images [33, 54] is a well-explored task in 3D computer vision, however, most previous works show limited generation performance on open-world images. Shap-E [25] is a large model for 3D generation by training a 3D hyper diffusion model with millions of non-public 3D objects. Therefore, we leverage Shap-E to provide the deep prior to convert 2D instance images into 3D reconstructions. By this way, we obtain  $K$  3D reconstruction proposals  $\{s_i^k\}_{k=1}^K$  for each 2D instance  $t_i$  by employing Shap-E on the top  $K$  amended images.

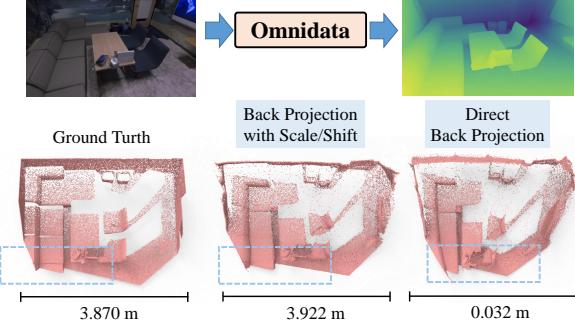


Figure 4: **Illustration of the depth transform.** The es-  
 175 timated depth maps from Omnidata is not scale-aware,  
 176 resulting in scale inaccuracies and distortion in the back-  
 177 projected depth point clouds. We achieve the accurate  
 178 depth point cloud by first transforming the depth maps  
 179 with the pre-solved scale and shift before back-projecting.

198 each dataset. Specifically, we leverage least-squares criterion [15, 41] which has a closed-form  
 199 solution to solve the depth scale and shift by matching the pair of predicted and real depth with a  
 200 specific scene camera intrinsic parameter  $\mathbb{C}_K$ . We justify that the scale and shift from the randomly  
 201 selected scene generalize well to different scenes or even different datasets where the real depth is  
 202 not available. Please refer to the appendix for more analysis and results. After transforming  $D$  with  
 203 the scale  $h$  and shift  $q$ , we can now back-project  $D$  into the 3D space with camera intrinsic parameter  
 204  $\mathbb{C}_K$ , achieving a 3D scene depth point cloud. The example in Fig. 4 shows that the depth point cloud  
 205 produced using the transformed depth maps precisely aligns with the ground truth scene. The depth  
 206 point cloud  $d_i$  for each instance  $t_i$  is further achieved by masking the back-projected 3D point cloud.

207 **Pose/Scale Optimization.** We further estimate the scale and pose of each 3D model proposal  $s_i^k$   
 208 to put them into the right position in the 3D scene. We propose to solve this problem with an  
 209 optimization-based approach on the location, rotation and size for  $s_i^k$  per the mask  $m_i$  from the  
 210 Grounded-SAM and the depth point cloud  $d_i$  from Omnidata.

211 We model this problem as a 7-DoF shape registration task with 3-DoF of translation ( $tx, ty, tz$ ),  
 212 3-DoF of rotation ( $rx, ry, rz$ ) and one DoF of object scale ( $v$ ). Specifically, we first sample a point  
 213 cloud  $p_i^k$  from the mesh of a 3D proposal  $s_i^k$  and initialize a 7-DoF transform as a transformation  
 214 function  $f_\phi$  with learnable 7-DoF parameters  $\phi$ . We then project  $p_i^k$  with  $f_\phi$  to achieve the transformed  
 215 prediction  $\hat{p}_i^k$  by:

$$\hat{p}_i^k = f_\phi(p_i^k). \quad (3)$$

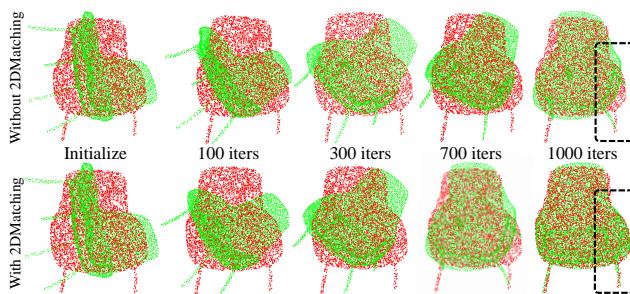
216 We obtain  $\hat{p}_i^k$  by optimizing the 7-Dof parameters  $\phi$  with supervisions. With the estimated depth  $d_i$ ,  
 217 we can draw the direct 3D matching supervision by minimizing the Chamfer Distance Loss between  
 218 the transformed  $\hat{p}_i^k$  and the depth points  $d_i$ . However, merely with the 3D matching constraint, the  
 219 pose/scale optimization do not always converge stably since the predicted depth  $d_i$  is usually with  
 220 noises in complex scenes, which significantly affects the registration on shapes.

221 To resolve this issue, we get inspirations from [9] to leverage the mask information predicted by Stable-  
 222 Diffusion as an extra matching supervision in 2D space. Specifically, we project the transformed  
 223 3D point cloud  $\hat{p}_i^k$  to the 2D space with the camera intrinsic parameters  $\mathbb{C}_K$ , resulting in a 2D point  
 224 cloud  $\tilde{p}_i^k$ . Meanwhile, we form another 2D point cloud  $\tilde{m}_i$  from the mask  $m_i$  by randomly sampling  
 225 dense 2D points on the occupied region of the mask  $m_i$ . We then use the 2D matching constraint  
 226 to minimize the 2D Chamfer Distance between  $\tilde{p}_i^k$  and  $\tilde{m}_i$ . We illustrate the effect of 2D Matching  
 227 constraint with an optimization example in Fig. 5, please also see the visualization of optimization  
 228 process in our video in supplementary materials. The final loss for pose/scale optimization of 3D  
 229 reconstruction  $s_i^k$  is then formulated as:

$$\mathcal{L} = \mathcal{L}_{CD}^{3D}(d_i, \hat{p}_i^k) + \mathcal{L}_{CD}^{2D}(\tilde{p}_i^k, \tilde{m}_i). \quad (4)$$

### 230 **Robust RANSAC-like Solution.**

231 With the optimization-based 7-DoF  
 232 registration, we are now able to put  
 233 the generated 3D object proposals to  
 234 the 3D scene. However, if the mis-  
 235 registration is quite large, especially  
 236 in the rotation, the optimization may  
 237 be trapped in a local optimum and  
 238 fail to produce accurate registrations.  
 239 We further introduce a RANSAC-like  
 240 solution to enhance the robustness of  
 241 pose/scale optimization. Specifically,  
 242 we repeat the optimization  $r$  times  
 243 with randomly initialized rotation  
 244 matrices for  $f_\phi$  each time. The final  
 245 transform for the 3D proposal  $s_i^k$  is  
 246 selected as the one with minimum  
 247 matching loss in Eq. (4) among  $r$   
 248 optimal optimizations, and we define  
 249 the matching error  $w_i^k$  of  $s_i^k$  as the  
 250 minimal matching loss.



