VideoHandles: Editing 3D Object Compositions in Videos Using Video **Generative Priors**

Anonymous CVPR submission

Paper ID 14057



Figure 1. VideoHandles edits 3D object composition in videos of static scenes. Solid axes represent the original 3D position and dotted axes the user-provided target position. The edit plausibly updates effects like the reflection of the wine glass and handles disocclusions like the lamp behind the book pile that is exposed by the edit. In addition to generated videos, we can also edit real (non-generated) videos by inverting the video into its corresponding latent, as shown on the right.

Abstract

001 Generative methods for image and video editing use generative models as priors to perform edits despite incomplete 002 information, such as changing the composition of 3D ob-003 jects shown in a single image. Recent methods have shown 004 promising composition editing results in the image setting, 005 006 but in the video setting, editing methods have focused on editing object's appearance and motion, or camera motion, 007 and as a result, methods to edit object composition in videos 008 are still missing. We propose VideoHandles as a method for 009 010 editing 3D object compositions in videos of static scenes 011 with camera motion. Our approach allows editing the 3D 012 position of a 3D object across all frames of a video in a temporally consistent manner. This is achieved by lifting inter-013 014 mediate features of a generative model to a 3D reconstruction that is shared between all frames, editing the recon-015 016 struction, and projecting the features on the edited reconstruction back to each frame. To the best of our knowledge, 017 this is the first generative approach to edit object composi-018 tions in videos. Our approach is simple and training-free, while outperforming state-of-the-art image editing baselines. 021

1. Introduction

Diffusion models and flow-based models are currently the 023 standard for high-quality text-to-image generation. Text-024 to-video diffusion/flow-based models lag behind in quality, 025 but have recently seen big improvements. The prevalent 026 text-based control is easy to use, but impractical for some 027 types of edits, such as edits of the object composition in a 028 scene: specifying the position of an object with text is in-029 accurate and iterative editing workflows are not supported. 030 Several recent methods address this issue in the image do-031 main by proposing different types of iterative image editing 032

019 020

022

090

091

092

093

094

095

096

097

098

099

033 methods. These either focus on editing the appearance of 034 objects [4, 13, 46], or their spatial composition [1, 3, 28]. 035 In the video domain, current methods support editing only the appearance [6, 22] while lacking methods to edit spa-036 037 tial object compositions, for example, editing the 3D position of objects in generated videos, as shown in Figure 1. 038 Editing the object composition in a video introduces several 039 challenges: a *plausible* editing output requires generating 040 041 details such as shadows and lighting that may have changed due to the edited composition; furthermore, the edited video 042 043 needs to preserve the identity of the original objects and should adhere to an edit control manipulated by the user. 044 Finally, the edit needs to be applied to all video frames in a 045 temporally consistent manner. 046

We propose VideoHandles as a generative approach to 047 edit the object composition in a video of a static scene. Our 048 approach allows editing the 3D position of a 3D object in 049 a video, resulting in a plausible, temporally consistent edit 050 051 that preserves the identity of the original object. To the best of our knowledge, ours is the first generative approach that 052 allows editing the object composition in a video. Given a 053 pretrained flow-based video generative model, we present a 054 novel method to edit the intermediate features from the gen-055 056 erative model's network in a temporally consistent manner. 057 Specifically, we lift the intermediate features of each frame to a common 3D reconstruction, effectively treating them 058 059 as latent textures. We then edit the 3D location of an object using 3D translations or rotations, and project the features 060 061 back to their corresponding frames. We use such projected 062 features as guidance during the generative process to create a plausible edited video. Editing of real (non-generated) 063 videos is supported by first inverting them into the random 064 noise. Our approach is simple and does not require any 065 066 training or finetuning that risks biasing the distribution of the generative model. 067

We evaluate our method on several generated and cap-068 tured videos. As there are no existing methods that are 069 specialized to editing the 3D object composition in videos, 070 071 we compare to several image editing baselines that can be 072 applied in a per-frame manner. We evaluate the results in terms of plausibility, temporal consistency, identity preser-073 074 vation, and adherence to the target edit. In addition to a 075 large number of qualitative comparisons, we also conduct a 076 user study. The users have a clear preference for our method in terms of plausibility and temporal consistency, while our 077 method is at least on par or slightly better than image editing 078 baselines in terms of identity preservation and edit adher-079 080 ence. Finally, we perform a quantitative evaluation which confirms these findings. 081

We summarize our contributions as follows:

082

We introduce a zero-shot method for editing object composition in videos using video generative priors, for the first time to our knowledge.

- For the feature-based generative editing process, we describe the optimal approach for feature extraction from the video generative model network (Section 4.2).
- We also demonstrate that self-attention-map-based weighting (Section 4.4) and null-text prediction in the foreground region (Section 4.5) further improve the edit-ing quality.
- We demonstrate the effectiveness of our method with both generated and real videos.

2. Related Work

In the context of diffusion/flow-based generative models, several methods have been proposed for image and video editing that can be roughly grouped by the type of edits they perform.

