

---

# Neural Universal Scene Descriptors

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Although recent progress in generative modeling has produced models capable of  
2 generating high-quality images conditioned on multiple modalities, there exists  
3 no common portable representation format for specifying conditioning signals.  
4 Instead, conditioning techniques are usually tailor-made for specific model archi-  
5 tectures, and limit the user to a small set of control signals. In addition, common  
6 approaches are not *object centric*, meaning the user is not able to control individual  
7 objects in the image, and changing the conditioning signal leads to global changes.  
8 In contrast, the computer graphics community has developed standards like the Uni-  
9 versal Scene Descriptor (USD), which represents scenes and objects in a structured,  
10 hierarchical manner. Inspired by USD, we propose the “Neural Universal Scene  
11 Descriptor” (Neural USD), a flexible conditioning structure that accommodates di-  
12 verse signals, minimizes model-specific constraints, and enables per-object control  
13 over appearance, geometry, and pose. We further apply a fine-tuning approach that  
14 ensures disentangled control signals and evaluate key design considerations for a  
15 universal conditioning format, demonstrating how Neural USD enables iterative  
16 and incremental workflows.

## 1 Introduction

17  
18 The surge in relevance of visual generative models has led to the development of a wide range  
19 of conditioning approaches. These approaches enable control over generated outputs by allow-  
20 ing users to guide the generation process using textual prompts, reference images, or other forms  
21 of input. However, conditioning choices are often tailor-made for specific model architectures,  
22 or limit the user to global scene edits, restricting portability across models and limiting users from  
23 performing object-level, incremental updates to their content.  
24  
25 Several approaches have been proposed to address the challenge of conditioning and control-  
26 ling visual generative models. One area of study investigates how models can be conditioned on  
27 depth maps, edge maps, segmentation maps, and other conditioning signals [Zhang et al., 2023,  
28 Mou et al., 2023, Avrahami et al., 2022, Li et al., 2023b]. These approaches allow the user to control  
29 one aspect of the scene at a time, and can result in global scene changes as a result of local

30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
Several approaches have been proposed to address the challenge of conditioning and control-  
ling visual generative models. One area of study investigates how models can be conditioned on  
depth maps, edge maps, segmentation maps, and other conditioning signals [Zhang et al., 2023,  
Mou et al., 2023, Avrahami et al., 2022, Li et al., 2023b]. These approaches allow the user to control  
one aspect of the scene at a time, and can result in global scene changes as a result of local

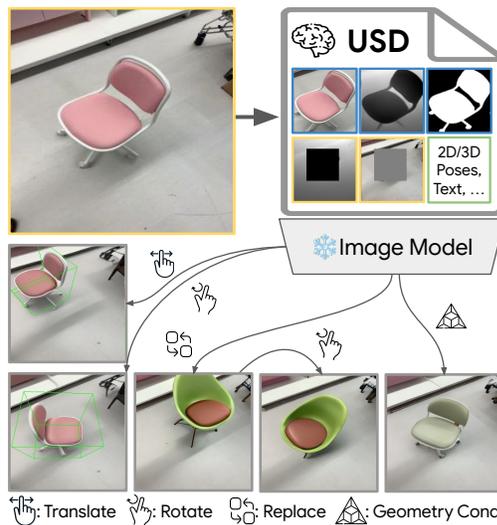


Figure 1: Neural USD enables computer graphics-style control of image models. A Neural USD represents an image as assets with appearance, geometry, and pose. Fine-tuning adapts pre-trained models to these signals while keeping appearance and geometry pose-invariant.

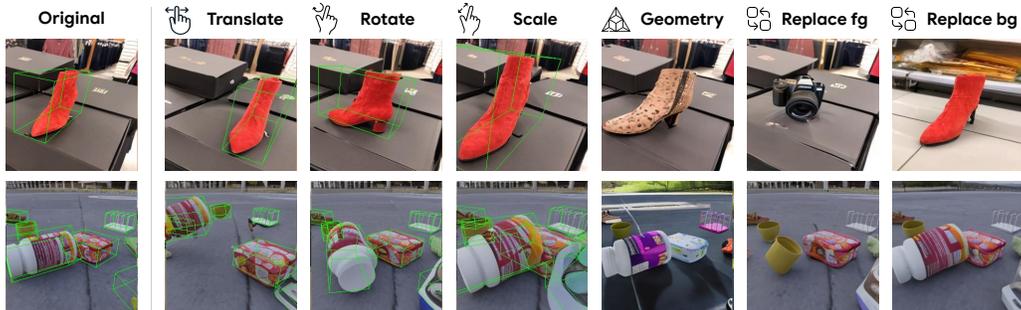


Figure 2: Neural USD enables users to perform a variety of pose, appearance, and geometry modifications to foreground objects and the background, which is simply an additional asset in the Neural USD. Object replacement only requires modifying the appearance and/or geometry for an existing object.

40 conditioning signal edits. Another line of work uses text prompts to guide image generation and  
 41 editing [Brooks et al., 2022, Rombach et al., 2022]. While these approaches generate impressive  
 42 results, they often limit what a user is able to express due to the challenge of describing complex  
 43 scene layouts with text. Recent approaches propose object-centric conditioning formats to guide  
 44 image generation [Bhat et al., 2023b, Wu et al., 2024, Liu et al., 2023a, Michel et al., 2023]. Despite  
 45 their innovations, these methods often struggle with handling multiple objects, or are limited in  
 46 supporting conditioning modalities beyond reference images.

47 To address these challenges, we introduce Neural Universal Scene Descriptors (Neural USD): an  
 48 object-centric conditioning standard that enables precise control of geometry, appearance, and object  
 49 poses within generative models. The Neural USD is defined as an XML-style format consisting of  
 50 per-object attributes ranging from images to geometry to text. Each of these modalities is tokenized  
 51 into a sequence of conditioning vectors and passed to downstream generative models for conditioning.  
 52 After fine-tuning the image models with the Neural USD data, we can manually edit the objects and  
 53 backgrounds in the image using a simple recipe:  $\mathcal{I} = \text{Decode}(\text{Edit}(\text{USD}))$ . However, naive training  
 54 on such a representation would cause challenges, as the model would have difficulty disentangling  
 55 pose and appearance attributes of the conditioning signal, empirically resulting in poor object control.  
 56 We solve this by training the model to reconstruct target images in a video sequence, conditioning on  
 57 geometry and appearance from source images in the same video sequence, and the target pose from  
 58 the corresponding target image. This results in robust object pose, geometry, and appearance control.

59 In summary, our contributions are as follows:

- 60 1. Neural USD: an object-centric conditioning format for generative models that provides  
 61 precise control over object position, appearance, and geometry.
- 62 2. A compact and efficient representation compatible with any model architecture supporting  
 63 vector-based conditioning - e.g. diffusion models [Sohl-Dickstein et al., 2015, Ho et al.,  
 64 2020], DiTs [Peebles and Xie, 2023], transformers [Vaswani et al., 2017, Yu et al., 2022] -  
 65 thereby facilitating cross-model portability.
- 66 3. Demonstrations of the applicability of Neural USD to real-world and synthetic data sets,  
 67 showcasing its ability to embed arbitrary objects sourced from the Internet for rapid scene  
 68 testing.

69 By addressing these limitations and introducing a unified framework for object-centric conditioning,  
 70 Neural USD sets the stage for the next generation of controllable and interoperable generative models.

## 71 2 Related work

72 Recent work in object-centric learning decomposes visual scenes into distinct object representations  
 73 for structured, interpretable image generation. Slot-based methods such as Slot Attention [Locatello  
 74 et al., 2020], SLATE [Singh et al., 2021], and STEVE [Singh et al., 2022b] model scenes as in-  
 75 dependent entities, with refinements like LSD [Jiang et al., 2023] and Slot Diffusion [Wu et al.,  
 76 2023] improving disentanglement. These representations support tasks including attribute manipu-  
 77 lation [Singh et al., 2022a], motion modeling [Seitzer et al., 2023], and 3D pose estimation [Jabri

78 et al., 2023]; OSRT [Sajjadi et al., 2022] further addresses global camera pose. Yet such models  
79 struggle with real-world data. Our approach leverages self-supervised visual encoders with large-scale  
80 pre-trained diffusion models, enabling scalable object-centric learning in realistic settings.

81 Personalized image generation has progressed from test-time fine-tuning (DreamBooth [Ruiz et al.,  
82 2022], Textual Inversion [Gal et al., 2022]) to zero-shot personalization (InstantBooth [Shi et al.,  
83 2023a], ZeroShotBooth [Jia et al., 2023], BLIP-Diffusion [Li et al., 2023a], ELITE [Wei et al., 2023],  
84 InstantID [Wang et al., 2024b], FastComposer [Xiao et al., 2023]). While effective, most produce  
85 single-subject images without spatial control. Exceptions such as VisualComposer [Parmar et al.,  
86 2025], TokenVerse [Garibi et al., 2025], and Video Alchemist [Chen et al., 2025] allow multi-entity  
87 composition, but with limited control. Subject-Diffusion [Ma et al., 2023] introduces 2D bounding-  
88 box conditioning but lacks 3D pose control. We extend controllability by incorporating object poses,  
89 enabling structured multi-object generation and manipulation.

90 Spatial control in diffusion models is pursued through bounding boxes or segmentation masks.  
91 Strategies include prompt manipulation [Kawar et al., 2022, Ge et al., 2023, Brooks et al., 2022],  
92 attention adjustments [Xie et al., 2023, Kim et al., 2023, Chen et al., 2023, Chefer et al., 2023, Feng  
93 et al., 2022, Hertz et al., 2022, Cao et al., 2023], and latent editing [Epstein et al., 2023, Shi et al.,  
94 2023c, Luo et al., 2023]. Fine-tuned approaches add spatial conditioning [Gafni et al., 2022, Avrahami  
95 et al., 2022, Yang et al., 2022, Hu et al., 2023, Xu et al., 2023, Goel et al., 2023]. GLIGEN [Li  
96 et al., 2023b] introduces attention layers for box conditioning, InstanceDiffusion [Wang et al., 2024c]  
97 supports masks and scribbles, and ControlNet [Zhang et al., 2023] incorporates depth and normals;  
98 Boxinator [Wang et al., 2024a] extends these ideas to video.