251 **Figure 5: Effect of the 2D Matching.** An example of optimiz-  
 252 ing the pose and scale for a chair. We visualize the  
 253 optimization in 2D space. The red 2D points indicate the  
 254 dense 2D point cloud sampled in the mask, which is the target.  
 255 And the green 2D points donate the 2D projection of trans-  
 256 formed 3D point clouds sampled from the generated shape  
 257 of this chair instance. More robust registration is achieved  
 258 with the proposed 2D matching constraint. The total 1,000  
 259 iterations take 9.2 seconds on a single 3090 GPU.

| Method          | 3D-Front     |              |              | BlendSwap    |              |              | Replica      |              |              |
|-----------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
|                 | CDL1-S ↓     | CDL1 ↓       | F-Score↑     | CDL1-S ↓     | CDL1 ↓       | F-Score↑     | CDL1-S ↓     | CDL1 ↓       | F-Score↑     |
| Mesh R-CNN [19] | 0.449        | 0.471        | 23.90        | 0.265        | 0.406        | 21.87        | 0.268        | 0.408        | 25.42        |
| Total3D [35]    | 0.198        | 0.520        | 18.44        | 0.133        | 0.400        | 26.93        | 0.390        | 0.780        | 24.01        |
| PanoRecon [11]  | 0.120        | <b>0.125</b> | 31.94        | 0.355        | 0.417        | 17.11        | 0.326        | 0.440        | 17.13        |
| Ours            | <b>0.109</b> | 0.134        | <b>35.67</b> | <b>0.106</b> | <b>0.089</b> | <b>73.19</b> | <b>0.113</b> | <b>0.110</b> | <b>70.48</b> |

Table 1: Comparisons on scene reconstruction from single images. Lower is better for CDL1 (i.e., Chamfer Distance), higher is better for F-Score. CDL1-S is the single-direction Chamfer Distance from the generated objects to the ground truth meshes.

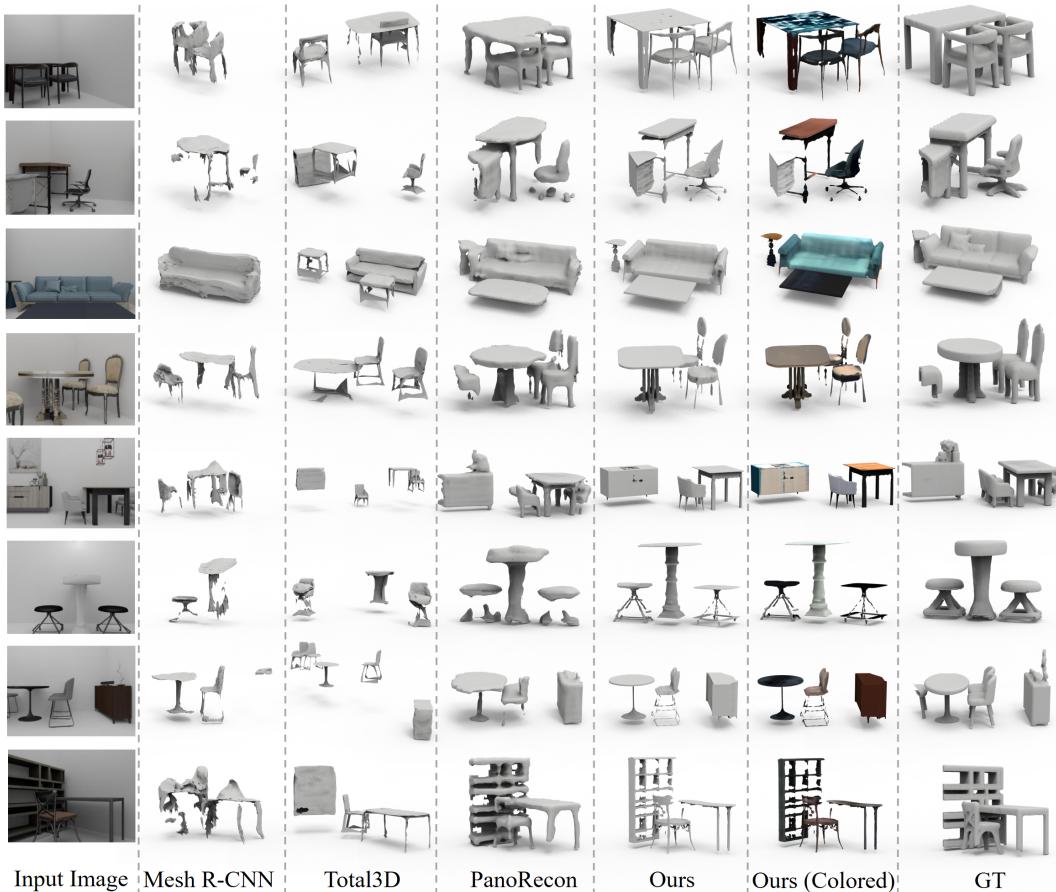


Figure 6: Comparisons on scene generation from single images under the 3D-Front dataset.

251 We repeat the above procedure for each one of the  $K$  3D proposals  $\{s_i^k\}_{k=1}^K$ . We select the 3D  
252 proposal with the minimum  $w_i^k$  as the final reconstruction  $\hat{s}_i$  of  $t_i$ . The final scene generation from  
253 the single image  $I$  is achieved by combining the transformed  $\{\hat{s}_i\}_{i=1}^N$  together.

## 254 4 Experiments

### 255 4.1 Setup

256 **Implement Details.** The number  $M$  of samples generated by Stable-Diffusion for each instance is  
257 set to 6, where we select the Top  $K = 3$  samples with Open-CLIP. The pose/scale optimization is  
258 repeated for  $r = 10$  times for each instance with RANSAC-like solution.

259 **Datasets.** We evaluate deep prior assembly under three widely-used 3D scene reconstruction  
260 benchmarks 3D-Front [18], Replica [48] and BlendSwap [3]. Please refer to the appendix for more  
261 reconstruction results on real-captured scene images from other challenging datasets (e.g. ScanNet  
262 [12]).



Figure 7: Comparisons on scene generation from single images under Replica and BlendSwap dataset.

263 3D-Front [18] is a synthetic 3D dataset of indoor 3D scenes. We adopt the data pre-processed by  
 264 PanoRecon [11] and randomly select 1,000 scene images from the test set as the single-image dataset.  
 265 Note that for the 3D-Front dataset, all the images are captured parallel to the ground with camera  
 266 locations at 0.75m height above the floor. We follow PanoRecon to achieve the corresponding ground  
 267 truth mesh for each image by only keeping the geometry at the same view and cull anything outside  
 268 of the view frustum.

269 The Replica [48] dataset is an indoor scene dataset which contains 8 scanned 3D indoor scene  
 270 with highly photo-realistic 3D indoor scene reconstruction at both room and flat scale. We adopt  
 271 the pre-processed data provided by MonoSDF [56] and sample one image for each scene as the  
 272 single-image dataset. The ground truth meshes are obtained with the same way as 3D-Front.

273 The BlendSwap [3] dataset is an high-fidelity synthesis 3D scene dataset collected by Neural-RGBD  
 274 [3], containing 9 scenes with complex geometries. We collect single-view images and corresponding  
 275 ground truth meshes as the same way as Replica dataset.

276 **Baselines.** We mainly compare our method with the state-of-the-art methods in scene reconstruction  
 277 from single images, i.e., Mesh R-CNN [19], Total3D [35] and PanoRecon [11]. Note that all these  
 278 methods are data-driven methods and trained under 3D datasets with ground truth 3D annotations,  
 279 while our method solves the task in a zero-shot manner. This means that we do not require any  
 280 data-driven training on any 3D or 2D datasets, which is a much more flexible and general solution for  
 281 the single image reconstruction task. We direct evaluate these methods with the official codes and the  
 282 pre-trained models for numerical and visual comparisons.

283 **Metrics.** We use Chamfer Distance and F-Score with the default threshold of 0.1 following [35, 38]  
 284 as metrics. Since Mesh R-CNN and Total3D only predicts the 3D objects and do not generate the  
 285 backgrounds (e.g. wall and floor), we further report the single-direction Chamfer Distance from the  
 286 generated objects to the ground truth meshes, i.e., CDL1-S, to only evaluate the accuracy of generated  
 287 objects. Note that Total3D can generate the scene layout which can roughly represent the background,  
 288 however, we find that Total3D generates layouts with large errors on all the three datasets. Therefore  
 289 we do not keep the layout of Total3D for evaluation. While we achieve the background points by  
 290 back-projecting the segmented background depth maps. We sample 10k points on the ground truth  
 291 meshes and the generated scenes of each methods for evaluation. Please refer to the appendix for  
 292 more details on evaluation.