Image Appearance Editing. There has been a series of 100 work that focus on manipulating intermediate features or at-101 tention maps of pre-trained image diffusion models to edit 102 the appearance of objects within an image [4, 5, 10, 13, 41, 103 46] in a zero-shot setting. While effective, such methods 104 often do not focus on editing the composition of objects, 105 which requires control over object positions and strict iden-106 tity preservation. To tackle identity preservation, various 107 customization approaches have been proposed that enable 108 the generation of images of a particular object or subject 109 in different compositions. However, such methods do not 110 provide edit controls and typically require finetuning of the 111 base model [18, 32]. The prior of image diffusion mod-112 els has been further utilized to enable editing of 3D static 113 scenes represented as 3D neural assets via iterative opti-114 mization approaches [16, 17, 29]. These methods, however, 115 also focus on changing the appearance of objects, rather 116 than our goal of composition editing. 117

Image Composition Editing. Several recent methods 118 aim at editing the composition of objects in an image [1-119 3, 7-9, 25, 28, 47, 49]. Another line of work aims at in-120 serting an object from a source image into a new target im-121 age [38, 39, 45], which can be repurposed as image editing 122 tools by using the same image as source and target. An-123 other popular editing workflow provides control points that 124 can be dragged by a user to deform objects or edit 2D object 125 positions [21, 27, 34–36]. Additionally, a few more general 126 image editing methods have been proposed that can be used 127 for either image appearance editing or image composition 128 editing [24, 48]. All of these methods can be applied to 129 videos by separately editing each frame, but this loses tem-130 poral consistency, as we show in our experiments in Sec-131 tion 5. Most related to our work is Diffusion Handles [28] 132 which inspired our approach of editing intermediate fea-133 tures using a 3D reconstruction. We show how to modify 134 this approach so it can be applied to non-depth-conditioned 135 video priors, including which features to pick, which 3D 136

201

202

203

204

205

reconstruction method to use, how to avoid artifacts fromhard object masks, and how to effectively remove the origi-

139 nal object from the edited video.

Video Appearance Editing. With the increasing quality 140 of video generators, various works have focused on editing 141 the appearance of objects within videos. A key issue these 142 works aim to tackle is to maintain temporal consistency be-143 tween frames while changing the appearance. To address 144 145 this, some works [6, 7, 31] have proposed techniques to maintain consistency using only image diffusion models, 146 147 while others [14, 15] have leveraged priors from video diffusion models to tackle this challenge. While successful in 148 149 preserving temporal consistency, these approaches are lim-150 ited to appearance changes, and no prior work addresses changing object compositions in videos. 151

152 Video Motion Control. Recently, another line of work in the video domain focuses on controlling the motion of ob-153 jects or cameras during generation [19, 33, 42, 44]. Al-154 though these methods allow specifying how a particular ob-155 ject should move in a video, they are designed for gen-156 eration rather than editing tasks, thus they do not allow 157 modifying the compositions of static object arrangements 158 in videos. Moreover, unlike these approaches that require 159 training on task-specific datasets to learn motion control, 160 VideoHandles is training-free. 161

162 3. Preliminary: Flow-Based Latent Video163 Model

In this section, we briefly discuss the video prior we usein our experiments, which is the flow-based latent videomodel, OpenSora [51].

167 Flow-Based Generative Model. Similar to diffusion models [12, 37], flow-based generative models [20, 168 23] model high-dimensional data distributions through a 169 learned iterative process. Given a data sample $Z_1 \sim p_{data}$ 170 and random noise $oldsymbol{Z}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$ a linear trajectory is de-171 172 fined as $Z_t = tZ_1 + (1-t)Z_0$. Based on the linear trajectory, a veloicty prediction network v_{θ} is trained to estimate 173 the derivative dZ_t/dt : 174

175
$$v_{\theta}(\boldsymbol{Z}_t, t, y) \approx \frac{d}{dt} \boldsymbol{Z}_t = \boldsymbol{Z}_1 - \boldsymbol{Z}_0, \qquad (1)$$

176 where y encodes the text prompt corresponding to Z_1 . 177 Given a trained velocity prediction network v_{θ} , a new data 178 sample can be generated through the generative process, 179 starting from Z_0 :

$$\mathbf{Z}_{t+\Delta t} = \mathbf{Z}_t + \Delta t \cdot v_{\theta}^{\omega}(\mathbf{Z}_t, t, y), \qquad (2)$$

181 where $v_{\theta}^{\omega}(\mathbf{Z}_t, t, y) = v_{\theta}(\mathbf{Z}_t, t, \emptyset) + \omega(v_{\theta}(\mathbf{Z}_t, t, y) - v_{\theta}(\mathbf{Z}_t, t, \emptyset))$ denotes a prediction using classifier-free

guidance [11] with null-text embedding \varnothing and guidance183scale ω . The step size Δt can be chosen at inference time184to balance quality with speed.185

DiT-Based Architecture for Latent Video Model. A 186 video $\boldsymbol{X} \in \mathbb{R}^{n imes h imes w imes 3}$ with n frames is encoded into a 187 latent representation $\boldsymbol{Z}_1 \in \mathbb{R}^{M imes H imes W imes D}$ by a pre-trained 188 encoder, where all dimensions except the feature dimension 189 D are reduced. Each pixel of the latent representation en-190 codes a spatio-temporal patch of X. The velocity prediction 191 network v_{θ} is implemented as a DiT [30] that operates on 192 this latent representation, with alternating blocks of spatial 193 self-attention, temporal self-attention, and cross-attention to 194 the text prompt. A total of 24 blocks of each type are used. 195 A latent sampled from the generative process is decoded by 196 a pre-trained decoder to produce a video sample. 197

4. VideoHandles: A 3D-Aware Video Editing Method 198

Consider a static input video $X_{src} \in \mathbb{R}^{n \times h \times w \times 3}$, where objects remain stationary and only the camera moves. Our goal is to apply a 3D transformation to an object selected by the user in the first frame while preserving the identity of the input video, realism, and temporal consistency. See Figure 2 for an architecture overview.