99 3D-aware image generation pursues structured scene synthesis. GAN-based methods use explicit 3D  
100 representations such as radiance fields [Chan et al., 2020, Gu et al., 2021, Chan et al., 2021, Schwarz  
101 et al., 2020, Niemeyer and Geiger, 2020, Xu et al., 2022] and meshes [Chen et al., 2019, 2021, Gao  
102 et al., 2022, Pavllo et al., 2020, 2021]. Diffusion-based approaches [Shi et al., 2023b, Liu et al., 2023b,  
103 Poole et al., 2022, Wang et al., 2022, Lin et al., 2022, Kant et al., 2024, Melas-Kyriazi et al., 2023,  
104 Watson et al., 2022] transfer 2D knowledge into 3D. 3DiM [Watson et al., 2022] and Zero-1-to-3 [Liu  
105 et al., 2023a] leverage multiview training but remain object-centric. More recent methods address  
106 multi-object real-world scenes [Sargent et al., 2023, Pandey et al., 2023, Yenphraphai et al., 2024,  
107 Alzayer et al., 2024]; OBJect-3DIT [Michel et al., 2023] enables language-guided editing but is  
108 synthetic-data limited. LooseControl [Bhat et al., 2023b] uses 3D bounding boxes as depth maps for  
109 pose control, but is not directly applicable to editing.

110 We unify these directions by enabling both 2D and 3D spatial conditioning in pre-trained diffusion  
111 models, with support for appearance and geometry inputs. Using bounding boxes as control signals,  
112 our approach enables fine-grained object pose manipulation—including rotation and occlusion—while  
113 scaling to multiple modalities and offering a general recipe for incorporating new ones.

## 114 **3 Method**

115 We propose Neural USD as an object-centric representation of a scene, composed of appearance,  
116 geometry, and pose representations. Image models are trained to reconstruct target objects defined  
117 in the Neural USD by using paired images extracted from video sequences. We additionally apply  
118 modality dropout. This allows the model to learn disentangled appearance, geometry, and pose  
119 representations. The resulting model allows for fine-grained control of objects in the scene. Such  
120 a conditioning format draws parallels to the intuitive object-centric workflows used in computer  
121 graphics programs such as Blender [Blender Online Community, 2018].

### 122 **3.1 Data**

123 In this work, we explore datasets with 2D and 3D annotations readily available. Obtaining tracked 2D  
124 bounding boxes for video is a problem with promising solutions [Li et al., 2024]. However, obtaining  
125 3D bounding box annotations at scale is still an open challenge, but may be addressed in the near  
126 future given improvements in SLAM, point tracking, and depth estimation [Zhang et al., 2024, Bhat  
127 et al., 2023a, Yang et al., 2024, Doersch et al., 2023].

128 We compose the Neural USD dataset by applying separate annotation models to the original datasets.  
129 We acquire depth annotations by applying ZoeDepth [Bhat et al., 2023a] then crop and normalize

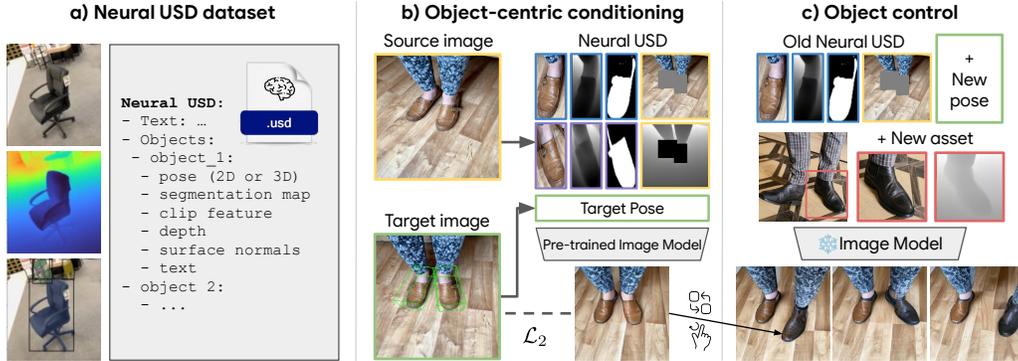


Figure 3: Neural USD Overview. a) A Neural USD consists of assets with multiple modalities: appearance, geometry, and pose. b) Pre-trained image models fine-tune on Neural USD, encoding appearance and geometry from a source image and pose from a target image to reconstruct the target. c) At inference, objects’ poses, geometry, and appearance can be modified, including the background.

130 per-object depth maps. Object masks are computed by applying SAM [Kirillov et al., 2023] with  
 131 bounding box conditioning. Additional conditioning signals such as surface normals and pointclouds  
 132 can be extracted with open source models [Yang et al., 2023], though we leave this for future work.

### 133 3.2 Assets in Neural USD

134 We borrow the nomenclature of Neural Assets [Wu et al., 2024] to describe the components of the  
 135 Neural USD. A Neural USD is defined as a list of  $N$  assets  $\{\hat{a}_1, \dots, \hat{a}_N\}$ , where each asset  $\hat{a}_i$  is  
 136 defined as a tuple of attributes such as 2D or 3D bounding box coordinates  $\mathcal{P}_i^{2D}, \mathcal{P}_i^{3D}$ , appearance  
 137 descriptors such as image crops or clip embeddings  $\mathcal{A}_i$ , geometry signals from depth images, masks,  
 138 pointclouds, and surface normals  $\mathcal{G}_i$ , or even text  $\mathcal{T}_i$ . The resulting asset can be defined as the tuple  
 139  $\hat{a}_i = (\mathcal{P}_i^{2D}, \mathcal{P}_i^{3D}, \mathcal{A}_i, \mathcal{G}_i, \mathcal{T}_i, \dots)$ .

### 140 3.3 Encoding assets

141 To make the Neural USD a compatible conditioning format for arbitrary downstream models, it  
 142 must first be encoded into a continuous vector representation such that the encoded appearance and  
 143 geometry descriptors can be defined as continuous vectors  $\mathcal{A}_i^{\text{emb}}, \mathcal{G}_i^{\text{emb}} \in \mathbb{R}^{K \times D}$ , and pose embeddings  
 144 as  $\mathcal{P}_i^{\text{emb}, 2D}, \mathcal{P}_i^{\text{emb}, 3D} \in \mathbb{R}^{D'}$ . This token representation enables the fine-tuning of arbitrary model  
 145 architectures with the Neural USD encoding, as well as separate control of pose, geometry, and  
 146 appearance. We now describe how appearance, geometry, and poses are encoded in this format.

#### 147 3.3.1 Appearance and geometry encoding

148 Obtaining  $\mathbb{R}^{K \times D}$  encodings of appearance and geometry can vary depending on the modality being  
 149 handled. Often times modalities can have varying dimensions (images vs. pointclouds) or different  
 150 semantic meaning (surface normals vs. depth). As such, we utilize separate encoders for each modality  
 151 in the Neural USD. In the case of appearance signals in the form of images  $\mathcal{I} \in H \times W \times C$  we  
 152 apply a pre-trained DINOv2 [Caron et al., 2021] model to obtain output features  $\mathcal{F} = \text{Encoder}(\mathcal{I}_i)$ ,  
 153 where  $\mathcal{F} \in h \times w \times D$ . The first two dimensions are then flattened to obtain the resulting embedding  
 154  $\mathcal{A}_i^{\text{emb}} \in \mathbb{R}^{K \times D}$ . Similarly, geometry features such as surface normals and depth can be processed  
 155 using a separate pretrained DINOv2 backbone. We find that normalizing depth features on a per-  
 156 object basis leads to improved generalization performance, as raw metric depth signals constrain  
 157 the object to certain locations in the scene. Preliminary experiments using a shared backbone for  
 158 both image and depth yielded suboptimal results, as the model struggled to accurately represent both  
 159 geometry and appearance.

160 Approaches such as Neural Assets [Wu et al., 2024] first embed an image with a pre-trained backbone  
 161 and slice the resulting feature map using the corresponding 2D bounding box locations for each  
 162 object. While this results in fewer forward passes, it leads to challenges when replacing objects in  
 163 the scene, or conditioning on modalities such as depth or points, since object features are globally  
 164 correlated. Instead, we process each object appearance and depth feature separately, removing global



Figure 4: Object replacement examples with appearance and geometry conditioning (top) and geometry conditioning (bottom).

165 correlations. This can be done efficiently by pre-computing these features and storing them in the  
 166 Neural USD.

### 167 3.3.2 2D and 3D pose encoding

168 Utilizing a separate encoding for 2D and 3D pose signals provides the user access to two different  
 169 interfaces for controlling the position of the object in the scene. The 2D pose allows for simple  
 170 dragging of the object around the scene, while the 3D pose allows for more sophisticated control  
 171 of properties such as distance from the camera and rotation. The 2D bounding box is defined as  
 172 the image-normalized coordinates of the top left corner of the bounding box, as well as the image-  
 173 normalized height and width  $p_i^{2D} = (x_i, y_i, h_i, w_i)$ . We represent the 3D bounding box by projecting  
 174 four corners that span the bounding box to the image plane, arriving at  $\{p_i^{3D,j} = (h_i^j, w_i^j, d_i^j)\}_{j=1}^4$ ,  
 175 with projected image-normalized 2D coordinates  $(h_i^j, w_i^j)$  and 3D depth  $d_i^j$ . We project the 2D  
 176 bounding box and the concatenated corners of the 3D bounding box with a simple multi-layer  
 177 perceptron to obtain:

$$\mathcal{P}_i^{\text{emb}, 2D} = \text{MLP}(p_i^{2D}), \quad (1)$$

$$\mathcal{P}_i^{\text{emb}, 3D} = \text{MLP}(p_i^{3D}), \quad (2)$$

$$p_i^{3D} = \text{Concat}[p_i^{3D,1}, p_i^{3D,2}, p_i^{3D,3}, p_i^{3D,4}]. \quad (3)$$

### 178 3.4 Combining encodings

179 In this section, we describe how we combine the defined encodings so that they can be used to  
 180 condition downstream models via cross-attention [Vaswani et al., 2017], FiLM layers [Perez et al.,  
 181 2017], or in place of text embeddings. To do so, we simply concatenate the appearance, geometry,  
 182 and pose tokens channel-wise.

$$\tilde{a}_i = \text{Concat}[\mathcal{A}_i^{\text{emb}}, \mathcal{G}_i^{\text{emb}}, \mathcal{P}_i^{\text{emb}, 3D}, \mathcal{P}_i^{\text{emb}, 2D}], \quad (4)$$

$$\tilde{a}_i \in \mathbb{R}^{K \times M}. \quad (5)$$

183 where the Neural USD asset encoding  $\tilde{a}_i$  is projected via an MLP to obtain:

$$a_i = \text{MLP}(\tilde{a}_i), a_i \in \mathbb{R}^{K \times D}. \quad (6)$$

184 Finally all Neural USD asset encodings are concatenated along the token dimension to obtain the  
 185 final Neural USD encoding:

$$\mathcal{N} = \text{Concat}[\tilde{a}_1, \dots, \tilde{a}_N], \mathcal{N} \in \mathbb{R}^{(N \times K) \times D}. \quad (7)$$

186 Although approaches such as Neural Assets [Wu et al., 2024] require the background to be encoded  
 187 separately, the flexible structure of the Neural USD allows the background to simply be defined as  
 188 another asset, with its corresponding appearance and geometry signals. Masking foreground objects  
 189 in the provided background signals led to improved results. We provide an additional pose embedding  
 190 during training to represent the source-to-target transform. This helps the model learn not only the  
 191 movement of the objects, but also that of the background scene (i.e. camera movement).