293 **4.2 Scene Reconstruction on 3D-Front**

294 Tab. 1 reports numerical comparisons on the 3D-Front dataset. We achieve the best performance  
 295 among the state-of-the-art methods. Specifically, PanoRecon is trained under the 3D-Front dataset,  
 296 therefore it shows convincing results in this dataset. Mesh R-CNN and Total3D are trained under  
 297 Pix3D [50]/ShapeNet [8] and SUN-RGBD [47]/Pix3D datasets, respectively.

298 The qualitative comparison is shown in Fig. 6, where we remove the background geometries for  
 299 PanoRecon, ours and GT for a clear visual comparison on the generated instances among all the  
 300 methods. We further show the colored scene since the used object generator Shap-E is able to generate  
 301 textured 3D objects. The visualization demonstrates our superior performance to produce accurate  
 302 and visual-appealing scene reconstruction from merely a single image in a zero-shot manner.

303 **4.3 Scene Generation of Open-World Images**

304 We further evaluate our method under the open-world images from the BlendSwap dataset and the  
 305 indoor scene dataset Replica. The quantitative comparisons in these two datasets are shown in Tab. 1,  
 306 where we achieve the best performance over all the baseline methods. Note that the performance  
 307 of PanoRecon [11] largely degrades under open-world scene images compared to the performance  
 308 under 3D-Front dataset. The reason is that PanoRecon fails to generalize to the out-of-distribution  
 309 inputs and can only handle the specific image patterns in the trained 3D-Front dataset.

310 The visual comparison is shown in Fig. 7, where we significantly outperform the previous works in  
 311 the generation accuracy and completeness. Specifically, as shown in the 3-nd and 5-th row in Fig. 7,  
 312 our method generates accurate geometries for the table with thin legs and the chair with a complex  
 313 back. While Mesh R-CNN and Total3D can only generate the coarse 3D shapes and also fail to  
 314 estimate accurate layout.

315 **4.4 Ablation Study**

316 **Framework Design.** To evaluate the major components in our methods, we conduct ablations  
 317 under the Replica dataset [48] and report the results in Tab. 2. We first justify the effectiveness  
 318 of introducing Stable-Diffusion for enhancing and inpainting images as shown in ‘W/o  
 319 Stable-Diffusion’, where we directly adopt the segmented instances for shape generation without  
 320 leveraging Stable-Diffusion for enhancing and inpainting them. We then report the performance  
 321 of removing the 2D or 3D matching constraints as shown in ‘W/o 2D-Matching’ and ‘W/o 3D-  
 322 Matching’. The ablation studies demonstrate the effect of each design by significantly improving the  
 323 generation performance. Note that the pose/scale optimization is broken without 3D-Matching since  
 324 the only 2D-Matching does not involve depth information.

325 **Effect of Open-CLIP and RANSAC-like solution.** We further evaluate the effectiveness of  
 326 filtering bad samples with Open-CLIP and the RANSAC-like solution for robust pose / scale  
 327 optimization. The results is shown in Tab. 3, where both components improve the scene re-  
 328 construction accuracy.

329 **5 Conclusion**

330 We introduce deep prior assembly, a novel framework that assembles diverse deep priors from large  
 331 models for scene reconstruction from single images in a zero-shot manner. This approach involves  
 332 breaking down the task into several sub-tasks, each of which is handled by a deep prior. We do not  
 333 rely on any 3D or 2D data-driven training, and provide the key solution on layout estimation and  
 334 occlusion parsing to make all deep priors work together robustly. We report analysis, numerical and  
 335 visual comparisons to show remarkable performance over the latest methods.

| Ablation             | CDL1-S ↓     | CDL1 ↓       | F-Score↑     |
|----------------------|--------------|--------------|--------------|
| W/o Stable-Diffusion | 0.128        | 0.125        | 67.22        |
| W/o 2D-Matching      | 0.124        | 0.121        | 68.42        |
| W/o 3D-Matching      | 0.199        | 0.168        | 56.08        |
| Full                 | <b>0.113</b> | <b>0.110</b> | <b>70.48</b> |

Table 2: Ablation studies on framework designs.

| Open-CLIP | RANSAC | CDL1-S ↓     | CDL1 ↓       | F-Score↑     |
|-----------|--------|--------------|--------------|--------------|
| ✓         | ✗      | 0.121        | 0.118        | 69.28        |
| ✗         | ✓      | 0.129        | 0.123        | 68.92        |
| ✓         | ✓      | <b>0.113</b> | <b>0.110</b> | <b>70.48</b> |

Table 3: Ablation studies on the effect of Open-CLIP filtering and RANSAC-like solution.

343 **References**

- [1] Samir Agarwala, Linyi Jin, Chris Rockwell, and David F Fouhey. Planeformers: From sparse view planes to 3d reconstruction. In *European Conference on Computer Vision*, pages 192–209. Springer, 2022.
- [2] Motilal Agrawal and Larry S Davis. A probabilistic framework for surface reconstruction from multiple images. In *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001*, volume 2, pages II–II. IEEE, 2001.
- [3] Dejan Azinović, Ricardo Martin-Brualla, Dan B Goldman, Matthias Nießner, and Justus Thies. Neural rgb-d surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6290–6301, 2022.
- [4] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. In *International Conference on Learning Representations*, 2021.
- [5] Michael Bleyer, Christoph Rhemann, and Carsten Rother. Patchmatch stereo-stereo matching with slanted support windows. In *Bmvc*, volume 11, pages 1–11, 2011.
- [6] Adrian Broadhurst, Tom W Drummond, and Roberto Cipolla. A probabilistic framework for space carving. In *Proceedings eighth IEEE international conference on computer vision. ICCV 2001*, volume 1, pages 388–393. IEEE, 2001.
- [7] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [8] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3D model repository. *arXiv preprint arXiv:1512.03012*, 2015.
- [9] Chao Chen, Zhizhong Han, Yu-Shen Liu, and Matthias Zwicker. Unsupervised learning of fine structure generation for 3D point clouds by 2D projections matching. In *Proceedings of the ieee/cvf international conference on computer vision*, pages 12466–12477, 2021.
- [10] Tao Chu, Pan Zhang, Qiong Liu, and Jiaqi Wang. Buol: A bottom-up framework with occupancy-aware lifting for panoptic 3d scene reconstruction from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4937–4946, 2023.
- [11] Manuel Dahnert, Ji Hou, Matthias Nießner, and Angela Dai. Panoptic 3d scene reconstruction from a single rgb image. *Advances in Neural Information Processing Systems*, 34:8282–8293, 2021.
- [12] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5828–5839, 2017.
- [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [14] Ainaz Eftekhar, Alexander Sax, Jitendra Malik, and Amir Zamir. Omnidata: A scalable pipeline for making multi-task mid-level vision datasets from 3d scans. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10786–10796, 2021.
- [15] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. *Advances in neural information processing systems*, 27, 2014.
- [16] Yuxin Fang, Wen Wang, Binhuai Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning at scale. *arXiv preprint arXiv:2211.07636*, 2022.
- [17] Huan Fu, Bowen Cai, Lin Gao, Ling-Xiao Zhang, Jiaming Wang, Cao Li, Qixun Zeng, Chengyue Sun, Rongfei Jia, Binqiang Zhao, et al. 3d-front: 3d furnished rooms with layouts and semantics. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10933–10942, 2021.
- [18] Huan Fu, Rongfei Jia, Lin Gao, Mingming Gong, Binqiang Zhao, Steve Maybank, and Dacheng Tao. 3d-future: 3d furniture shape with texture. *International Journal of Computer Vision*, 129:3313–3337, 2021.
- [19] Georgia Gkioxari, Jitendra Malik, and Justin Johnson. Mesh r-cnn. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9785–9795, 2019.
- [20] Eleni Gregoromichelaki, Arash Eshghi, Christine Howes, Gregory J Mills, Ruth Kempson, Julian Hough, Patrick GT Healey, Matthew Purver, et al. Language and cognition as distributed process interactions. In *Proceedings of the 26th Workshop on the Semantics and Pragmatics of Dialogue*, pages 160–171, 2022.