To ensure that transformations in each frame of a video 206 align with those in other frames, we define a 3D space in 207 which a point cloud $P_{src} = {\mathbf{p}^{(j)}}_{j=1}^{J}$ represents the 3D 208 scene in the video with a shared coordinate system across 209 all frames. A transformation is performed in this shared 210 3D space, denoted by \mathcal{T} : $\mathbb{R}^3 \to \mathbb{R}^3$, with each input 211 frame $\mathbf{x}_{\text{src}}^{(i)}$ modeled as a 2D rendering of $\boldsymbol{P}_{\text{src}}$ from the *i*-th 212 view. Specifically, we reconstruct $\boldsymbol{P}_{\mathrm{src}}$ and estimate a cam-213 era pose for each frame from $X_{\rm src}$ using DUST3R [43]. By 214 leveraging the reconstructed 3D scene from $X_{
m src}$, we define 215 a 3D-aware warping function in the 2D space of each frame. 216

However, due to inaccuracies in warping caused by er-217 rors in reconstructing the 3D scene, directly warping pixel 218 colors often leads to unrealistic videos. Moreover, this ap-219 proach fails to appropriately adjust the video according to 220 the 3D scene and the transformed object, such as new shad-221 ows, reflections, and relighting effects. Therefore, inspired 222 by Diffusion Handles [28], we perform warping in the fea-223 ture space of a pre-trained video generative model and use 224 the warped features as guidance during the generative pro-225 cess. This ensures that the generative prior of the video 226 model adapts the scene with appropriate context changes 227 according to the new object composition while maintaining 228 temporal consistency. 229

In the following sections, we first introduce how to compute the warping function for each 2D frame based on the transformation of an object in the 3D scene (Section 4.1). 232

257



Figure 2. VideoHandles Architecture. We use the intermediate features $\Psi_{\rm src}$ of a video generative model to represent the identity of objects in a source video. Given a 3D transformation of an object, we can use a 3D reconstruction of the scene to warp the intermediate features consistently across frames. Guiding the video generator with these warped features $\Psi_{\rm tgt}$ gives us a an edited video where the object is transformed, while also maintaining the plausibility of effects like shadows and reflections.

Next, we describe the features of the pretrained flow-based
latent video model and how these features are warped (Section 4.2). Lastly, we explain how the warped video model
features serve as guidance in the energy-based guided generative process (Section 4.3).

238 4.1. 3D-Aware Warping Function

239 We first describe how to obtain a 3D-aware warping function in the 2D space of each frame. Given a set of 2D co-240 ordinates $\Omega_{H,W} = \{(v, u) \mid v \in [0, H), u \in [0, W)\}$, the 241 connection between the 3D space and the *i*-th 2D frame is 242 established through the projection function $f^{(i)}$: $\mathbb{R}^3 \rightarrow$ 243 244 $\Omega_{H,W}$, which is defined by the *i*-th camera pose. Let $\mathcal{B}_{
m src}^{(1)}:\Omega_{H,W} o \{0,1\}$ denote the 2D binary mask of an 245 object selected by users in the first frame. Based on the 2D 246 object mask in the first frame $\mathcal{B}_{src}^{(1)}$, we first partition \boldsymbol{P}_{src} , 247 the point cloud reconstructed from the input video $X_{\rm src}$, as 248 249 follows:

250
$$P_f = \{ \mathbf{p} \in \mathbf{P}_{src} \mid \mathcal{B}_{src}^{(1)}(f^{(1)}(\mathbf{p})) = 1 \},\$$

$$\boldsymbol{P}_b = \boldsymbol{P}_{\rm src} \setminus \boldsymbol{P}_f, \tag{4}$$

where P_f consists of points whose projections lie within the 2D masked region defined by $\mathcal{B}_{src}^{(1)}$, and P_b denotes the remaining points representing the background. By applying a 3D transformation \mathcal{T} to P_f alone, we construct a rough target 3D scene represented as a point cloud:

$$\boldsymbol{P}_{\text{tgt}} = \mathcal{T} \boldsymbol{P}_f \cup \boldsymbol{P}_b. \tag{5}$$

The lifting function $g_{\text{src}}^{(i)} : \Omega_{H,W} \to \mathbb{R}^3$ takes a 2D coordinate $\mathbf{u} = (v, u)$ as input and returns the 3D point in P_{src} closest to the *i*-th camera from among the points projected close to \mathbf{u} :

262
$$g_{\text{src}}^{(i)}(\mathbf{u}) = \operatorname*{arg\,min}_{\mathbf{p}\in \boldsymbol{P}_{\text{src},\mathbf{u}}^{(i)}} z^{(i)}(\mathbf{p}), \qquad (6)$$

where $P_{\text{src},\mathbf{u}}^{(i)} = \{\mathbf{p} \in P_{\text{src}} \mid ||f^{(i)}(\mathbf{p}) - \mathbf{u}||_1 < \epsilon\}$ represents the set of 3D points that are projected close to \mathbf{u} and $z^{(i)}(\mathbf{p})$ denotes the distance of point \mathbf{p} from the *i*-th camera. Similarly, $g_{\text{tgt}}^{(i)}(\mathbf{u})$ returns the 3D point in P_{tgt} closest to the *i*-th camera from among the points projected close to \mathbf{u} . 266 the *i*-th camera from among the points projected close to \mathbf{u} . 267 Using the functions $g_{\text{src}}^{(i)}$ and $g_{\text{tgt}}^{(i)}$, we define an occlusionaware foreground point cloud $P_f^{(i)} \subseteq P_f$ for each frame as follows: 270

$$\boldsymbol{P}_{f}^{(i)} = \{g_{\rm src}^{(i)}(\mathbf{u})\} \cap \{\mathcal{T}^{-1}g_{\rm tgt}^{(i)}(\mathbf{u})\} \cap \boldsymbol{P}_{f}, \qquad (7) \qquad \mathbf{271}$$

where $\mathbf{u} \in \Omega_{H \times W}$. It consists of foreground points that are not occluded by the background either before or after the transformation. Using this 3D information, we compute a 2D warping function $\mathcal{W}^{(i)} : \Omega_{H,W} \to \Omega_{H,W}$ as follows: 275

$$\mathcal{W}^{(i)}(\mathbf{u}) = \begin{cases} f^{(i)} \left(\mathcal{T}^{-1} g^{(i)}_{\text{tgt}}(\mathbf{u}) \right), & \text{if } g^{(i)}_{\text{tgt}}(\mathbf{u}) \in \boldsymbol{P}_{f}^{(i)} \\ \mathbf{u}, & \text{otherwise.} \end{cases}$$
(8) 276