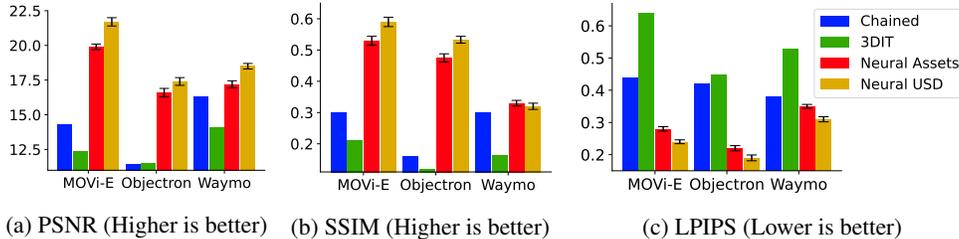


Figure 5: Object control performance on MOVi-E, Objectron, and Waymo. Values measure quality of target reconstruction for single and multi-object scenes.

### 192 3.5 Fine-tuning models with Neural USD

193 Neural USD makes few assumptions about the downstream  
 194 image model to be fine-tuned with the Neural USD. Given  
 195 that a Neural USD is simply a sequence of tokens, it is  
 196 amenable to conditioning via cross-attention or FiLM layers;  
 197 techniques supported by nearly all architectures currently  
 198 used in generative modeling. In this work, we use Stable  
 199 Diffusion v2.1 [Rombach et al., 2022], an exemplary open  
 200 source generative model which is widely used. Given the relatively  
 201 poor performance of Stable Diffusion v2.1 compared to SOTA  
 202 models, we do not expect SOTA-level prediction quality  
 203 and instead demonstrate new conditioning capabilities that  
 204 can be applied to image models. Both the encoders and the  
 205 image model are fine-tuned end to end using the training  
 206 objective defined in the following section. During training,  
 207 we randomly zero out tokens for the entire asset, for modalities  
 208 of an asset, or for 2D and 3D pose signals. This helps  
 209 individual modality features to be invariant of other modality  
 210 features, allowing for precise control. Modality dropout  
 211 also allows the use of Classifier Free Guidance [Ho and Salimans,  
 212 2022] during evaluation.

### 213 3.6 Learning from pairs of images

214 The naïve approach of using individual images with Neural USD  
 215 annotations leads to the entanglement of conditioning signals,  
 216 whereby the model only uses the appearance and geometry  
 217 encodings and entirely disregards the pose encodings, thereby  
 218 limiting pose control of objects in the scene. The use of  
 219 video sequences yields a promising solution to this challenge.  
 220 Video sequences offer multiple views of objects in the scene,  
 221 granting us access to a variety of sources information when  
 222 constructing our Neural USD encoding. Specifically, we extract  
 223 appearance, geometry, and other spatial conditioning signals  
 224 from a source image  $I_{src}$  by cropping out these elements with  
 225 the corresponding 2D bounding-box annotations. The 2D and 3D  
 226 target poses are referenced from a target image  $I_{tgt}$ . The  
 Neural USD encoding is composed using the source spatial  
 modalities and target poses and provided to the image model.  
 The training objective is to reconstruct  $I_{tgt}$  using the  
 denoising diffusion loss of Stable Diffusion v2.1 [Rombach  
 et al., 2022]. This training recipe encourages the model to  
 learn appearance and geometry encodings that are not  
 correlated with pose encodings. The resulting model can be  
 controlled via multiple independent control signals.

## 227 4 Experiments

228 Through our experiments, we try to answer the following  
 229 questions: 1) Does Neural USD allow for precise object  
 230 pose control? 2) Does Neural USD allow for object  
 appearance and geometry control? 3) How does Neural USD  
 compare to other generative model conditioning approaches?

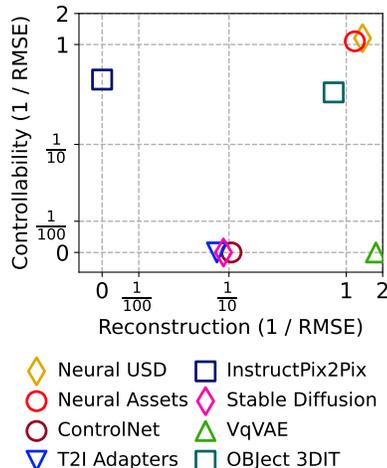


Figure 6: Reconstruction vs. Controllability performance on MOVi-E, Objectron, and Waymo. Axes are log-scale and values are reported as  $1/\text{RMSE}$ .

## 231 4.1 Experiments and datasets

232 We utilized four datasets containing video sequences of scenes containing various objects with 2D  
233 and 3D bounding box annotations. *MOVi-E* [Greff et al., 2022] is a synthetic generated using Blender  
234 scenes with up to 23 objects and consists of object and camera movement. *Objectron* [Ahmadyan  
235 et al., 2020] increases visual complexity by capturing single- and multi-object real world scenes. It  
236 consists of 15,000 videos of objects that cover nine categories. *Waymo Open* [Sun et al., 2020] is  
237 a dataset of real-world self-driving cars captured by a car-mounted camera from multiple angles.  
238 Finally, *EgoTracks* [Tang et al., 2023] is an annotated version of *Ego4D* [Grauman et al., 2022]  
239 consisting of 22.5k object tracks derived from 5.9k videos. The dataset contains a vast number of  
240 objects, many of which are only seen once, and challenging egocentric movement. Unlike the other  
241 datasets, *EgoTracks* only contains 2D bounding boxes. We filter out objects with small bounding  
242 boxes and random flip images. We list additional dataset information in Appendix A.

## 243 4.2 Baselines

244 We compare our work to methods that can a) perform 3D-aware object-centric editing of images,  
245 or b) allow for modification of images using other conditioning signals such as text, geometry, or  
246 appearance for image editing. 3D-aware object-centric editing approaches include *Neural Assets* [Wu  
247 et al., 2024], *Object 3DIT* [Michel et al., 2023], and *Chained* [Wu et al., 2024]. *Object 3DIT* is  
248 limited in its ability to render large viewpoint changes as it does not encode camera poses, while  
249 *Neural Assets* only supports 3D bounding box and RGB appearance conditioning. We also study  
250 non 3D-aware baselines such as *InstructPix2Pix* [Brooks et al., 2022], which uses text descriptions  
251 to modify a source image, *ControlNet* [Zhang et al., 2023], which supports various control signals  
252 during image generation, *T2I-Adapter* [Mou et al., 2023], which learns various adapters to support  
253 additional spatial control signals, and *Stable Diffusion v2.1* [Rombach et al., 2022], a large pre-trained  
254 text-to-image model. Additionally, we include a VqVAE baseline, consisting of a convolutional  
255 encoder, finite scalar quantization [Mentzer et al., 2023], and a UViT [Hooigeboom et al., 2023]  
256 decoder, and trained with a multi-scale denoising diffusion loss [Hooigeboom et al., 2023]. A more  
257 comprehensive discussion of baselines is included in Appendix B.

## 258 4.3 Evaluation criteria

259 We evaluated neural USD and baselines using two different criteria. The first simply measures the  
260 model’s ability to correctly reconstruct the target image from the provided source encodings and  
261 target pose. We use SSIM [Wang et al., 2004], LPIPS [Zhang et al., 2018], and PSNR. We use these  
262 criteria to compare the model to other 3D-aware image editing approaches.

263 We also introduce a novel metric to determine the completeness of a neural scene descriptor. Inspired  
264 by the computer graphics community, we argue that a useful descriptor of a scene must perform well  
265 along two axes; First, the descriptor must provide a full description of the scene’s underlying state,  
266 such that the scene can be reliably reconstructed. Second, the descriptor must expose control signals  
267 that enable a user to reliably convert a source scene to a target scene. In the neural setting, we measure  
268 these two axes as follows: To determine the model’s reconstruction performance, we supply it with all  
269 available conditioning signals extracted from a source image and measure the similarity between the  
270 source image and the model prediction. Controllability is measured by providing a control input that  
271 describes the difference between the source scene and the target scene, and measuring the similarity  
272 between the model prediction and the target image. *InstructPix2Pix* is not capable of performing pure  
273 image reconstruction and can only modify source images given a text prompt. In contrast, *ControlNet*,  
274 the *T2I-Adapter*, *Stable Diffusion*, and the VqVAE are reconstruction-only methods, as they do not  
275 support incremental editing of input control signals.

276 Across all experiment images are resized to  $256 \times 256$ . We use DINO ViT-B/8 [Caron et al., 2021]  
277 as the pre-trained appearance and geometry encoder, and *Stable Diffusion v2.1* [Rombach et al.,  
278 2022] as the pre-trained image generation model. Both the conditioning encoders and the image  
279 models are individually fine-tuned end-to-end using the Adam optimizer [Kingma and Ba, 2017].  
280 The experiments were carried out on 128 TPU v6e chips. Additional hyperparameters can be found  
281 in Appendix C.