- 403 [21] Haoyu Guo, Sida Peng, Haotong Lin, Qianqian Wang, Guofeng Zhang, Hujun Bao, and Xiaowei  
 404 Zhou. Neural 3d scene reconstruction with the manhattan-world assumption. In *Proceedings*  
 405 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5511–5520,  
 406 2022.
- 407 [22] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked  
 408 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on*  
 409 *computer vision and pattern recognition*, pages 16000–16009, 2022.
- 410 [23] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of*  
 411 *the IEEE international conference on computer vision*, pages 2961–2969, 2017.
- 412 [24] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan  
 413 Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi,  
 414 Ali Farhadi, and Ludwig Schmidt. Openclip, July 2021. If you use this software, please cite it  
 415 as below.
- 416 [25] Heewoo Jun and Alex Nichol. Shap-e: Generating conditional 3d implicit functions. *arXiv*  
 417 *preprint arXiv:2305.02463*, 2023.
- 418 [26] Oğuzhan Fatih Kar, Teresa Yeo, Andrei Atanov, and Amir Zamir. 3d common corruptions  
 419 and data augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
 420 *Pattern Recognition*, pages 18963–18974, 2022.
- 421 [27] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson,  
 422 Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv*  
 423 *preprint arXiv:2304.02643*, 2023.
- 424 [28] Zizhang Li, Xiaoyang Lyu, Yuanyuan Ding, Mengmeng Wang, Yiyi Liao, and Yong Liu.  
 425 Rico: Regularizing the unobservable for indoor compositional reconstruction. *arXiv preprint*  
 426 *arXiv:2303.08605*, 2023.
- 427 [29] Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Mukund Varma T, Zexiang Xu, and Hao  
 428 Su. One-2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization.  
 429 *Advances in Neural Information Processing Systems*, 36, 2024.
- 430 [30] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei  
 431 Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for  
 432 open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023.
- 433 [31] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike  
 434 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining  
 435 approach. *arXiv preprint arXiv:1907.11692*, 2019.
- 436 [32] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano  
 437 Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv*  
 438 *preprint arXiv:2108.01073*, 2021.
- 439 [33] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger.  
 440 Occupancy networks: Learning 3D reconstruction in function space. In *Proceedings of the*  
 441 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4460–4470, 2019.
- 442 [34] B Mildenhall, PP Srinivasan, M Tancik, JT Barron, R Ramamoorthi, and R Ng. NeRF:  
 443 Representing scenes as neural radiance fields for view synthesis. In *European Conference on*  
 444 *Computer Vision*, 2020.
- 445 [35] Yinyu Nie, Xiaoguang Han, Shihui Guo, Yujian Zheng, Jian Chang, and Jian Jun Zhang.  
 446 Total3dunderstanding: Joint layout, object pose and mesh reconstruction for indoor scenes from  
 447 a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
 448 *Recognition*, pages 55–64, 2020.
- 449 [36] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit  
 450 surfaces and radiance fields for multi-view reconstruction. In *Proceedings of the IEEE/CVF*  
 451 *International Conference on Computer Vision*, pages 5589–5599, 2021.
- 452 [37] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove.  
 453 DeepSDF: Learning continuous signed distance functions for shape representation. In *Proceed-*  
 454 *ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 165–174,  
 455 2019.
- 456 [38] Songyou Peng, Chiyu Jiang, Yiyi Liao, Michael Niemeyer, Marc Pollefeys, and Andreas Geiger.  
 457 Shape as points: A differentiable poisson solver. *Advances in Neural Information Processing*  
 458 *Systems*, 34, 2021.
- 459 [39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
 460 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
 461 models from natural language supervision. In *International conference on machine learning*,  
 462 pages 8748–8763. PMLR, 2021.

- 463 [40] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,  
 464 Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified  
 465 text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- 466 [41] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards  
 467 robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE  
 468 transactions on pattern analysis and machine intelligence*, 44(3):1623–1637, 2020.
- 469 [42] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards  
 470 robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE  
 471 Transactions on Pattern Analysis and Machine Intelligence*, 44(3), 2022.
- 472 [43] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
 473 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF  
 474 conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- 475 [44] Szymon Rusinkiewicz and Marc Levoy. Efficient variants of the icp algorithm. In *Proceedings  
 476 third international conference on 3-D digital imaging and modeling*, pages 145–152. IEEE,  
 477 2001.
- 478 [45] Radu Bogdan Rusu, Nico Blodow, and Michael Beetz. Fast point feature histograms (fpfh)  
 479 for 3d registration. In *2009 IEEE international conference on robotics and automation*, pages  
 480 3212–3217. IEEE, 2009.
- 481 [46] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceed-  
 482 ings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113,  
 483 2016.
- 484 [47] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene under-  
 485 standing benchmark suite. In *Proceedings of the IEEE conference on computer vision and pattern  
 486 recognition*, pages 567–576, 2015.
- 487 [48] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J  
 488 Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, et al. The replica dataset: A digital replica of  
 489 indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019.
- 490 [49] Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training  
 491 techniques for clip at scale. *arXiv preprint arXiv:2303.15389*, 2023.
- 492 [50] Xingyuan Sun, Jiajun Wu, Xiuming Zhang, Zhoutong Zhang, Chengkai Zhang, Tianfan Xue,  
 493 Joshua B Tenenbaum, and William T Freeman. Pix3d: Dataset and methods for single-image  
 494 3d shape modeling. In *Proceedings of the IEEE conference on computer vision and pattern  
 495 recognition*, pages 2974–2983, 2018.
- 496 [51] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timo-  
 497 thée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open  
 498 and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 499 [52] Jiepeng Wang, Peng Wang, Xiaoxiao Long, Christian Theobalt, Taku Komura, Lingjie Liu, and  
 500 Wenping Wang. NeuRIS: Neural reconstruction of indoor scenes using normal priors. In *17th  
 501 European Conference on Computer Vision*, pages 139–155. Springer, 2022.
- 502 [53] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang.  
 503 NeuS: Learning neural implicit surfaces by volume rendering for multi-view reconstruction.  
 504 *Advances in Neural Information Processing Systems*, 34:27171–27183, 2021.
- 505 [54] Xin Wen, Junsheng Zhou, Yu-Shen Liu, Hua Su, Zhen Dong, and Zhizhong Han. 3D shape  
 506 reconstruction from 2D images with disentangled attribute flow. In *Proceedings of the IEEE/CVF  
 507 Conference on Computer Vision and Pattern Recognition*, pages 3803–3813, 2022.
- 508 [55] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit  
 509 surfaces. *Advances in Neural Information Processing Systems*, 34:4805–4815, 2021.
- 510 [56] Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. Monosdf:  
 511 Exploring monocular geometric cues for neural implicit surface reconstruction. In *Advances in  
 512 Neural Information Processing Systems*.
- 513 [57] Xiang Zhang, Zeyuan Chen, Fangyin Wei, and Zhuowen Tu. Uni-3d: A universal model for  
 514 panoptic 3d scene reconstruction. In *Proceedings of the IEEE/CVF International Conference  
 515 on Computer Vision*, pages 9256–9266, 2023.
- 516 [58] Junsheng Zhou, Jinsheng Wang, Baorui Ma, Yu-Shen Liu, Tiejun Huang, and Xinlong Wang.  
 517 Uni3d: Exploring unified 3d representation at scale. *arXiv preprint arXiv:2310.06773*, 2023.