This warping function gives us the corresponding coordi-277 nate in the source image for any coordinate in the target im-278 age. All coordinates that do not project to the edited fore-279 ground point cloud remain unchanged. We denote warp-280 ing a 2D signal $\mathcal{X} : \Omega_{H,W} \to \mathbb{R}^{\breve{C}}$ as $(\mathcal{W}^{(i)} * \mathcal{X})(\mathbf{u}) :=$ 281 $\mathcal{X}(\mathcal{W}^{(i)}(\mathbf{u}))$. Similarly, we denote its application to a ten-282 sor $X \in \mathbb{R}^{\dots \times H \times W \times \dots}$ as $\mathcal{W}^{(i)} * X$. Here H and W are the 283 two spatial tensor dimensions that the warping is applied to 284 and the ellipses denote arbitrary additional dimensions. The 285 tensor is sampled at non-integer coordinates using linear in-286 terpolation. 287

As we will show in our evaluation, directly warping RGB frames results in a noisy video, due to inaccuracies in camera predictions and 3D reconstructions, and since this direct warping does not update effects like reflections and shadows that may have changed due to the edit. Therefore, we propose warping the *features* instead of the frames in the 293

(3)

333

339

347

348

368

video and synthesizing the edited video through a generative process that guides the features of the edited video to
match the warped features. In the next section, we introduce
our choice of features for the guided generative process.

4.2. Warping Video Features

In this section, we describe our choice of features ex-299 tracted from OpenSora [51] and explain how these features 300 301 are warped using the warping function introduced in Sec-302 tion 4.1. The DiT architecture [30] of OpenSora alternates 303 layers that perform spatial self-attention, temporal selfattention, cross-attention to the prompt, and feed-forward 304 computations. Spatial attention operates within each frame, 305 while temporal attention is performed among pixels at the 306 same spatial position across frames. We empirically found 307 that the features from the temporal self-attention layers tend 308 to produce global changes; since each temporal attention 309 layer follows a spatial one, its features tend to affect all pix-310 els in each frame globally. This global spatial context is un-311 suitable for our local editing tasks, where only the selected 312 object needs to be transformed. Therefore, we use only ex-313 tract features from the spatial layers for guidance, as these 314 retain more localized information. 315

316 Let $Q_l(Z_t), K_l(Z_t), V_l(Z_t) \in \mathbb{R}^{M \times H \times W \times d}$ be the **317** query, key, and value features of the *l*-th self-attention layer **318** extracted from $v_{\theta}^{\omega}(Z_t, t, y)$, where *M* denotes the number **319** of frames and *d* is the feature dimension. We use their con-**320** catenation from all layers as our extracted feature Ψ :

$$\Psi(\boldsymbol{Z}_t) = [\boldsymbol{Q}_l(\boldsymbol{Z}_t) \parallel \boldsymbol{K}_l(\boldsymbol{Z}_t) \parallel \boldsymbol{V}_l(\boldsymbol{Z}_t)]_{l=1}^L.$$
(9)

322 Let $\Psi^{(i)}(Z_t) \in \mathbb{R}^{H \times W \times D}$ denote the feature for frame *i*, 323 where *D* is the total dimensionality of the feature. Applying 324 the previosuly defined warping function, given the latent of 325 the input video Z_t^{src} , its warped feature is defined as $\Psi_{\text{tgt}}^{(i)} :=$ 326 $\mathcal{W}^{(i)} * \Psi^{(i)}(Z_t^{\text{src}})$.

327 4.3. Warping-Based Guided Generative Process

To guide the generation process of Z_t with $\Psi_{tgt}^{(i)}(Z_t)$, we use an energy-guided generative process [8], similar to classifier-free guidance. Given an energy function $\mathcal{G}(Z_t)$, the gradient of \mathcal{G} is injected at each step of the generative process, steering it towards minimizing the energy function:

$$\boldsymbol{Z}_{t+\Delta t} = \boldsymbol{Z}_t + \Delta t \cdot \boldsymbol{v}_{\theta}^{\omega}(\boldsymbol{Z}_t, t, y) + \rho \nabla_{\boldsymbol{Z}_t} \mathcal{G}(\boldsymbol{Z}_t), \quad (10)$$

where ρ is a hyperparameter to control the step size of $\nabla_{Z_t} \mathcal{G}$. Below, we describe our specific design of \mathcal{G} to edit object compositions in videos.

337 **Object transformation energy.** Let $M_{\text{src}}^{(i)}, M_{\text{tgt}}^{(i)} \in \mathbb{R}^{H \times W}$ denote the occlusion-aware 2D masks of the se-

lected object before and after the transformation:

$$\boldsymbol{M}_{\mathrm{src}}^{(i)}(\mathbf{u}) := \begin{cases} 1, \text{ if } \mathbf{u} \in \{f^{(i)}(\mathbf{p}) \mid \mathbf{p} \in \boldsymbol{P}_{f}^{(i)}\}, \\ 0, \text{ otherwise.} \end{cases}, \quad (11) \qquad 340$$

$$M_{\text{tgt}}^{(i)} := \mathcal{W}^{(i)} * M_{\text{src}}^{(i)},$$
 (12) 341

where $\mathbf{u} \in \Omega_{H \times W}$. Note that $M_{\text{src}}^{(i)}$, which marks the region 342 where the occlusion-aware foreground point cloud $P_f^{(i)}$ is 343 projected, is a subset of the object selection mask $\mathcal{B}_{\text{src}}^{(i)}$ since 344 $M_{\text{src}}^{(i)}$ only includes the object region visible before and after the transformation. 346