## 282 4.4 Results

### 283 4.4.1 Reconstruction quality

284 In Figure 5 we compare Neural USD, Neural Assets, Object 3DIT, and Chained. All methods require  
285 roughly equal amounts of compute for inference (2x2 v6e TPUs) and produce outputs with similar  
286 speeds (10-50ms). We extract source and target frames from the dataset which contain multi-object  
287 changes and compare the model’s ability to accurately reconstruct the target image. We find that  
288 Neural USD outperforms Object 3DIT, Chained, InstructPix2Pix and improves over Neural Assets  
289 while introducing a more flexible conditioning format with additional control inputs. Figure 7 in  
290 Appendix D displays qualitative examples. We notice that Object 3DIT is unable to accurately place  
291 the objects in the target scene, and struggles especially with camera perspective changes. Figure 2  
292 showcases additional Neural USD qualitative examples<sup>1</sup>.

### 293 4.4.2 Scene representation completeness

294 Figure 6 shows the performance of various models across the axes of controllability and reconstruction.  
295 ControlNet, Stable Diffusion, T2I adapters, and VqVAE only expose reconstruction interfaces, and  
296 don’t allow for object-centric or global edits of an input image. As such, measuring their controllability  
297 is challenging, since proposing modifications to input conditioning signals like depth, edge maps, or  
298 masks is non-trivial. Alternatively, InstructPix2Pix does provide an interface for image-level edits via  
299 text, but does not allow for reconstruction of an image from input conditioning signals. Approaches  
300 that allow for both include Object 3DIT, Neural Assets, and Neural USD, which can serve to either  
301 fully reconstruct a desired image, or modify a source image using object-centric controls.

## 302 4.5 Object-centric image editing

303 Neural USD exposes various ways for the user to interact with the image model that were previously  
304 unavailable. In Figure 2, we demonstrate how Neural USD can let the user translate, rotate, and  
305 scale objects as desired. Additionally, users can choose to condition solely on geometry, which leads  
306 to novel appearances that satisfy the 3D structure of the original object. Neural USD also exposes  
307 the ability to replace objects with other desired objects, or to replace the background in an image.  
308 Figure 4 showcases more examples of object replacement, in which the user replaces the original  
309 object in the image with the appearance or geometry of a new image obtained from the Internet.  
310 Geometry information can be easily acquired from arbitrary images by annotating them with depth  
311 annotation models such as ZoeDepth [Bhat et al., 2023a]. Additional editing examples can be found  
312 in Appendix E.

## 313 5 Conclusion

314 Neural USD introduces an object-centric conditioning framework for generative models, enabling  
315 precise control over appearance, geometry, and pose. Inspired by the Universal Scene Descriptor  
316 (USD), it encodes structured per-object attributes into conditioning vectors, ensuring cross-model  
317 compatibility. Using a fine-tuning approach with paired video frames, Neural USD disentangles  
318 control signals for independent object manipulation. Experiments show superior performance in  
319 structured scene synthesis and object control, establishing Neural USD as a flexible and portable  
320 standard for generative modeling.

321 A limitation of the current approach is the reliance on 3D bounding-box annotations. Given the  
322 limited amount of data with 3D bounding box annotations, we find that introducing novel objects  
323 sometimes fails (Appendix Figure 11). We expect this problem to be resolved by co-training on 2D  
324 bounding box datasets or with large-scale datasets with 3D annotations, which may soon be within  
325 reach given recent advancements in 3D bounding box prediction [Krishnan et al., 2024].

---

<sup>1</sup>Additional visuals can be found on the project website: [www.neural-usd.com](http://www.neural-usd.com).

## 326 References

- 327 Adel Ahmadyan, Liangkai Zhang, Jianing Wei, Artsiom Ablavatski, and Matthias Grundmann.  
328 Objectron: A large scale dataset of object-centric videos in the wild with pose annotations, 2020.  
329 URL <https://arxiv.org/abs/2012.09988>.
- 330 Hadi Alzayer, Zhihao Xia, Xuaner Zhang, Eli Shechtman, Jia-Bin Huang, and Michael Gharbi. Magic  
331 fixup: Streamlining photo editing by watching dynamic videos. *ArXiv*, abs/2403.13044, 2024.  
332 URL <https://api.semanticscholar.org/CorpusID:268537350>.
- 333 Omri Avrahami, Thomas Hayes, Oran Gafni, Sonal Gupta, Yaniv Taigman, Devi Parikh, Dani  
334 Lischinski, Ohad Fried, and Xiaoyue Yin. Spatext: Spatio-textual representation for controllable  
335 image generation. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition*  
336 *(CVPR)*, pages 18370–18380, 2022. URL <https://api.semanticscholar.org/CorpusID:254018089>.  
337
- 338 Shariq Farooq Bhat, Reiner Birkel, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth:  
339 Zero-shot transfer by combining relative and metric depth, 2023a. URL <https://arxiv.org/abs/2302.12288>.  
340
- 341 Shariq Farooq Bhat, Niloy Jyoti Mitra, and Peter Wonka. Loosecontrol: Lifting controlnet for general-  
342 ized depth conditioning. *ArXiv*, abs/2312.03079, 2023b. URL <https://api.semanticscholar.org/CorpusID:265693942>.  
343
- 344 Blender Online Community. *Blender - a 3D modelling and rendering package*. Blender Foundation,  
345 Stichting Blender Foundation, Amsterdam, 2018. URL <http://www.blender.org>.
- 346 Tim Brooks, Aleksander Holynski, and Alexei A. Efros. Instructpix2pix: Learning to follow image  
347 editing instructions. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition*  
348 *(CVPR)*, pages 18392–18402, 2022. URL <https://api.semanticscholar.org/CorpusID:253581213>.  
349
- 350 Ming Cao, Xintao Wang, Zhongang Qi, Ying Shan, Xiaohu Qie, and Yinqiang Zheng. Masac-  
351 trl: Tuning-free mutual self-attention control for consistent image synthesis and editing. *2023*  
352 *IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 22503–22513, 2023. URL  
353 <https://api.semanticscholar.org/CorpusID:258179432>.
- 354 Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and  
355 Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of*  
356 *the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9650–9660, October  
357 2021.
- 358 Eric Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic  
359 implicit generative adversarial networks for 3d-aware image synthesis. *2021 IEEE/CVF Conference*  
360 *on Computer Vision and Pattern Recognition (CVPR)*, pages 5795–5805, 2020. URL <https://api.semanticscholar.org/CorpusID:227247980>.  
361
- 362 Eric Chan, Connor Z. Lin, Matthew Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio  
363 Gallo, Leonidas J. Guibas, Jonathan Tremblay, S. Khamis, Tero Karras, and Gordon Wetzstein.  
364 Efficient geometry-aware 3d generative adversarial networks. *2022 IEEE/CVF Conference on*  
365 *Computer Vision and Pattern Recognition (CVPR)*, pages 16102–16112, 2021. URL <https://api.semanticscholar.org/CorpusID:245144673>.  
366
- 367 Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite:  
368 Attention-based semantic guidance for text-to-image diffusion models. *ACM Transactions on*  
369 *Graphics (TOG)*, 42:1 – 10, 2023. URL <https://api.semanticscholar.org/CorpusID:256416326>.  
370
- 371 Minghao Chen, Iro Laina, and Andrea Vedaldi. Training-free layout control with cross-attention  
372 guidance. *2024 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages  
373 5331–5341, 2023. URL <https://api.semanticscholar.org/CorpusID:258041377>.

- 374 Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, Yuwei Fang, Kwot Sin Lee, Ivan Sko-  
375 rokhodov, Kfir Aberman, Jun-Yan Zhu, Ming-Hsuan Yang, and Sergey Tulyakov. Multi-subject  
376 open-set personalization in video generation. *arXiv preprint arXiv:2501.06187*, 2025.
- 377 Wenzheng Chen, Jun Gao, Huan Ling, Edward James Smith, Jaakko Lehtinen, Alec Jacobson, and  
378 Sanja Fidler. Learning to predict 3d objects with an interpolation-based differentiable renderer.  
379 In *Neural Information Processing Systems*, 2019. URL [https://api.semanticscholar.org/  
380 CorpusID:199442423](https://api.semanticscholar.org/CorpusID:199442423).
- 381 Wenzheng Chen, Joey Litalien, Jun Gao, Zian Wang, Clément Fuji Tsang, S. Khamis, Or Litany, and  
382 Sanja Fidler. Dib-r++: Learning to predict lighting and material with a hybrid differentiable ren-  
383 derer. In *Neural Information Processing Systems*, 2021. URL [https://api.semanticscholar.  
384 org/CorpusID:240353816](https://api.semanticscholar.org/CorpusID:240353816).
- 385 Carl Doersch, Yi Yang, Mel Vecerik, Dilara Gokay, Ankush Gupta, Yusuf Aytar, Joao Carreira,  
386 and Andrew Zisserman. Tapir: Tracking any point with per-frame initialization and temporal  
387 refinement, 2023. URL <https://arxiv.org/abs/2306.08637>.
- 388 Laura Downs, Anthony Francis, Nate Koenig, Brandon Kinman, Ryan Hickman, Krista Reymann,  
389 Thomas B. McHugh, and Vincent Vanhoucke. Google scanned objects: A high-quality dataset of  
390 3d scanned household items, 2022. URL <https://arxiv.org/abs/2204.11918>.
- 391 Dave Epstein, A. Jabri, Ben Poole, Alexei A. Efros, and Aleksander Holynski. Diffusion self-  
392 guidance for controllable image generation. *ArXiv*, abs/2306.00986, 2023. URL [https://api.  
393 semanticscholar.org/CorpusID:258999106](https://api.semanticscholar.org/CorpusID:258999106).
- 394 Weixi Feng, Xuehai He, Tsu-Jui Fu, Varun Jampani, Arjun Reddy Akula, P. Narayana, Sugato  
395 Basu, Xin Eric Wang, and William Yang Wang. Training-free structured diffusion guidance  
396 for compositional text-to-image synthesis. *ArXiv*, abs/2212.05032, 2022. URL [https://api.  
397 semanticscholar.org/CorpusID:254535649](https://api.semanticscholar.org/CorpusID:254535649).
- 398 Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-  
399 scene: Scene-based text-to-image generation with human priors. *ArXiv*, abs/2203.13131, 2022.  
400 URL <https://api.semanticscholar.org/CorpusID:247628171>.
- 401 Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel  
402 Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inver-  
403 sion. *ArXiv*, abs/2208.01618, 2022. URL [https://api.semanticscholar.org/CorpusID:  
404 251253049](https://api.semanticscholar.org/CorpusID:251253049).
- 405 Jun Gao, Tianchang Shen, Zian Wang, Wenzheng Chen, K. Yin, Daiqing Li, Or Litany, Zan Gojic,  
406 and Sanja Fidler. Get3d: A generative model of high quality 3d textured shapes learned from  
407 images. *ArXiv*, abs/2209.11163, 2022. URL [https://api.semanticscholar.org/CorpusID:  
408 252438648](https://api.semanticscholar.org/CorpusID:252438648).
- 409 Daniel Garibi, Shahar Yadin, Roni Paiss, Omer Tov, Shiran Zada, Ariel Ephrat, Tomer Michaeli, Inbar  
410 Mosseri, and Tali Dekel. Tokenverse: Versatile multi-concept personalization in token modulation  
411 space. *arXiv preprint arXiv:2501.12224*, 2025.
- 412 Songwei Ge, Taesung Park, Jun-Yan Zhu, and Jia-Bin Huang. Expressive text-to-image generation  
413 with rich text. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages  
414 7511–7522, 2023. URL <https://api.semanticscholar.org/CorpusID:258108187>.
- 415 Vedit Goel, Elia Peruzzo, Yifan Jiang, Dejjia Xu, Niculae Sebe, Trevor Darrell, Zhangyang Wang, and  
416 Humphrey Shi. Pair-diffusion: Object-level image editing with structure-and-appearance paired  
417 diffusion models. *ArXiv*, abs/2303.17546, 2023. URL [https://api.semanticscholar.org/  
418 CorpusID:257834185](https://api.semanticscholar.org/CorpusID:257834185).
- 419 Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Gird-  
420 har, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, Miguel Martin, Tushar Nagarajan,  
421 Ilija Radosavovic, Santhosh Kumar Ramakrishnan, Fiona Ryan, Jayant Sharma, Michael Wray,  
422 Mengmeng Xu, Eric Zhongcong Xu, Chen Zhao, Siddhant Bansal, Dhruv Batra, Vincent Car-  
423 tillier, Sean Crane, Tien Do, Morrie Doulaty, Akshay Erapalli, Christoph Feichtenhofer, Adriano