518 **A Source Codes**

519 We provide our demonstration code as a part of our supplementary materials. We will release the  
 520 source code, data and instructions upon acceptance.

521 **B Scene Reconstruction from Real Images**

522 **Dataset.** We further evaluate deep prior assembly under the ScanNet [12] using real images. For  
 523 a qualitative comparison with other methods, we select 7 scenes from the test set of ScanNet and  
 524 sample one image from each scene as the input.

525 **Depth Scale and Shift.** We directly adopt the depth scale and shift determined in Replica dataset to  
 526 the scenes of real-world ScanNet dataset, demonstrating that the scale and shift from the randomly  
 527 selected scene generalize well to different scenes or even different datasets where the real depth is  
 528 not available.

529 **Visual Comparisons.** We compare deep prior assembly with the state-of-the-art methods in scene  
 530 reconstruction from single images, e.g., Mesh R-CNN [19] and Total3D [35]. We do not compare  
 531 with PanoRecon [11] here since it fails to generalize to the out-of-distribution inputs and can only  
 532 handle the specific image patterns in the trained 3D-Front dataset as we discussed in Sec. 4.3, Tab. 1  
 533 and Fig. 7 in our paper.

534 We show the visual comparisons in Fig. 9, where we successfully reconstruct scenes from real images  
 535 and significantly outperform the previous works in the reconstruction accuracy and completeness.  
 536 This demonstrates the huge potentials of the assembled deep priors in reconstructing real-world 3D  
 537 scenes. Note that the real-world images are often blurred and corrupted when the camera doesn't  
 538 focus well, e.g., the blurred input image shown in the 5-th row in Fig. 9. While our proposed deep  
 539 prior assembly can also handle these challenging situations due to the powerful and robust deep priors  
 540 from the large visual models.

541 **C More Ablation Studies and Analysis**

542 **C.1 Alternatives on sub-tasks.**

543 We explore the effectiveness of our chosen solutions in each sub-task by comparing them with the  
 544 alternatives. Specifically, we conduct ablations to replace Shap-E [25] with One-2-3-45 [29], replace  
 545 Open-CLIP [39, 24] with EVA-CLIP [49] and replace Omnidata [14] with MiDaS [42] in Tab. 4.  
 546 We visually compare Shap-E with One-2-3-45 for shape generation in Fig. 8, where the results  
 547 demonstrate that Shap-E is a more robust solution for generating 3D models from 2D instances.

548 **C.2 The Effect of Instance Scale**

549 We further conduct ablation studies to explore the effect of the instance scales to the generation  
 550 qualities of Shap-E as described in “**Enhance and Inpaint 2D instances**” of Sec. 3.2 in our paper.  
 551 We provide an visual comparison of the generations with different instance scales in Fig. 10. The  
 552 results show that Shap-E is quite sensitive to the scale of instances in the images, where a too small  
 553 or too large scale will lead to inaccurate generations with unreliable geometries and appearances. We  
 554 set the scale to 6 where the Shap-E performs the best in shape generation from instance images.

| Ablation   | CDL1-S ↓     | CDL1 ↓       | F-Score↑     |
|------------|--------------|--------------|--------------|
| One-2-3-45 | 0.123        | 0.122        | 67.98        |
| EVA-CLIP   | <b>0.113</b> | 0.111        | 70.41        |
| MiDaS      | 0.120        | 0.119        | 68.71        |
| Ours       | <b>0.113</b> | <b>0.110</b> | <b>70.48</b> |



Table 4: Ablation studies on the sub-task alternatives. Figure 8: Ablation on shape generation alternatives.

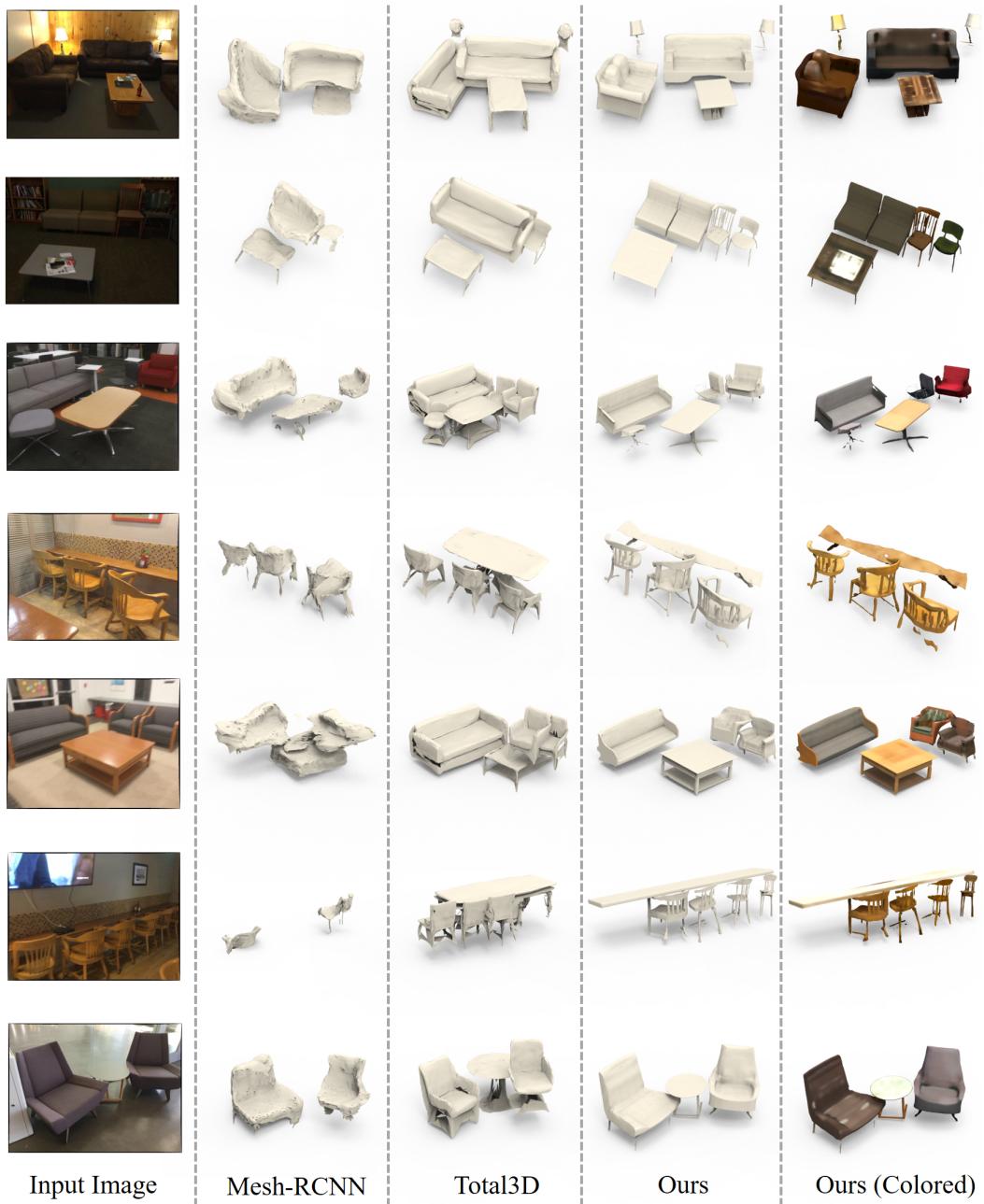


Figure 9: Comparisons of the Scene reconstructions under the real-captured scene images from ScanNet dataset.