To transform the selected object in the video, we define the *object transformation energy* $\mathcal{G}_o(\mathbf{Z}_t)$ as follows:

$$\sum_{i=1}^{M} \left\| \boldsymbol{M}_{tgt}^{(i)} \odot \left(\Psi_{tgt}^{(i)} - \Psi^{(i)}(\boldsymbol{Z}_{t}) \right) \right\|_{2}^{2}, \quad (13) \quad \mathbf{349}$$

where \odot is the element-wise product (broadcasting to additional dimensions where needed). This function measures the discrepancy between the current features $\Psi^{(i)}$ and the target features $\Psi^{(i)}_{tgt}$ within the region of the edited object $M_{tgt}^{(i)}$. 354

Background preservation energy. To further preserve 355 background details, we define an additional energy function called the *background preservation energy* $\mathcal{G}_b(\mathbf{Z}_t)$ as 357 follows: 358

$$\left\|\psi_{\mathrm{MHW}}\left(\boldsymbol{M}_{b}\odot\Psi_{\mathrm{tgt}}\right)-\psi_{\mathrm{MHW}}\left(\boldsymbol{M}_{b}\odot\Psi(\boldsymbol{Z}_{t})\right)\right\|_{2}^{2},\quad(14)$$
359

where ψ_{MHW} denotes the average over time and spatial dimensions, and $M_b^{(i)} = \max((1 - M_{\text{src}}^{(i)} - M_{\text{tgt}}^{(i)}, 0))$ is the background mask. This function measures the discrepancy between the sums of the features in the background region. Unlike \mathcal{G}_o , \mathcal{G}_b compares only the averages of the features, allowing the guidance of \mathcal{G}_b to facilitate appropriate context changes according to the new object position, such as new shadows or reflections.

4.4. Weighted Guidance with Self-Attention Maps

When applying the gradients of the energy function above, 369 the inaccurate 3D reconstruction and camera paths result 370 in guidance sometimes being applied inaccurately to back-371 ground regions, for example at incorrect spatial positions. 372 This sometimes results in hallucinated objects in the back-373 ground regions or other artifacts. To address this, we weight 374 the gradients of the guidance energy using an attention map 375 based on self-attention from the foreground object to other 376 image regions. Intuitively, this includes regions that an edit 377 of the foreground object should affect, including regions 378 that receive updated shadows or reflections, but not regions 379 of the background that should remain unaffected by the edit. 380

391

421

422

423

424

425

426

427



Figure 3. Visualization of our self-attention-based masks. The masks do not only include the the edited object, but also regions requiring semantic adjustments, such as a new reflection under the wine glass and newly disoccluded lamp.

We denote the query and key features of the *i*-the frame, stacked across all spatial self-attention layers and flattened as $Q^{(i)}(Z_t) \in \mathbb{R}^{HW \times 1 \times D}$ and $K^{(i)}(Z_t) \in \mathbb{R}^{1 \times HW \times D}$, both with *H* and *W* are flattened into a single spatial dimension. Then, we define the spatial self-attention map $A^{(i)}(Z_t) \in \mathbb{R}^{HW \times HW}$ for the *i*-the frame as:

$$\boldsymbol{A}^{(i)}(\boldsymbol{Z}_t) \coloneqq \boldsymbol{Q}^{(i)}(\boldsymbol{Z}_t) \ \boldsymbol{K}^{(i)}(\boldsymbol{Z}_t), \tag{15}$$

388 We then find regions that the transformed object pays at-589 tention to by multiplying with the transformed object mask 590 $M_{tgt}^{(i)}$ and normalizing:

$$\Lambda^{(i)} \coloneqq \operatorname{norm}_{[0,1]} \left(\boldsymbol{M}_{\mathsf{tgt}}^{(i)} \boldsymbol{A}^{(i)}(\boldsymbol{Z}_t) \right), \qquad (16)$$

where norm_[0,1] denotes normalization of the value range to [0, 1], $M_{tgt}^{(i)}$ is flattened to $\mathbb{R}^{1 \times HW}$ and the resulting selfattention-based mask $\Lambda^{(i)}$ is unflattened to $\mathbb{R}^{H \times W}$.

The final masked and aggregated self-attention map $\Lambda \in \mathbb{R}^{M \times H \times W}$ is obtained by stacking $\Lambda^{(i)}$ along the temporal dimension. Figure 3 shows that the target self-attention map locally highlights not only the target position of the selected object but also regions requiring adjustments for context changes, such as areas for a new reflection or the disoccluded lamp.

With this mask, a guided step of our generative processis defined as:

404
$$\boldsymbol{Z}_{t+\Delta t} = \boldsymbol{Z}_t + \Delta t \cdot \boldsymbol{v}_{\theta}^{\omega} + \Lambda \odot \nabla_{\boldsymbol{Z}_t} \big(\rho_o \mathcal{G}_o + \rho_b \mathcal{G}_b \big), \quad (17)$$

405 where ρ_o and ρ_b are the step sizes for the gradients of \mathcal{G}_o 406 and \mathcal{G}_b , respectively.

4.5. Null-Text Prediction on Original Object Region 407

When transforming an object, it is undesirable for the ob-408 ject to remain in its original position while being duplicated 409 in the target position. To avoid this object issue, we employ 410 two techniques. First, at the beginning of the generative 411 process of the target, we randomly initialize the original ob-412 ject area of $Z_0^{\rm src}$, as highlighted by the source masks $M_{\rm src}$, 413 and start the generative process from this partially random-414 ized noise. Then, during the generative process, to reduce 415 the influence of text guidance in the original object area and 416 prevent the introduction of a new object in that region, we 417 apply the null-text prediction $v_{\theta}(\mathbf{Z}_t, t, \emptyset)$ within the origi-418 nal object area $M_{\rm src}$ instead of a prediction with classifier-419 free guidance [11]. 420

5. Experiments

Dataset. For quantitative and qualitative comparisons, we generate 27 input videos to be edited, each with a resolution of 320×320 and 51 frames. To enhance the realism of the generated videos, we lightly finetune OpenSora [51] on 71,556 indoor scene videos from the RealEstate10K dataset [52] for 14,000 iterations.