- 424 Fragomeni, Qichen Fu, Abrham Gebreselasie, Cristina Gonzalez, James Hillis, Xuhua Huang,  
425 Yifei Huang, Wenqi Jia, Weslie Khoo, Jachym Kolar, Satwik Kottur, Anurag Kumar, Federico  
426 Landini, Chao Li, Yanghao Li, Zhenqiang Li, Kartikeya Mangalam, Raghava Modhugu, Jonathan  
427 Munro, Tullie Murrell, Takumi Nishiyasu, Will Price, Paola Ruiz Puentes, Merey Ramazanova,  
428 Leda Sari, Kiran Somasundaram, Audrey Southerland, Yusuke Sugano, Ruijie Tao, Minh Vo,  
429 Yuchen Wang, Xindi Wu, Takuma Yagi, Ziwei Zhao, Yunyi Zhu, Pablo Arbelaez, David Crandall,  
430 Dima Damen, Giovanni Maria Farinella, Christian Fuegen, Bernard Ghanem, Vamsi Krishna  
431 Ithapu, C. V. Jawahar, Hanbyul Joo, Kris Kitani, Haizhou Li, Richard Newcombe, Aude Oliva,  
432 Hyun Soo Park, James M. Rehg, Yoichi Sato, Jianbo Shi, Mike Zheng Shou, Antonio Torralba,  
433 Lorenzo Torresani, Mingfei Yan, and Jitendra Malik. Ego4d: Around the world in 3,000 hours of  
434 egocentric video, 2022. URL <https://arxiv.org/abs/2110.07058>.
- 435 Klaus Greff, Francois Belletti, Lucas Beyer, Carl Doersch, Yilun Du, Daniel Duckworth, David J.  
436 Fleet, Dan Gnanapragasam, Florian Golemo, Charles Herrmann, Thomas Kipf, Abhijit Kundu,  
437 Dmitry Lagun, Issam Laradji, Hsueh-Ti, Liu, Henning Meyer, Yishu Miao, Derek Nowrouzezahrai,  
438 Cengiz Oztireli, Etienne Pot, Noha Radwan, Daniel Rebain, Sara Sabour, Mehdi S. M. Sajjadi,  
439 Matan Sela, Vincent Sitzmann, Austin Stone, Deqing Sun, Suhani Vora, Ziyu Wang, Tianhao Wu,  
440 Kwang Moo Yi, Fangcheng Zhong, and Andrea Tagliasacchi. Kubric: A scalable dataset generator,  
441 2022. URL <https://arxiv.org/abs/2203.03570>.
- 442 Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. Stylenerf: A style-based 3d-aware  
443 generator for high-resolution image synthesis. *ArXiv*, abs/2110.08985, 2021. URL <https://api.semanticscholar.org/CorpusID:239016913>.
- 444 Amir Hertz, Ron Mokady, Jay M. Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or.  
445 Prompt-to-prompt image editing with cross attention control. *ArXiv*, abs/2208.01626, 2022. URL  
446 <https://api.semanticscholar.org/CorpusID:251252882>.
- 447 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance, 2022. URL <https://arxiv.org/abs/2207.12598>.
- 448 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
449 *neural information processing systems*, 33:6840–6851, 2020.
- 450 Emiel Hooeboom, Jonathan Heek, and Tim Salimans. Simple diffusion: End-to-end diffusion for  
451 high resolution images, 2023. URL <https://arxiv.org/abs/2301.11093>.
- 452 Liucheng Hu, Xin Gao, Peng Zhang, Ke Sun, Bang Zhang, and Liefeng Bo. Animate anyone:  
453 Consistent and controllable image-to-video synthesis for character animation. *2024 IEEE/CVF*  
454 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8153–8163, 2023. URL  
455 <https://api.semanticscholar.org/CorpusID:265499043>.
- 456 A. Jabri, Sjoerd van Steenkiste, Emiel Hooeboom, Mehdi S. M. Sajjadi, and Thomas Kipf. Dorsal:  
457 Diffusion for object-centric representations of scenes et al. *ArXiv*, abs/2306.08068, 2023. URL  
458 <https://api.semanticscholar.org/CorpusID:259190298>.
- 459 Xuhui Jia, Yang Zhao, Kelvin C. K. Chan, Yandong Li, Han-Ying Zhang, Boqing Gong, Tingbo Hou,  
460 H. Wang, and Yu-Chuan Su. Taming encoder for zero fine-tuning image customization with text-to-  
461 image diffusion models. *ArXiv*, abs/2304.02642, 2023. URL <https://api.semanticscholar.org/CorpusID:257952647>.
- 462 Jindong Jiang, Fei Deng, Gautam Singh, and Sung Tae Ahn. Object-centric slot diffusion. *ArXiv*,  
463 abs/2303.10834, 2023. URL <https://api.semanticscholar.org/CorpusID:257632090>.
- 464 Yash Kant, Ziyi Wu, Michael Vasilkovsky, Guocheng Qian, Jian Ren, Riza Alp Guler, Bernard  
465 Ghanem, S. Tulyakov, Igor Gilitschenski, and Aliaksandr Siarohin. Spad: Spatially aware  
466 multi-view diffusers. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition*  
467 *(CVPR)*, pages 10026–10038, 2024. URL <https://api.semanticscholar.org/CorpusID:267547881>.
- 468 Bahjat Kavar, Shiran Zada, Oran Lang, Omer Tov, Hui-Tang Chang, Tali Dekel, Inbar Mosseri,  
469 and Michal Irani. Imagic: Text-based real image editing with diffusion models. *2023 IEEE/CVF*  
470 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6007–6017, 2022. URL  
471 <https://api.semanticscholar.org/CorpusID:252918469>.

- 476 Yunji Kim, Jiyoung Lee, Jin-Hwa Kim, Jung-Woo Ha, and Jun-Yan Zhu. Dense text-to-image  
477 generation with attention modulation. *2023 IEEE/CVF International Conference on Computer Vi-*  
478 *sion (ICCV)*, pages 7667–7677, 2023. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:261101003)  
479 261101003.
- 480 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL  
481 <https://arxiv.org/abs/1412.6980>.
- 482 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete  
483 Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick.  
484 Segment anything, 2023.
- 485 Akshay Krishnan, Abhijit Kundu, Kevis-Kokitsi Maninis, James Hays, and Matthew Brown.  
486 Omninos: A unified nocs dataset and model for 3d lifting of 2d objects, 2024. URL  
487 <https://arxiv.org/abs/2407.08711>.
- 488 Dongxu Li, Junnan Li, and Steven C. H. Hoi. Blip-diffusion: Pre-trained subject representation  
489 for controllable text-to-image generation and editing. *ArXiv*, abs/2305.14720, 2023a. URL  
490 <https://api.semanticscholar.org/CorpusID:258865473>.
- 491 Siyuan Li, Lei Ke, Martin Danelljan, Luigi Piccinelli, Mattia Segu, Luc Van Gool, and Fisher  
492 Yu. Matching anything by segmenting anything, 2024. URL [https://arxiv.org/abs/2406.](https://arxiv.org/abs/2406.04221)  
493 04221.
- 494 Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan  
495 Li, and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. *2023 IEEE/CVF*  
496 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22511–22521, 2023b.  
497 URL <https://api.semanticscholar.org/CorpusID:255942528>.
- 498 Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten  
499 Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content  
500 creation. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages  
501 300–309, 2022. URL <https://api.semanticscholar.org/CorpusID:253708074>.
- 502 Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick.  
503 Zero-1-to-3: Zero-shot one image to 3d object. *2023 IEEE/CVF International Conference on*  
504 *Computer Vision (ICCV)*, pages 9264–9275, 2023a. URL [https://api.semanticscholar.](https://api.semanticscholar.org/CorpusID:257631738)  
505 [org/CorpusID:257631738](https://api.semanticscholar.org/CorpusID:257631738).
- 506 Yuan Liu, Chu-Hsing Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping  
507 Wang. Syncdreamer: Generating multiview-consistent images from a single-view image. *ArXiv*,  
508 abs/2309.03453, 2023b. URL <https://api.semanticscholar.org/CorpusID:261582503>.
- 509 Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold,  
510 Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot atten-  
511 tion. *ArXiv*, abs/2006.15055, 2020. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:220127924)  
512 220127924.
- 513 Grace Luo, Trevor Darrell, Oliver Wang, Dan B Goldman, and Aleksander Holynski. Read-  
514 out guidance: Learning control from diffusion features. *2024 IEEE/CVF Conference on*  
515 *Computer Vision and Pattern Recognition (CVPR)*, pages 8217–8227, 2023. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:265608773)  
516 265608773.
- 517 Jiancang Ma, Junhao Liang, Chen Chen, and H. Lu. Subject-diffusion: Open domain personalized  
518 text-to-image generation without test-time fine-tuning. *ArXiv*, abs/2307.11410, 2023. URL  
519 <https://api.semanticscholar.org/CorpusID:260091569>.
- 520 Luke Melas-Kyriazi, Iro Laina, C. Rupprecht, and Andrea Vedaldi. Realfusion 360° reconstruction  
521 of any object from a single image. *2023 IEEE/CVF Conference on Computer Vision and Pattern*  
522 *Recognition (CVPR)*, pages 8446–8455, 2023. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:261092262)  
523 [CorpusID:261092262](https://api.semanticscholar.org/CorpusID:261092262).
- 524 Fabian Mentzer, David Minnen, Eirikur Agustsson, and Michael Tschannen. Finite scalar quantization:  
525 Vq-vae made simple, 2023. URL <https://arxiv.org/abs/2309.15505>.