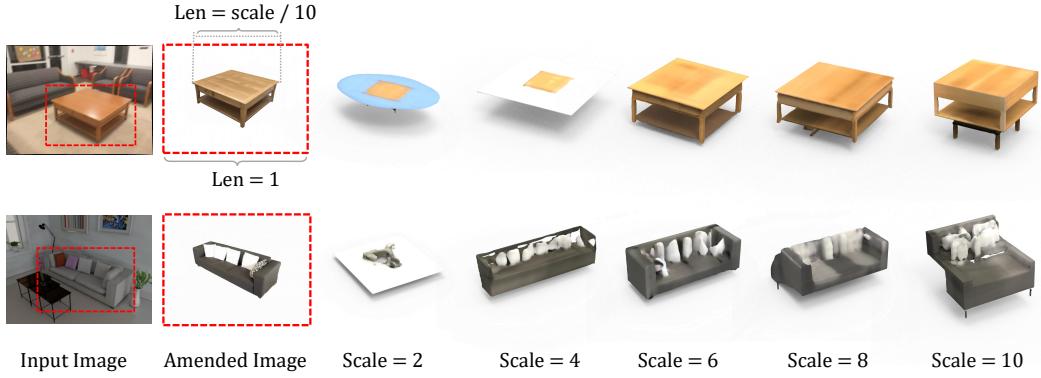


Figure 10: The ablation study on the instance scale. We select one instance for each input image and show the amended instance images. The generations obtained with Shap-E under different instance scales are visualized on the right.

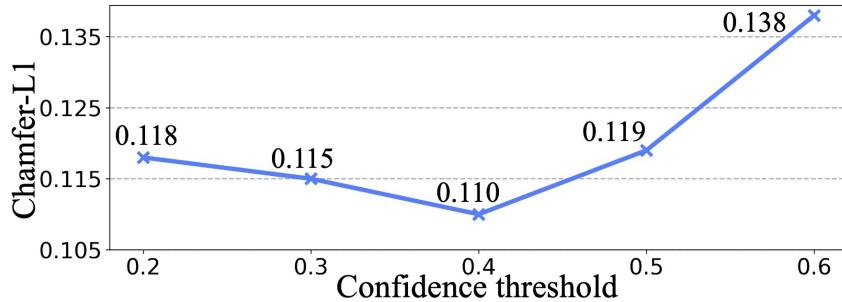


Figure 11: The ablation study on the confidence threshold.

### 555 C.3 The Effect of Confidence Threshold.

556 We further conduct ablations to evaluate the effect of confidence threshold  $\sigma$  as described in “**Detect**  
 557 **and Segment 2D instances.**” of Sec.3.2 in our paper. As shown in Fig. 11, a too large  $\sigma$  will drop too  
 558 many instances and a too small  $\sigma$  struggles to filter bad instances. We choose  $\sigma = 0.4$  as the suitable  
 559 confidence threshold.

### 560 C.4 The Analysis of Depth Scale and Shift

561 For the 3D-Front, Replica and BlendSwap dataset, we leverage one randomly selected pair of  
 562 predicted and real (or scanned) depth as the reference to solve the scale and shift which are then used  
 563 in other scenes in the dataset. We also found that the scale and shift generalize well across datasets,  
 564 since these scale and shift work well on the real-captured images where real depth are not available.  
 565 For instance, the results under ScanNet dataset in Fig. 9 are directly obtained with the depth scale  
 566 and shift determined in Replica dataset.

## 567 D More Comparisons with Data-Driven Reconstruction Methods

568 We additionally compare our method with SOTA data-driven scene reconstruction works PanoRecon  
 569 [11], BUOL [10] and Uni-3D [57]. We show the visual comparisons under 3D-Front and ScanNet  
 570 datasets in Fig. 12, where our method achieves better and more visual-appealing results under both  
 571 3D-Front and ScanNet datasets. Specifically, our method significantly outperforms other methods  
 572 using real-world images in ScanNet. The reason is that all the three methods are trained under  
 573 3D-Front, and struggle to generalize on real-world images.

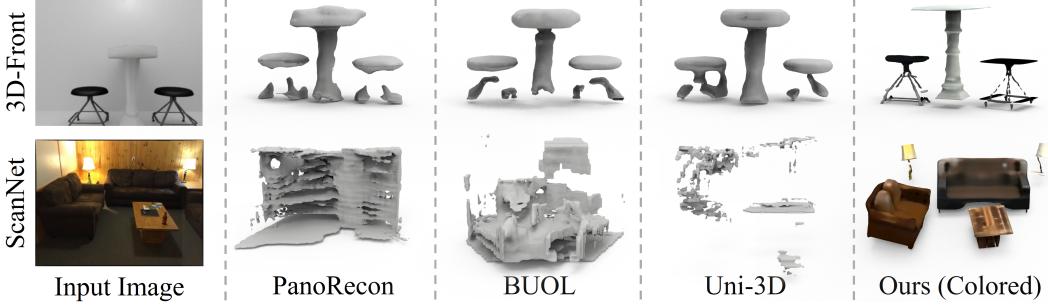


Figure 12: Visual comparisons under 3D-Front and ScanNet dataset.

## 574 E Background Reconstruction

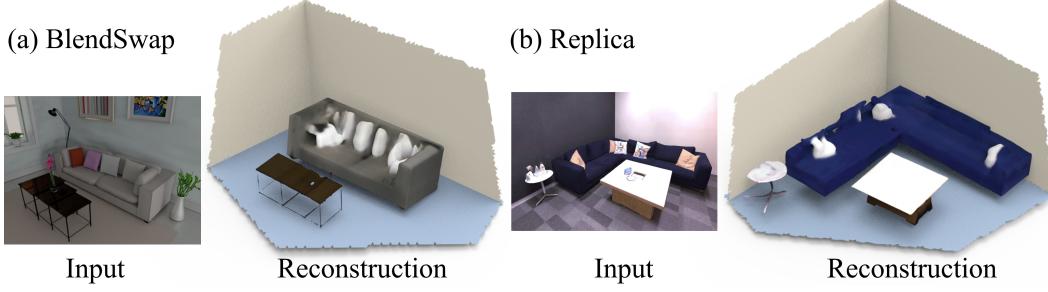


Figure 13: Scene reconstructions with backgrounds.

575 We demonstrate that deep prior assembly can also reconstruct the background geometries from the  
 576 scene images. We show two scene reconstructions with backgrounds (i.e. floor, wall) in BlendSwap  
 577 and Replica datasets in Fig. 13. The backgrounds are achieved by fitting planes to the projected  
 578 background depth points in a similar way as our pose/scale optimization algorithm. We then cull the  
 579 geometries outside of the view frustum for a clear visualization.

## 580 F Efficiency Analysis

581 We further evaluate the efficiency of our proposed deep prior assembly by reporting the average  
 582 runtime of each sub-pipeline in our framework. The results in Tab. 5 show that reconstructing a  
 583 scene from a single image takes less than 3 minutes in total, where the inference of Grounded-SAM,  
 584 Open-CLIP and Omnidata takes only about 1 second. The most time consuming parts include the  
 585 StableDiffusion, Shap-E and the RANSAC-like pose/scale optimization. For these three parts, we  
 586 process all instances of the scene in parallel, resulting in significant time savings.

|          | G-SAM | Sta.-Diff. | Open-CLIP | Omnidata | Shap-E | RANSAC-Opti | Total |
|----------|-------|------------|-----------|----------|--------|-------------|-------|
| Time (s) | 0.93  | 33.6       | 0.05      | 0.21     | 39.2   | 97.2        | 171.2 |

Table 5: Runtime of each sub-pipeline on a RTX3090 GPU.