Baselines. In the absence of prior work on modifying 428 3D object composition in videos, we compare our method 429 to Diffusion Handles [28], the state-of-the-art method for 430 composition editing in 2D images, applying the editing pro-431 cess frame by frame. To further demonstrate the effective-432 ness of our feature-guided generative process, we also com-433 pare it to direct frame warping. Specifically, we first re-434 move the selected object from all frames using an existing 435 inpainting technique [40] and then render the transformed 436 foreground point cloud, $\mathcal{T}P_f$, onto the frames where the 437 selected object is removed. Additionally, we introduce an 438 improved version of the direct frame warping, where the 439 video is further refined using SDEdit [24]. SDEdit is per-440 formed for 15 out of the total 30 steps with OpenSora [51]. 441

Qualitative Results. Please refer to the supplementary 442 material for the edited video results. We also present snap-443 shots of the edited videos in Figure 1 and Figure 4. Qual-444 itatively, our method successfully edits object composition 445 in videos while making appropriate contextual adjustments, 446 such as the new reflection beneath the wine glass in Fig-447 ure 1 and the new shadows beneath the transformed car, 448 apple, and vase in rows 2, 3, and 5 of Figure 4, respec-449 tively. In comparison, Diffusion Handles [28] (the fourth 450 column in Figure 4) alters the identity of objects or the back-451 ground across different frames, as seen in the second row, 452 and frequently duplicates objects, as shown in the first row. 453 These failures are more evident in the videos shown in 454 the supplementary material. Direct frame warping (the 455 second column) and its refined one by SDEdit [24] (the third 456

CVPR 2025 Submission #14057. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

481

482

483

484

485

486

487

488



Figure 4. A qualitative comparison with other baselines. The examples show that ours best demonstrates plausibility by avoiding object duplication, adjusting shadows properly, and maintaining consistent outputs across frames, desipte warping errors, as illustrated in the direct frame warping outputs (column 2).

457 column) also typically produce visual seams (second row)
458 and implausible objects (fourth row) due to inaccuracies in
459 warping.

User Study Results. Proper quantitative evaluation for 460 video editing results is very challenging, as there are no es-461 tablished metrics for this task. Therefore, we conducted 462 a user study that included questions about the plausibil-463 ity, identity preservation, and edit coherence of the edited 464 videos. More details about the user study including the 465 queries and setup are provided in the supplementary ma-466 terial. Figure 5 shows human preferences when partici-467 468 pants were presented with two videos-one generated by our method and the other by a competing method-along with 469 the input video, and were asked to choose the better one 470 based on each criterion. The results show that our method 471 is preferred over all baselines across all criteria by signifi-472 cant margins. Notably, our method achieved a preference of 473 474 100% for plausibility compared to Diffusion Handles [28], and 75% and 57% for identity preservation and edit coher-475 ence compared to the SDEdit [24] output of the direct frame 476 warping. 477

478 Temporal Consistency Evaluation. The biggest advan479 tage of our method compared to per-frame-based editing
480 baselines is its ability to achieve temporal consistency. To

Table 1. A quantitative evaluation of Frame LPIPS. Frame LPIPS is scaled by 10^2 , with the best result highlighted in **bold**.

Per-Frame-Based Editing			Ablation Cases			Ours
Direct Warp.	Direct Warp. +SDEdit	Diffusion Handles	w/ Temp. Feature	w/o Self-Attn	w/o Null-Text	Video Handles
5.19	5.03	18.63	3.81	3.77	3.79	3.71

further evaluate this, we introduce a metric called *Frame LPIPS*, which is the average LPIPS [50] score measured between pairs of adjacent frames in the edited video. Frame LPIPS scores for all methods are presented in Table 1. Our method significantly outperforms the baselines, with a score of 3.71 compared to 18.6 for Diffusion Handles [28], demonstrating the superior temporal consistency achieved by leveraging a video prior.

Ablation Study Results. We demonstrate the effective-489 ness of each key aspect of our method through an abla-490 tion study involving three cases: using both spatial and 491 temporal self-attention layer features (w/ Temporal Feature, 492 Section 4.2), omitting self-attention-based weighting in the 493 guided generative process (w/o Self-Attn, Section 4.4), and 494 not using null-text prediction in the original object area (w/o 495 Null-Text, Section 4.5). The user study results in the sec-496 ond row of Figure 5 show that our full method outperforms 497

499

500

501

502

528

CVPR 2025 Submission #14057. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 5. User study results on the plausibility, identity preservation, and edit coherence of the edited videos. Each bar pair shows user preferences, with the green bar for our method and the other for the baseline, along with 95% confidence intervals. We also include a comparison with the input video to represent the upper bound of plausibility.



Figure 6. A qualitative comparison of the ablation study. We show the effect of each component in our method. As demonstrated, our full method avoids object duplication and unnecessary drastic changes in the background, while effectively preserving the identity of the selected object.

all three cases across all metrics by large margins. Moreover, the best temporal consistency is achieved with our full method, as indicated by the lowest Frame LPIPS score compared to the ablation cases, as shown at the bottom of Table 1. Qualitative comparisons are shown in Figure 6.