- 526 Oscar Michel, Anand Bhattad, Eli VanderBilt, Ranjay Krishna, Aniruddha Kembhavi, and Tanmay  
527 Gupta. Object 3dit: Language-guided 3d-aware image editing. *ArXiv*, abs/2307.11073, 2023. URL  
528 <https://api.semanticscholar.org/CorpusID:259991631>.
- 529 Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, Ying Shan, and  
530 Xiaohu Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image  
531 diffusion models, 2023. URL <https://arxiv.org/abs/2302.08453>.
- 532 Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative  
533 neural feature fields. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition*  
534 *(CVPR)*, pages 11448–11459, 2020. URL <https://api.semanticscholar.org/CorpusID:227151657>.  
535
- 536 Karran Pandey, Paul Guerrero, Matheus Gadelha, Yannick Hold-Geoffroy, Karan Singh, and Niloy J.  
537 Mitra. Diffusion handles enabling 3d edits for diffusion models by lifting activations to 3d. *2024*  
538 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7695–7704,  
539 2023. URL <https://api.semanticscholar.org/CorpusID:265659119>.
- 540 Gaurav Parmar, Or Patashnik, Kuan-Chieh Wang, Daniil Ostashev, Srinivasa Narasimhan, Jun-Yan  
541 Zhu, Daniel Cohen-Or, and Kfir Aberman. Object-level visual prompts for compositional image  
542 generation. *arXiv preprint arXiv:2501.01424*, 2025.
- 543 Dario Pavllo, Graham Spinks, Thomas Hofmann, Marie-Francine Moens, and Aurélien Lucchi.  
544 Convolutional generation of textured 3d meshes. *ArXiv*, abs/2006.07660, 2020. URL <https://api.semanticscholar.org/CorpusID:219687111>.  
545
- 546 Dario Pavllo, Jonas Köhler, Thomas Hofmann, and Aurélien Lucchi. Learning generative models  
547 of textured 3d meshes from real-world images. *2021 IEEE/CVF International Conference on*  
548 *Computer Vision (ICCV)*, pages 13859–13869, 2021. URL <https://api.semanticscholar.org/CorpusID:232404704>.  
549
- 550 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*  
551 *the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023.
- 552 Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual rea-  
553 soning with a general conditioning layer, 2017. URL <https://arxiv.org/abs/1709.07871>.
- 554 Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d dif-  
555 fusion. *ArXiv*, abs/2209.14988, 2022. URL <https://api.semanticscholar.org/CorpusID:252596091>.  
556
- 557 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
558 resolution image synthesis with latent diffusion models, 2022. URL <https://arxiv.org/abs/2112.10752>.  
559
- 560 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.  
561 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. *2023*  
562 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22500–22510,  
563 2022. URL <https://api.semanticscholar.org/CorpusID:251800180>.
- 564 Mehdi S. M. Sajjadi, Daniel Duckworth, Aravindh Mahendran, Sjoerd van Steenkiste, Filip Paveti’c,  
565 Mario Luvci’c, Leonidas J. Guibas, Klaus Greff, and Thomas Kipf. Object scene representa-  
566 tion transformer. *ArXiv*, abs/2206.06922, 2022. URL <https://api.semanticscholar.org/CorpusID:249642130>.  
567
- 568 Kyle Sargent, Zizhang Li, Tanmay Shah, Charles Herrmann, Hong-Xing Yu, Yunzhi Zhang, Eric Ryan  
569 Chan, Dmitry Lagun, Fei-Fei Li, Deqing Sun, and Jiajun Wu. Zeronvs: Zero-shot 360-degree  
570 view synthesis from a single real image. *ArXiv*, abs/2310.17994, 2023. URL <https://api.semanticscholar.org/CorpusID:264555531>.  
571
- 572 Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance  
573 fields for 3d-aware image synthesis. *ArXiv*, abs/2007.02442, 2020. URL <https://api.semanticscholar.org/CorpusID:220364071>.  
574

575 Maximilian Seitzer, Sjoerd van Steenkiste, Thomas Kipf, Klaus Greff, and Mehdi S. M. Sajjadi. Dyst:  
576 Towards dynamic neural scene representations on real-world videos. *ArXiv*, abs/2310.06020, 2023.  
577 URL <https://api.semanticscholar.org/CorpusID:263829437>.

578 Jing Shi, Wei Xiong, Zhe Lin, and Hyun Joon Jung. Instantbooth: Personalized text-to-image  
579 generation without test-time finetuning. *2024 IEEE/CVF Conference on Computer Vision and*  
580 *Pattern Recognition (CVPR)*, pages 8543–8552, 2023a. URL <https://api.semanticscholar.org/CorpusID:258041269>.

582 Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and X. Yang. Mvdream: Multi-view diffu-  
583 sion for 3d generation. *ArXiv*, abs/2308.16512, 2023b. URL <https://api.semanticscholar.org/CorpusID:261395233>.

585 Yujun Shi, Chuhui Xue, Jiachun Pan, Wenqing Zhang, Vincent Y. F. Tan, and Song Bai. Dragdiffusion:  
586 Harnessing diffusion models for interactive point-based image editing. *2024 IEEE/CVF Conference*  
587 *on Computer Vision and Pattern Recognition (CVPR)*, pages 8839–8849, 2023c. URL <https://api.semanticscholar.org/CorpusID:259252555>.

589 Gautam Singh, Fei Deng, and Sungjin Ahn. Illiterate dall-e learns to compose. *ArXiv*, abs/2110.11405,  
590 2021. URL <https://api.semanticscholar.org/CorpusID:239616181>.

591 Gautam Singh, Yeongbin Kim, and Sungjin Ahn. Neural systematic binder. In *International*  
592 *Conference on Learning Representations*, 2022a. URL <https://api.semanticscholar.org/CorpusID:255749563>.

594 Gautam Singh, Yi-Fu Wu, and Sungjin Ahn. Simple unsupervised object-centric learning for complex  
595 and naturalistic videos. *ArXiv*, abs/2205.14065, 2022b. URL <https://api.semanticscholar.org/CorpusID:249151816>.

597 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised  
598 learning using nonequilibrium thermodynamics. In *International conference on machine learning*,  
599 pages 2256–2265. PMLR, 2015.

600 Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui,  
601 James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, Jiquan Ngiam,  
602 Hang Zhao, Aleksei Timofeev, Scott Ettinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Sheng  
603 Zhao, Shuyang Cheng, Yu Zhang, Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov.  
604 Scalability in perception for autonomous driving: Waymo open dataset, 2020. URL <https://arxiv.org/abs/1912.04838>.

606 Hao Tang, Kevin Liang, Matt Feiszli, and Weiyao Wang. Egotracks: A long-term egocentric visual  
607 object tracking dataset, 2023. URL <https://arxiv.org/abs/2301.03213>.

608 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
609 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von  
610 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, edi-  
611 tors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates,  
612 Inc., 2017. URL [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).

614 Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A. Yeh, and Gregory Shakhnarovich. Score  
615 jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. *2023 IEEE/CVF*  
616 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12619–12629, 2022. URL  
617 <https://api.semanticscholar.org/CorpusID:254125253>.

618 Jiawei Wang, Yuchen Zhang, Jiaxin Zou, Yan Zeng, Guoqiang Wei, Liping Yuan, and Hang Li.  
619 Boximator: Generating rich and controllable motions for video synthesis. *ArXiv*, abs/2402.01566,  
620 2024a. URL <https://api.semanticscholar.org/CorpusID:267406297>.

621 Qixun Wang, Xu Bai, Haofan Wang, Zekui Qin, and Anthony Chen. Instantid: Zero-shot  
622 identity-preserving generation in seconds. *ArXiv*, abs/2401.07519, 2024b. URL <https://api.semanticscholar.org/CorpusID:266999462>.