## 587 G Evaluation Details

588 **PanoRecon.** For evaluating PanoRecon [11], we adopt the official code and pretrained models for  
 589 inference and directly report the performance under 3D-Front [17] dataset by evaluating the metrics  
 590 (e.g. Chamfer Distance and F-Score) between the reconstructions and the ground truth meshes. For

591 the Replica [48] and BlendSwap [3] datasets where the scene location and orientation do not match  
592 the 3D-Front dataset where the PanoRecon is trained, we first normalize the center of the predicted  
593 scenes to the ground truth scenes and then register the predicted scene reconstructions to the ground  
594 truth scenes. Specifically, we first predict an initial alignment by a global registration algorithm based  
595 on feature matching [45] with RANSAC and then leverage ICP (Iterative Closest Point) registration  
596 algorithm [44] to obtain the fine registration based on the initial alignment. The metrics are reported  
597 with the registered reconstructions and the ground truth meshes.

598 **Evaluating Total3D.** We leverage the official code and the pretrained models for predicting scene  
599 reconstructions with Total3D [35]. We evaluate Total3D under 3D-Front, Replica and BlendSwap  
600 with a similar way as we evaluating PanoRecon to first normalize the predicted scenes and register  
601 them to the ground truth ones before computing metrics. Total3D only predicts 3D objects and do not  
602 generate backgrounds (e.g. wall and floor). Therefore, we further report the single-direction chamfer  
603 distance from the generated objects to the ground truth meshes, i.e., Chamfer-L1 (S), to only evaluate  
604 the accuracy of the generated objects. Note that Total3D can generate the scene layout which can  
605 roughly represent the background, however, we find that Total3D generates layouts with large errors  
606 on all the three datasets. Therefore we do not keep the layout of Total3D for evaluation.

607 **Evaluating Mesh R-CNN.** We adopt the official code and the pretrained models for predicting scene  
608 reconstructions with Mesh R-CNN [19]. We notice that Mesh R-CNN produces reconstructions  
609 with larger scales than the predictions of other methods and the ground truths. Therefore, we first  
610 normalize both the center and scale of the predicted scenes to the ground truth scenes and then register  
611 the predicted scene reconstructions to the ground truth scenes with a similar way as we evaluate  
612 Total3D.

613 **Evaluating Deep Prior Assembly.** We evaluate our proposed deep prior assembly in the same  
614 way. We follow the same settings as we evaluate other baselines to first normalize the center of the  
615 predicted scenes to the ground truth scenes, and then register the predicted scenes to the ground  
616 truth scenes. The background points (e.g. wall and floor) of deep prior assembly are obtained by  
617 back-projecting the background depth maps, i.e., the areas where no instances exists.

## 618 H Limitation

619 One limitation of our method is that it may sometimes produce reconstructions with 3D instance  
620 models that are not perfectly aligned with the 2D instance segmentations in the input images. For  
621 example, the left chair generation in the last row of Fig. 9 exhibits a different color from the 2D  
622 instance of the input image. The reason is that we use Stable-Diffusion to enhance and inpaint the 2D  
623 instances, and then leverage Shap-E to generate 3D reconstructions. This process can introduce some  
624 randomness in the generated textures. The randomness primarily affects the appearances, while the  
625 geometries remain accurate. However, we justify that most reconstructions can faithfully recover the  
626 consistent scene appearances and geometries from the input images.

## 627 I Video

628 We provide a video containing the visualization of optimization process, the visualization of the scene  
629 reconstruction results, and comparisons under the synthetic and real-captured datasets as a part of our  
630 supplementary materials.

631 **NeurIPS Paper Checklist**

632 **1. Claims**

633 Question: Do the main claims made in the abstract and introduction accurately reflect the  
634 paper's contributions and scope?

635 Answer: [Yes]

636 Justification: Our main claims made in the abstract and introduction accurately reflect the  
637 paper's contributions and scope.

638 Guidelines:

- 639 • The answer NA means that the abstract and introduction do not include the claims  
640 made in the paper.
- 641 • The abstract and/or introduction should clearly state the claims made, including the  
642 contributions made in the paper and important assumptions and limitations. A No or  
643 NA answer to this question will not be perceived well by the reviewers.
- 644 • The claims made should match theoretical and experimental results, and reflect how  
645 much the results can be expected to generalize to other settings.
- 646 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
647 are not attained by the paper.

648 **2. Limitations**

649 Question: Does the paper discuss the limitations of the work performed by the authors?

650 Answer: [Yes]

651 Justification: We analysis the limitations of our method in Sec.H of the appendix.

652 Guidelines:

- 653 • The answer NA means that the paper has no limitation while the answer No means that  
654 the paper has limitations, but those are not discussed in the paper.
- 655 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 656 • The paper should point out any strong assumptions and how robust the results are to  
657 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
658 model well-specification, asymptotic approximations only holding locally). The authors  
659 should reflect on how these assumptions might be violated in practice and what the  
660 implications would be.
- 661 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
662 only tested on a few datasets or with a few runs. In general, empirical results often  
663 depend on implicit assumptions, which should be articulated.
- 664 • The authors should reflect on the factors that influence the performance of the approach.  
665 For example, a facial recognition algorithm may perform poorly when image resolution  
666 is low or images are taken in low lighting. Or a speech-to-text system might not be  
667 used reliably to provide closed captions for online lectures because it fails to handle  
668 technical jargon.
- 669 • The authors should discuss the computational efficiency of the proposed algorithms  
670 and how they scale with dataset size.
- 671 • If applicable, the authors should discuss possible limitations of their approach to  
672 address problems of privacy and fairness.
- 673 • While the authors might fear that complete honesty about limitations might be used by  
674 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
675 limitations that aren't acknowledged in the paper. The authors should use their best  
676 judgment and recognize that individual actions in favor of transparency play an impor-  
677 tant role in developing norms that preserve the integrity of the community. Reviewers  
678 will be specifically instructed to not penalize honesty concerning limitations.

679 **3. Theory Assumptions and Proofs**

680 Question: For each theoretical result, does the paper provide the full set of assumptions and  
681 a complete (and correct) proof?

682 Answer: [NA]

683 Justification: The paper does not include theoretical results.

684 Guidelines:

- 685 • The answer NA means that the paper does not include theoretical results.
- 686 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 688 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 689 • The proofs can either appear in the main paper or the supplemental material, but if 690 they appear in the supplemental material, the authors are encouraged to provide a short 691 proof sketch to provide intuition.
- 692 • Inversely, any informal proof provided in the core of the paper should be complemented 693 by formal proofs provided in appendix or supplemental material.
- 694 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 695 4. Experimental Result Reproducibility

696 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
697 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
698 of the paper (regardless of whether the code and data are provided or not)?

699 Answer: [Yes]

700 Justification: We provide the detailed information in reproducing our methods in Sec.3,  
701 Sec.4 of the main paper and the appendix. We also provide a demonstration code of our  
702 method in the supplementary materials.

703 Guidelines:

- 704 • The answer NA means that the paper does not include experiments.
- 705 • If the paper includes experiments, a No answer to this question will not be perceived  
706 well by the reviewers: Making the paper reproducible is important, regardless of  
707 whether the code and data are provided or not.
- 708 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
709 to make their results reproducible or verifiable.
- 710 • Depending on the contribution, reproducibility can be accomplished in various ways.  
711 For example, if the contribution is a novel architecture, describing the architecture fully  
712 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
713 be necessary to either make it possible for others to replicate the model with the same  
714 dataset, or provide access to the model. In general, releasing code and data is often  
715 one good way to accomplish this, but reproducibility can also be provided via detailed  
716 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
717 of a large language model), releasing of a model checkpoint, or other means that are  
718 appropriate to the research performed.
- 719 • While NeurIPS does not require releasing code, the conference does require all submissions  
720 to provide some reasonable avenue for reproducibility, which may depend on the  
721 nature of the contribution. For example
  - 722 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
723 to reproduce that algorithm.
  - 724 (b) If the contribution is primarily a new model architecture, the paper should describe  
725 the architecture clearly and fully.
  - 726 (c) If the contribution is a new model (e.g., a large language model), then there should  
727 either be a way to access this model for reproducing the results or a way to reproduce  
728 the model (e.g., with an open-source dataset or instructions for how to construct  
729 the dataset).
  - 730 (d) We recognize that reproducibility may be tricky in some cases, in which case  
731 authors are welcome to describe the particular way they provide for reproducibility.  
732 In the case of closed-source models, it may be that access to the model is limited in  
733 some way (e.g., to registered users), but it should be possible for other researchers  
734 to have some path to reproducing or verifying the results.