Please refer to the supplementary material for the edited 503 video results. In the first row, the results without null-text 504 505 (third column) exhibit object duplication, showing the arm-506 chair in both the original and target positions. In the second 507 row, the results without self-attention-based weighting (sec-508 ond column) drastically alter the background colors, and the results without null-text (third column) introduce a new 509 510 knob on the kettle. In contrast, our full method (last column) best preserves the identity of the kettle. In the third 511 row, the results without self-attention-based weighting (sec-512 ond column) and with temporal layer features (fourth col-513 514 umn) generate a new lamp next to the armchair and thus fail to preserve the background. Our method successfully 515 moves the selected armchair without changing the back-516 ground. 517

518 Editing Real Videos with Object Composition. We 519 also showcase the results of editing real videos using our

method, as seen in the rightmost image in Figure 1 and the 520 third row of Figure 6. In these examples, the apple in the 521 former and the armchair in the latter are moved to new posi-522 tions, with shading and shadows generated according to the 523 new composition while successfully preserving the back-524 ground. To edit the real videos, we mapped the videos to 525 their corresponding random latent noises using the null-text 526 inversion technique introduced by Mokady et al. [26]. 527

6. Conclusion

We have presented VideoHandles, the first method to our 529 knowledge that leverages the prior of video generative mod-530 els for editing object composition in video. Given the warp-531 ing function for each frame obtained from a 3D reconstruc-532 tion and transformation of an object in 3D space, Video-533 Handles applies temporally consistent warping to features 534 extracted from a pre-trained video generative model, rather 535 than to the frames themselves, using these features as guid-536 ance in the generative process. Experimental results, in-537 cluding a user study, demonstrate that VideoHandles out-538 performs per-frame editing methods in terms of plausibility, 539 identity preservation, and edit coherence. 540

551

552

553

554

555

556

557

558

559

560

567

568

569

570

571

578

579

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

References 541

- 542 [1] Hadi Alzayer, Zhihao Xia, Xuaner Zhang, Eli Shechtman, 543 Jia-Bin Huang, and Michael Gharbi. Magic fixup: Stream-544 lining photo editing by watching dynamic videos. arXiv 545 preprint arXiv:2403.13044, 2024. 2
- 546 [2] Omri Avrahami, Rinon Gal, Gal Chechik, Ohad Fried, Dani 547 Lischinski, Arash Vahdat, and Weili Nie. Diffuhaul: A 548 training-free method for object dragging in images. arXiv 549 preprint arXiv:2406.01594, 2024.
 - [3] Shariq Farooq Bhat, Niloy Mitra, and Peter Wonka. Loosecontrol: Lifting controlnet for generalized depth conditioning. In SIGGRAPH, 2024. 2
 - [4] Manuel Brack, Felix Friedrich, Katharia Kornmeier, Linoy Tsaban, Patrick Schramowski, Kristian Kersting, and Apolinário Passos. Ledits++: Limitless image editing using text-to-image models. In CVPR, 2024. 2
 - [5] Mingdeng Cao, Xintao Wang, Zhongang Qi, Ying Shan, Xiaohu Qie, and Yinqiang Zheng. Masactrl: Tuning-free mutual self-attention control for consistent image synthesis and editing. In CVPR, 2023. 2
- [6] Duygu Ceylan, Chun-Hao P Huang, and Niloy J Mitra. 561 Pix2video: Video editing using image diffusion. In ICCV, 562 563 2023. 2, 3
- 564 [7] Pascal Chang, Jingwei Tang, Markus Gross, and Vinicius C 565 Azevedo. How i warped your noise: a temporally-correlated 566 noise prior for diffusion models. In ICLR, 2024. 2, 3
 - [8] Dave Epstein, Allan Jabri, Ben Poole, Alexei Efros, and Aleksander Holynski. Diffusion self-guidance for controllable image generation. NeurIPS, 36, 2023. 5
- [9] Daniel Geng and Andrew Owens. Motion guidance: Diffusion-based image editing with differentiable motion es-572 timators. ICLR, 2024. 2
- [10] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, 573 574 Yael Pritch, and Daniel Cohen-or. Prompt-to-prompt image 575 editing with cross-attention control. In ICLR, 2023. 2
- [11] Jonathan Ho and Tim Salimans. Classifier-free diffusion 576 577 guidance. In NeurIPS, 2021. 3, 6
 - [12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In NeurIPS, 2020. 3
- [13] Inbar Huberman-Spiegelglas, Vladimir Kulikov, and Tomer 580 581 Michaeli. An edit friendly ddpm noise space: Inversion and manipulations. In CVPR, 2024. 2 582
- 583 [14] Hyeonho Jeong, Jinho Chang, Geon Yeong Park, and Jong Chul Ye. Dreammotion: Space-time self-similarity 584 585 score distillation for zero-shot video editing. In ECCV, 2024. 586
- 587 [15] Hyeonho Jeong, Geon Yeong Park, and Jong Chul Ye. Vmc: 588 Video motion customization using temporal attention adap-589 tion for text-to-video diffusion models. In CVPR, 2024. 3
- 590 [16] Jaihoon Kim, Juil Koo, Kyeongmin Yeo, and Minhyuk Sung. 591 Synctweedies: A general generative framework based on 592 synchronized diffusions. In NeurIPS, 2024. 2
- 593 [17] Juil Koo, Chanho Park, and Minhyuk Sung. Posterior distil-594 lation sampling. In CVPR, 2024. 2
- 595 [18] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli 596 Shechtman, and Jun-Yan Zhu. Multi-concept customization 597 of text-to-image diffusion. In CVPR, 2023. 2