- 624 Xudong Wang, Trevor Darrell, Sai Saketh Rambhatla, Rohit Girdhar, and Ishan Misra. In-  
625 stancediffusion: Instance-level control for image generation. *2024 IEEE/CVF Conference on*  
626 *Computer Vision and Pattern Recognition (CVPR)*, pages 6232–6242, 2024c. URL [https://](https://api.semanticscholar.org/CorpusID:267412534)  
627 [api.semanticscholar.org/CorpusID:267412534](https://api.semanticscholar.org/CorpusID:267412534).
- 628 Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. Image quality assessment: from error  
629 visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.  
630 doi: 10.1109/TIP.2003.819861.
- 631 Daniel Watson, William Chan, Ricardo Martin-Brualla, Jonathan Ho, Andrea Tagliasacchi, and  
632 Mohammad Norouzi. Novel view synthesis with diffusion models. *ArXiv*, abs/2210.04628, 2022.  
633 URL <https://api.semanticscholar.org/CorpusID:252780361>.
- 634 Yuxiang Wei, Yabo Zhang, Zhilong Ji, Jinfeng Bai, Lei Zhang, and Wangmeng Zuo. Elite: Encoding  
635 visual concepts into textual embeddings for customized text-to-image generation. *2023 IEEE/CVF*  
636 *International Conference on Computer Vision (ICCV)*, pages 15897–15907, 2023. URL [https://](https://api.semanticscholar.org/CorpusID:257219968)  
637 [api.semanticscholar.org/CorpusID:257219968](https://api.semanticscholar.org/CorpusID:257219968).
- 638 Ziyi Wu, Jingyu Hu, Wuyue Lu, Igor Gilitschenski, and Animesh Garg. Slotdiffusion: Object-  
639 centric generative modeling with diffusion models. *ArXiv*, abs/2305.11281, 2023. URL [https://](https://api.semanticscholar.org/CorpusID:258822805)  
640 [api.semanticscholar.org/CorpusID:258822805](https://api.semanticscholar.org/CorpusID:258822805).
- 641 Ziyi Wu, Yulia Rubanova, Rishabh Kabra, Drew A. Hudson, Igor Gilitschenski, Yusuf Aytar, Sjoerd  
642 van Steenkiste, Kelsey R. Allen, and Thomas Kipf. Neural assets: 3d-aware multi-object scene  
643 synthesis with image diffusion models, 2024. URL <https://arxiv.org/abs/2406.09292>.
- 644 Guangxuan Xiao, Tianwei Yin, William T. Freeman, Frédo Durand, and Song Han. Fastcomposer:  
645 Tuning-free multi-subject image generation with localized attention. *ArXiv*, abs/2305.10431, 2023.  
646 URL <https://api.semanticscholar.org/CorpusID:258740710>.
- 647 Jinheng Xie, Yuexiang Li, Yawen Huang, Haozhe Liu, Wentian Zhang, Yefeng Zheng, and  
648 Mike Zheng Shou. Boxdiff: Text-to-image synthesis with training-free box-constrained dif-  
649 fusion. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 7418–7427,  
650 2023. URL <https://api.semanticscholar.org/CorpusID:259991581>.
- 651 Yinghao Xu, Menglei Chai, Zifan Shi, Sida Peng, Ivan Skorokhodov, Aliaksandr Siarohin, Ceyuan  
652 Yang, Yujun Shen, Hsin-Ying Lee, Bolei Zhou, and S. Tulyakov. Discoscene: Spatially disentangled  
653 generative radiance fields for controllable 3d-aware scene synthesis. *2023 IEEE/CVF Conference*  
654 *on Computer Vision and Pattern Recognition (CVPR)*, pages 4402–4412, 2022. URL [https://](https://api.semanticscholar.org/CorpusID:254974555)  
655 [api.semanticscholar.org/CorpusID:254974555](https://api.semanticscholar.org/CorpusID:254974555).
- 656 Zhongcong Xu, Jianfeng Zhang, Jun Hao Liew, Hanshu Yan, Jia-Wei Liu, Chenxu Zhang, Jiashi  
657 Feng, and Mike Zheng Shou. Magicanimate: Temporally consistent human image animation using  
658 diffusion model. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,  
659 pages 1481–1490, 2023. URL <https://api.semanticscholar.org/CorpusID:265466012>.
- 660 Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang  
661 Zhao. Depth anything v2, 2024. URL <https://arxiv.org/abs/2406.09414>.
- 662 Xuan Yang, Liangzhe Yuan, Kimberly Wilber, Astuti Sharma, Xiuye Gu, Siyuan Qiao, Stephanie  
663 Debats, Huisheng Wang, Hartwig Adam, Mikhail Sirotenko, and Liang-Chieh Chen. Polymax:  
664 General dense prediction with mask transformer, 2023. URL [https://arxiv.org/abs/2311.](https://arxiv.org/abs/2311.05770)  
665 [05770](https://arxiv.org/abs/2311.05770).
- 666 Zhengyuan Yang, Jianfeng Wang, Zhe Gan, Linjie Li, Kevin Lin, Chenfei Wu, Nan Duan, Zicheng Liu,  
667 Ce Liu, Michael Zeng, and Lijuan Wang. Reco: Region-controlled text-to-image generation. *2023*  
668 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14246–14255,  
669 2022. URL <https://api.semanticscholar.org/CorpusID:254043880>.
- 670 Jiraphon Yenphraphai, Xichen Pan, Sainan Liu, Daniele Panozzo, and Saining Xie. Image  
671 sculpting: Precise object editing with 3d geometry control. *2024 IEEE/CVF Conference on*  
672 *Computer Vision and Pattern Recognition (CVPR)*, pages 4241–4251, 2024. URL [https://](https://api.semanticscholar.org/CorpusID:266741835)  
673 [api.semanticscholar.org/CorpusID:266741835](https://api.semanticscholar.org/CorpusID:266741835).

- 674 Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan,  
675 Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-  
676 rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2(3):5, 2022.
- 677 Junyi Zhang, Charles Herrmann, Junhwa Hur, Varun Jampani, Trevor Darrell, Forrester Cole, Deqing  
678 Sun, and Ming-Hsuan Yang. Monst3r: A simple approach for estimating geometry in the presence  
679 of motion, 2024. URL <https://arxiv.org/abs/2410.03825>.
- 680 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image  
681 diffusion models. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages  
682 3813–3824, 2023. URL <https://api.semanticscholar.org/CorpusID:256827727>.
- 683 Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable  
684 effectiveness of deep features as a perceptual metric, 2018. URL <https://arxiv.org/abs/1801.03924>.  
685

## 686 **A Datasets**

### 687 **A.1 EgoTracks**

688 EgoTracks [Tang et al., 2023] is a tracked bounding box dataset consisting of manually labeled 22.5k  
689 object tracks spanning 5.9k videos. Being a derivative dataset of Ego4D [Grauman et al., 2022],  
690 EgoTracks is ego-centric and features an extreme amount of foreground and background movement.  
691 Additionally, Ego4D object tracks often feature very small bounding boxes (e.g. for utensils). We  
692 filter the data in two ways: first, we filter out bounding boxes with small height or width, the threshold  
693 being 1/10th the normalized height and width of the screen. To prevent object tracks from leaving the  
694 field of view, we sample source and target frames from within 15 frames of each other, and discard  
695 samples for which no object is present. Finally, during evaluation, we filter out results with motion  
696 blur or extreme background shift.

### 697 **A.2 Objectron**

698 Objectron [Ahmadyan et al., 2020] consists of 15,000 object-centric video clips featuring everyday  
699 objects across nine categories. Each video includes object pose tracking, allowing us to extract  
700 3D bounding boxes. Since the dataset lacks 2D bounding box annotations, we generate them by  
701 projecting the eight corners of the 3D boxes onto the image and computing the tightest bounding box  
702 around the projected points.

### 703 **A.3 Waymo Open**

704 Waymo Open [Sun et al., 2020] comprises 1,000 video clips of self-driving scenes captured by  
705 car-mounted cameras. Following previous studies [Wu et al., 2024], we use the front-view camera  
706 and car bounding box annotations. The 3D bounding boxes include only the heading angle (yaw-axis  
707 rotation), so we set the other two rotation angles to zero. Additionally, the provided 2D and 3D boxes  
708 are misaligned, making paired frame training unfeasible. To address this, we project the 3D boxes to  
709 obtain corresponding 2D boxes, similar to the approach used for Objectron.

### 710 **A.4 MOVi-E**

711 MOVi-E [Greff et al., 2022] includes 10,000 videos simulated using Kubric [Greff et al., 2022],  
712 with each scene featuring 11 to 23 real-world objects from the Google Scanned Objects (GSO)  
713 repository [Downs et al., 2022]. At the beginning of each video, multiple objects are dropped onto the  
714 ground, causing them to collide. The scene’s lighting comes from a randomly sampled environment  
715 map. The camera follows a simple linear motion.

## 716 **B Baselines**

### 717 **B.1 Object 3DIT**

718 Object 3DIT [Michel et al., 2023] fine-tunes Zero-1-to-3 [Liu et al., 2023a] for scene-level 3D object  
719 editing. We derive editing instructions from the target object pose, including translation and rotation.  
720 However, this lacks support for significant viewpoint changes as it does not encode camera poses.  
721 We use the official code and pre-trained weights of the Multitask variant.

### 722 **B.2 InstructPix2Pix**

723 InstructPix2Pix enables text-guided image editing by fine-tuning a diffusion model to follow editing  
724 instructions. It conditions on both an input image and a text prompt, learning to predict pixel changes  
725 based on the instruction. However, it lacks explicit 3D control and struggles with complex multi-  
726 object edits. We construct a dataset of 100 source target pairs and their differences describes as text  
727 prompts. InstructPix2Pix is then conditioned on the source image and text prompt, and we evaluate  
728 how accurate it is at reconstructing the target image. We find that the model struggles to elicit the fine  
729 changes in object pose described in the text.

730 **B.3 T2I Adapters**

731 T2I-Adapters [Mou et al., 2023] enable additional conditioning mechanisms for pre-trained diffusion  
732 models, allowing control beyond text prompts. They integrate spatial signals like depth maps or  
733 segmentation masks to guide image generation while preserving the original model’s structure.  
734 These adapters typically introduce lightweight modules, such as attention layers or zero-initialized  
735 convolutions, that fuse external control signals with the model’s latent space. We condition the  
736 T2I-adapters model on the spatial modalities corresponding to the source image, such as masks and  
737 depth, and evaluate its performance in reconstructing the source image.

738 **B.4 Neural Assets**

739 Neural Assets [Wu et al., 2024] introduces a per-object representation for 3D-aware multi-object  
740 control in image diffusion models. It encodes appearance and pose separately, allowing object  
741 manipulation, including translation, rotation, and rescaling. We evaluate Neural Assets using the  
742 same criteria used for Neural USD: we extract source modalities from a source image and condition  
743 on the target poses derived from the target image. We then measure the reconstruction loss with the  
744 target image.