#### 735 5. Open access to data and code

736 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
737 tions to faithfully reproduce the main experimental results, as described in supplemental  
738 material?

739 Answer: [Yes]

740 Justification: We provide our demonstration code as a part of our supplementary materials.  
741 We will release the source code, data and instructions upon acceptance.

742 Guidelines:

- 743 • The answer NA means that paper does not include experiments requiring code.
- 744 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 745 • While we encourage the release of code and data, we understand that this might not be  
746 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
747 including code, unless this is central to the contribution (e.g., for a new open-source  
748 benchmark).
- 749 • The instructions should contain the exact command and environment needed to run to  
750 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 751 • The authors should provide instructions on data access and preparation, including how  
752 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 753 • The authors should provide scripts to reproduce all experimental results for the new  
754 proposed method and baselines. If only a subset of experiments are reproducible, they  
755 should state which ones are omitted from the script and why.
- 756 • At submission time, to preserve anonymity, the authors should release anonymized  
757 versions (if applicable).
- 758 • Providing as much information as possible in supplemental material (appended to the  
759 paper) is recommended, but including URLs to data and code is permitted.

## 762 6. Experimental Setting/Details

763 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
764 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
765 results?

766 Answer: [Yes]

767 Justification: We provide the training and testing details in the experiment section (Sec.4)  
768 and the appendix.

769 Guidelines:

- 770 • The answer NA means that the paper does not include experiments.
- 771 • The experimental setting should be presented in the core of the paper to a level of detail  
772 that is necessary to appreciate the results and make sense of them.
- 773 • The full details can be provided either with the code, in appendix, or as supplemental  
774 material.

## 775 7. Experiment Statistical Significance

776 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
777 information about the statistical significance of the experiments?

778 Answer: [No]

779 Justification: We report the average performance as the experimental results.

780 Guidelines:

- 781 • The answer NA means that the paper does not include experiments.
- 782 • The authors should answer “Yes” if the results are accompanied by error bars, confi-  
783 dence intervals, or statistical significance tests, at least for the experiments that support  
784 the main claims of the paper.
- 785 • The factors of variability that the error bars are capturing should be clearly stated (for  
786 example, train/test split, initialization, random drawing of some parameter, or overall  
787 run with given experimental conditions).

- 788           • The method for calculating the error bars should be explained (closed form formula,  
 789            call to a library function, bootstrap, etc.)  
 790           • The assumptions made should be given (e.g., Normally distributed errors).  
 791           • It should be clear whether the error bar is the standard deviation or the standard error  
 792            of the mean.  
 793           • It is OK to report 1-sigma error bars, but one should state it. The authors should  
 794            preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
 795            of Normality of errors is not verified.  
 796           • For asymmetric distributions, the authors should be careful not to show in tables or  
 797            figures symmetric error bars that would yield results that are out of range (e.g. negative  
 798            error rates).  
 799           • If error bars are reported in tables or plots, The authors should explain in the text how  
 800            they were calculated and reference the corresponding figures or tables in the text.

801           **8. Experiments Compute Resources**

802           Question: For each experiment, does the paper provide sufficient information on the com-  
 803           puter resources (type of compute workers, memory, time of execution) needed to reproduce  
 804           the experiments?

805           Answer: [Yes]

806           Justification: The computer resources needed to reproduce the experiments are provided in  
 807           Sec.F of the appendix.

808           Guidelines:

- 809           • The answer NA means that the paper does not include experiments.  
 810           • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
 811            or cloud provider, including relevant memory and storage.  
 812           • The paper should provide the amount of compute required for each of the individual  
 813            experimental runs as well as estimate the total compute.  
 814           • The paper should disclose whether the full research project required more compute  
 815            than the experiments reported in the paper (e.g., preliminary or failed experiments that  
 816            didn't make it into the paper).

817           **9. Code Of Ethics**

818           Question: Does the research conducted in the paper conform, in every respect, with the  
 819           NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

820           Answer: [Yes]

821           Justification: The research conducted in the paper conform, in every respect, with the  
 822           NeurIPS Code of Ethics.

823           Guidelines:

- 824           • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.  
 825           • If the authors answer No, they should explain the special circumstances that require a  
 826            deviation from the Code of Ethics.  
 827           • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
 828            eration due to laws or regulations in their jurisdiction).

829           **10. Broader Impacts**

830           Question: Does the paper discuss both potential positive societal impacts and negative  
 831           societal impacts of the work performed?

832           Answer: [Yes]

833           Justification: We discuss the applications and potential impacts of our method in the  
 834           introduction.

835           Guidelines:

- 836           • The answer NA means that there is no societal impact of the work performed.  
 837           • If the authors answer NA or No, they should explain why their work has no societal  
 838            impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We use the open-sourced datasets under their licenses.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- 893           • For existing datasets that are re-packaged, both the original license and the license of  
894           the derived asset (if it has changed) should be provided.  
895           • If this information is not available online, the authors are encouraged to reach out to  
896           the asset's creators.

897           **13. New Assets**

898           Question: Are new assets introduced in the paper well documented and is the documentation  
899           provided alongside the assets?

900           Answer: [NA]

901           Justification: The paper does not release new assets.

902           Guidelines:

- 903           • The answer NA means that the paper does not release new assets.  
904           • Researchers should communicate the details of the dataset/code/model as part of their  
905           submissions via structured templates. This includes details about training, license,  
906           limitations, etc.  
907           • The paper should discuss whether and how consent was obtained from people whose  
908           asset is used.  
909           • At submission time, remember to anonymize your assets (if applicable). You can either  
910           create an anonymized URL or include an anonymized zip file.

911           **14. Crowdsourcing and Research with Human Subjects**

912           Question: For crowdsourcing experiments and research with human subjects, does the paper  
913           include the full text of instructions given to participants and screenshots, if applicable, as  
914           well as details about compensation (if any)?

915           Answer: [NA]

916           Justification: The paper does not involve crowdsourcing nor research with human subjects.

917           Guidelines:

- 918           • The answer NA means that the paper does not involve crowdsourcing nor research with  
919           human subjects.  
920           • Including this information in the supplemental material is fine, but if the main contribu-  
921           tion of the paper involves human subjects, then as much detail as possible should be  
922           included in the main paper.  
923           • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
924           or other labor should be paid at least the minimum wage in the country of the data  
925           collector.

926           **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human  
927           Subjects**

928           Question: Does the paper describe potential risks incurred by study participants, whether  
929           such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
930           approvals (or an equivalent approval/review based on the requirements of your country or  
931           institution) were obtained?

932           Answer: [NA]

933           Justification: The paper does not involve crowdsourcing nor research with human subjects.

934           Guidelines:

- 935           • The answer NA means that the paper does not involve crowdsourcing nor research with  
936           human subjects.  
937           • Depending on the country in which research is conducted, IRB approval (or equivalent)  
938           may be required for any human subjects research. If you obtained IRB approval, you  
939           should clearly state this in the paper.  
940           • We recognize that the procedures for this may vary significantly between institutions  
941           and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
942           guidelines for their institution.  
943           • For initial submissions, do not include any information that would break anonymity (if  
944           applicable), such as the institution conducting the review.