- [19] Yaowei Li, Xintao Wang, Zhaoyang Zhang, Zhouxia Wang, 598 Zivang Yuan, Liangbin Xie, Yuexian Zou, and Ying Shan. 599 Image conductor: Precision control for interactive video syn-600 thesis. arXiv preprint arXiv:2406.15339, 2024. 3 601
- [20] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In ICLR, 2023. 3
- [21] Haofeng Liu, Chenshu Xu, Yifei Yang, Lihua Zeng, and Shengfeng He. Drag your noise: Interactive point-based editing via diffusion semantic propagation. In CVPR, 2024. 2
- [22] Shaoteng Liu, Yuechen Zhang, Wenbo Li, Zhe Lin, and Jiaya Jia. Video-p2p: Video editing with cross-attention control. In CVPR, 2024. 2
- [23] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In ICLR, 2023. 3
- [24] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. SDEdit: Guided image synthesis and editing with stochastic differential equations. In ICLR, 2021. 2, 6, 7
- [25] Oscar Michel, Anand Bhattad, Eli VanderBilt, Ranjay Krishna, Aniruddha Kembhavi, and Tanmay Gupta. Object 3dit: Language-guided 3d-aware image editing. NeurIPS, 36, 2024. 2
- [26] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models. In CVPR, 2023. 8
- [27] Chong Mou, Xintao Wang, Jiechong Song, Ying Shan, and Jian Zhang. Dragondiffusion: Enabling drag-style manipulation on diffusion models. In ICLR, 2024. 2
- [28] Karran Pandey, Paul Guerrero, Matheus Gadelha, Yannick Hold-Geoffroy, Karan Singh, and Niloy J Mitra. Diffusion handles enabling 3d edits for diffusion models by lifting activations to 3d. In CVPR, 2024. 2, 3, 6, 7
- [29] Jangho Park, Gihyun Kwon, and Jong Chul Ye. Ed-nerf: Efficient text-guided editing of 3d scene using latent space nerf. In ICLR, 2024. 2
- [30] William Peebles and Saining Xie. Scalable diffusion models with transformers. In ICCV, 2023. 3, 5
- [31] Chenyang Qi, Xiaodong Cun, Yong Zhang, Chenyang Lei, Xintao Wang, Ying Shan, and Qifeng Chen. Fatezero: Fusing attentions for zero-shot text-based video editing. In ICCV, 2023. 3
- [32] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. arXiv preprint arxiv:2208.12242, 2022. 2
- [33] Xiaoyu Shi, Zhaoyang Huang, Fu-Yun Wang, Weikang Bian, Dasong Li, Yi Zhang, Manyuan Zhang, Ka Chun Cheung, Simon See, Hongwei Qin, et al. Motion-i2v: Consistent and controllable image-to-video generation with explicit motion modeling. In SIGGRAPH, 2024. 3
- [34] Yujun Shi, Jun Hao Liew, Hanshu Yan, Vincent YF Tan, and Jiashi Feng. Instadrag: Lightning fast and accurate dragbased image editing emerging from videos. arXiv preprint arXiv:2405.13722, 2024. 2

- [35] Yujun Shi, Chuhui Xue, Jun Hao Liew, Jiachun Pan, Hanshu Yan, Wenqing Zhang, Vincent YF Tan, and Song Bai.
 Dragdiffusion: Harnessing diffusion models for interactive
 point-based image editing. In *CVPR*, 2024.
- [36] Joonghyuk Shin, Daehyeon Choi, and Jaesik Park. InstantDrag: Improving Interactivity in Drag-based Image Editing.
 In *SIGGRAPH*, 2024. 2
- [37] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based
 generative modeling through stochastic differential equations. In *ICLR*, 2021. 3
- [38] Yizhi Song, Zhifei Zhang, Zhe Lin, Scott Cohen, Brian Price,
 Jianming Zhang, Soo Ye Kim, and Daniel Aliaga. Objectstitch: Object compositing with diffusion model. In *CVPR*,
 2023. 2
- [39] Yizhi Song, Zhifei Zhang, Zhe Lin, Scott Cohen, Brian Price,
 Jianming Zhang, Soo Ye Kim, He Zhang, Wei Xiong, and
 Daniel Aliaga. Imprint: Generative object compositing by
 learning identity-preserving representation. In *CVPR*, 2024.
 2
- [40] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin,
 Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov,
 Naejin Kong, Harshith Goka, Kiwoong Park, and Victor
 Lempitsky. Resolution-robust large mask inpainting with
 fourier convolutions. In WACV, 2022. 6
- [41] Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali
 Dekel. Plug-and-play diffusion features for text-driven
 image-to-image translation. In *CVPR*, 2023. 2
- [42] Jiawei Wang, Yuchen Zhang, Jiaxin Zou, Yan Zeng, Guoqiang Wei, Liping Yuan, and Hang Li. Boximator: Generating rich and controllable motions for video synthesis. In *ICML*, 2024. 3
- [43] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris
 Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vision made easy. In *CVPR*, 2024. 3
- [44] Zhouxia Wang, Ziyang Yuan, Xintao Wang, Yaowei Li, Tianshui Chen, Menghan Xia, Ping Luo, and Ying Shan.
 Motionctrl: A unified and flexible motion controller for video generation. In *SIGGRAPH*, 2024. 3
- [45] Daniel Winter, Matan Cohen, Shlomi Fruchter, Yael Pritch,
 Alex Rav-Acha, and Yedid Hoshen. Objectdrop: Bootstrapping counterfactuals for photorealistic object removal and insertion. arXiv preprint arXiv:2403.18818, 2024. 2
- [46] Zongze Wu, Nicholas Kolkin, Jonathan Brandt, Richard
 Zhang, and Eli Shechtman. Turboedit: Instant text-based
 image editing. *ECCV*, 2024. 2
- [47] Jiraphon Yenphraphai, Xichen Pan, Sainan Liu, Daniele
 Panozzo, and Saining Xie. Image sculpting: Precise object
 editing with 3d geometry control. In *CVPR*, 2024. 2
- [48] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding
 conditional control to text-to-image diffusion models. In *ICCV*, 2023. 2
- [49] Qihang Zhang, Yinghao Xu, Chaoyang Wang, Hsin-Ying
 Lee, Gordon Wetzstein, Bolei Zhou, and Ceyuan Yang.
 3ditscene: Editing any scene via language-guided disentangled gaussian splatting. *arXiv preprint arXiv:2405.18424*,
 2024. 2

- [50] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 7
 713
- [51] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all, 2024. 3, 5, 6
 714
 715
 716
 717
- [52] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: learning view synthesis using multiplane images. *ACM TOG*, 2018. 6
 720