745 **B.5 ControlNet**

746 ControlNet [Zhang et al., 2023] enables spatial conditioning in diffusion models by introducing  
747 trainable layers that process external control signals, such as edge maps, depth maps, or pose  
748 keypoints. It retains the original model’s weights while adding zero-initialized convolution layers.  
749 This allows for control over image generation. However, it does not allow for object-centric image  
750 editing, as changes to the conditioning signals can lead to global changes in the image. We condition  
751 the Control model on the spatial modalities corresponding to the source image, such as masks and  
752 depth, and evaluate its performance in reconstructing the source image.

753 **C Hyperparameters**

754 Here we outline the hyperparameters used to implement and train Neural USD.

Table 1: Hyperparameters for Neural USD.

| PARAMETER                            | VALUE                |
|--------------------------------------|----------------------|
| STABLE DIFFUSION VARIANT             | v2.1                 |
| DINO VARIANT                         | ViT-B/8              |
| DINO FEATURE MAP SIZE                | 28 × 28              |
| INPUT IMAGE SIZE                     | 256 × 256            |
| TOKEN DIMENSION                      | 1024                 |
| BATCH SIZE                           | 512                  |
| OPTIMIZER                            | ADAM                 |
| STABLE DIFFUSION LR                  | 1 × 10 <sup>-4</sup> |
| IMAGE ENCODER (DINO) LR              | 5 × 10 <sup>-4</sup> |
| WARMUP STEPS                         | 2000                 |
| DECAY SCHEDULE                       | LINEAR               |
| FINE-TUNING STEPS                    | 50000                |
| GRADIENT CLIP VALUE                  | 1.0                  |
| MODALITY DROPOUT PROBABILITY         | 0.25                 |
| POSE DROPOUT PROBABILITY             | 0.25                 |
| ALL CONDITIONING DROPOUT PROBABILITY | 0.1                  |
| CFG WEIGHT                           | 3.0                  |

755 **D Qualitative baselines**



Figure 7: Object pose conditioning performance on MOVi-E, Objectron, Waymo Open, and EgoTracks. Models generate the target image provided a source image and the 3D bounding box targets (Neural USD, 3DIT) or textual prompts (InstructPix2Pix). Our method satisfies the desired pose while preserving the foreground and background appearance. InstructPix2Pix fails to elicit object movement.

756 **E Additional Experimental Results**

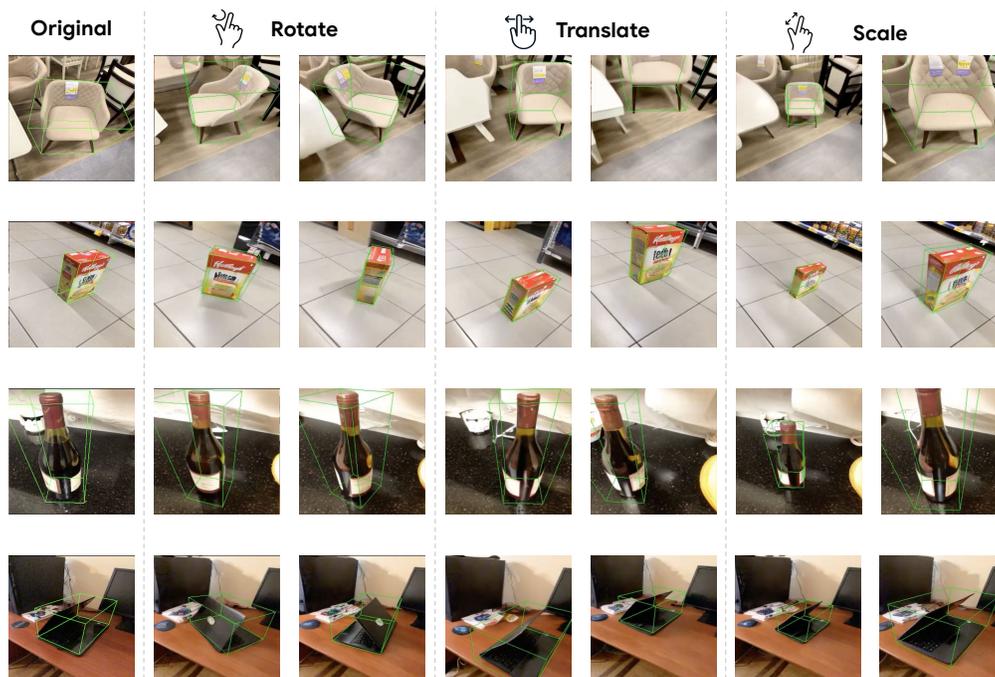


Figure 8: Additional Objectron 3D pose control examples.

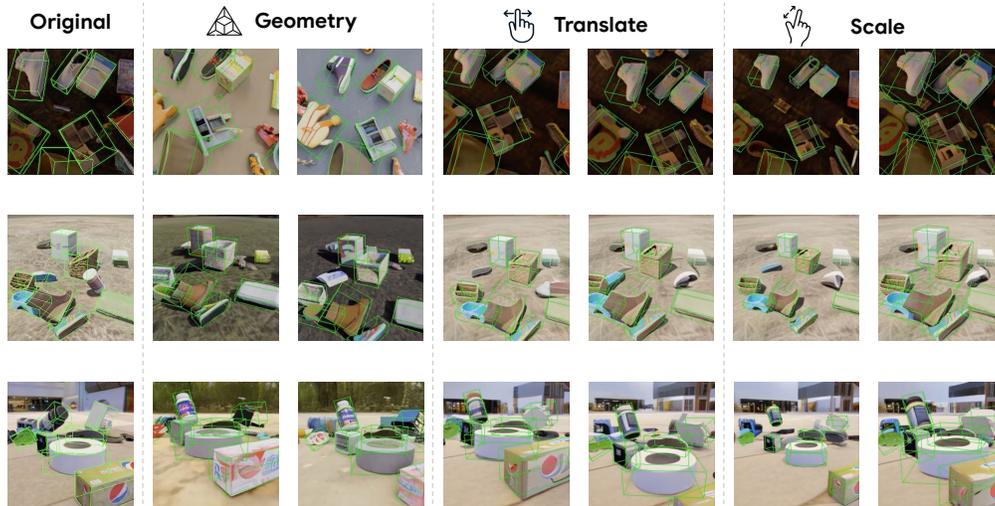


Figure 9: Additional MOVi-E 3D pose control examples.

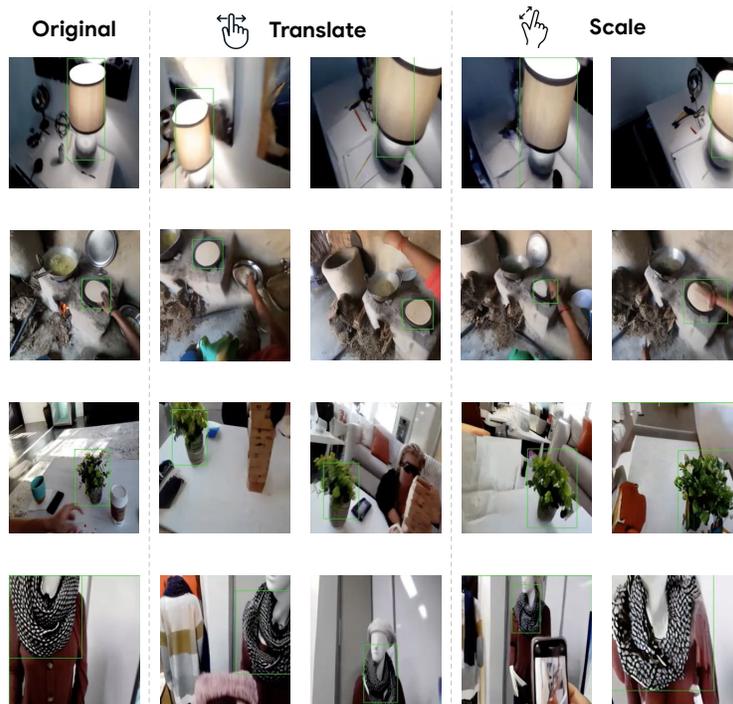


Figure 10: Additional Egotracks 2D bounding box control examples.



Figure 11: When trained on limited data - a handful of labeled datasets constituting  $<50,000$  sequences - Neural USD fails to generalize to new object categories. This lack of generalization can likely be remedied by co-training on more readily-available 2D bounding box datasets.

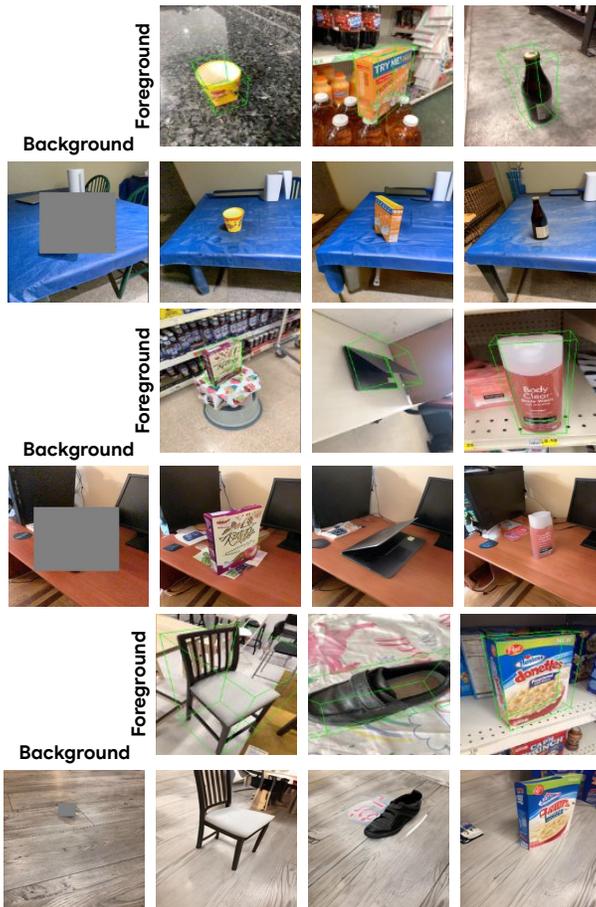


Figure 12: Foreground background replacement. Neural USD allows for easy swapping of assets in the scene. The background, simply being another asset, can be replaced with reference modalities.