# O-DisCo-Edit: Object Distortion Control for Unified Realistic Video Editing

Yuqing Chen[1,3*]    Junjie Wang[1✉]    Lin Liu[2✉†]    Ruihang Chu[1]

Xiaopeng Zhang[2]    Qi Tian[2]    Yujiu Yang[1]

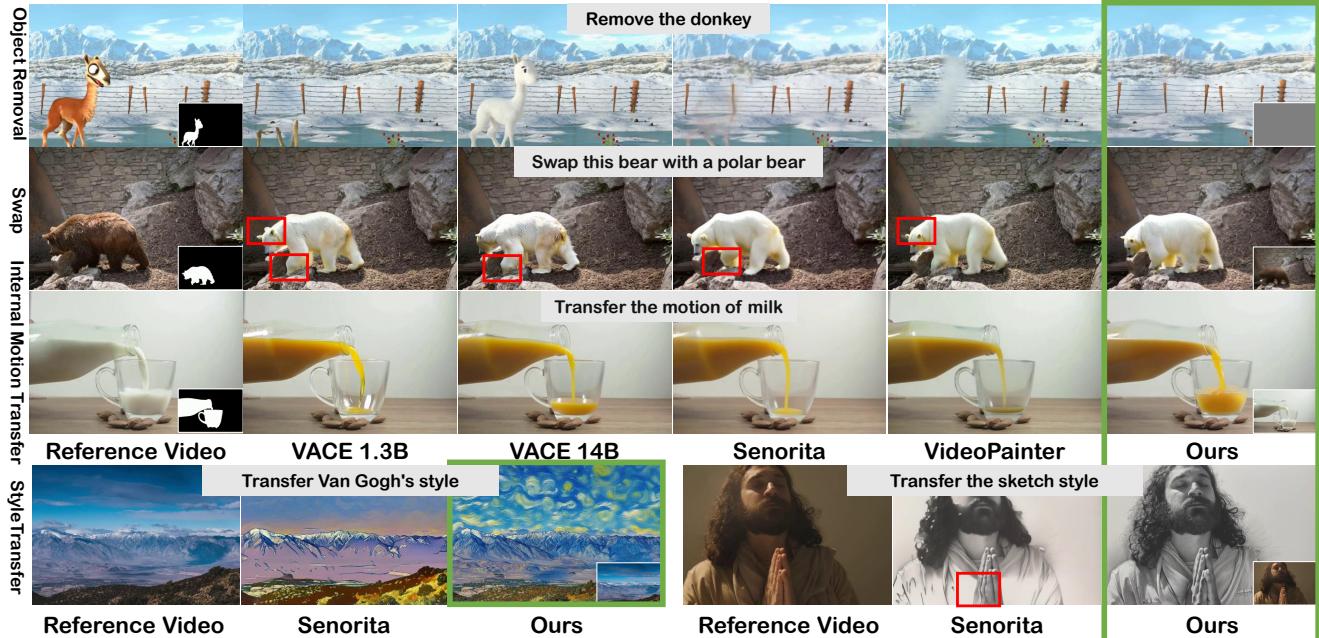[1] Tsinghua University.    [2] Huawei Inc.    [3] Pengcheng National Laboratory.

Figure 1. Given a reference video and image (typically the edited first frame), our method generates more realistic edited videos than SOTA approaches (VACE, Senorita and VideoPainter) across various tasks, including object removal, swap, object inside motion transfer, and style transfer. Zoom in to examine the visualization results. The bottom right of the reference video shows the input masks for all models, while the bottom right of our result displays our proposed novel control signal.

*This work was done during an internship at Huawei Inc.

†Project Leader.

✉Corresponding authors: wangjunjie@sz.tsinghua.edu.cn, ll0825@mail.ustc.edu.cn

## Abstract

*Diffusion models have recently advanced video editing, yet controllable editing remains challenging due to the need for precise manipulation of diverse object properties. Current methods require different control signal for diverse editing tasks, which complicates model design and demands significant training resources. To address this, we propose O-DisCo-Edit, a unified framework that incorporates a novel object distortion control (O-DisCo). This signal, based on random and adaptive noise, flexibly encapsulates a wide range of editing cues within a single representation. Paired with a "copy-form" preservation module for preserving non-edited regions, O-DisCo-Edit enables efficient, high-fidelity editing through an effective training paradigm. Extensive experiments and comprehensive human evaluations consistently demonstrate that O-DisCo-Edit surpasses both specialized and multitask state-of-the-art methods across various video editing tasks. https://cyqii.github.io/O-DisCo-Edit.github.io/.*

## 1. Introduction

Recent years have witnessed remarkable advancements in diffusion-based video generation [1–4]. Beyond pure generation, video editing has emerged as a crucial extension, which enables modifications to reference videos based on

Figure 2. Comparisons of different object properties, control signals, and models.

Table 1. Comparison of training configurations for different models. * indicates that the majority of the module is used for training. "Block" refers to a DiT block.

| Model | Dataset | Trainable Module | Steps | GPUs |
|---|---|---|---|---|
| VACE | Mutil-Task | 8 Blocks | 200K | 128 A100 |
| Senorita | Mutil-Task | 102 Blocks* | 4 epoch | \ |
| VideoPainter | 390.3k | 2 Blocks, 1 LoRA | 82K | 64 V100 |
| Ours | 180k | Two LoRAs | 7.55K | 8 A800 |

user instructions. Specifically, effective video editing necessitates precise control over the content within edited regions, while flawlessly preserving unedited areas.

For controllable video editing, single-task editing models incorporate additional control signals such as 2D bounding boxes [5–7], masks [8, 9], optical flow [10–12], and tracking points [13–15] to improve control precision. However, as shown in Fig. 2, conditions like bounding boxes and masks provide limited information, thus hindering fine-grained control for complex editing scenarios. Furthermore, video datasets with optical flow and tracking are scarce, and their extraction is often complex and prone to inaccuracies. These two issues make precise and intricate controllable editing difficult.

Single-task editing models, as discussed above, are no longer sufficient to meet user diverse demands. Consequently, unified multi-task video editing approaches [16–20] have emerged, which can accomplish diverse editing tasks by introducing various signals. However, they typically demand complex training pipelines. This complexity necessitates the construction of specialized multi-task datasets and the design of task-specific modules (e.g., multiple DiT blocks), resulting in a large number of trainable parameters, as shown in Tab. 1. Furthermore, it requires integrating various conditions across numerous training stages, often demanding tens of thousands of steps.

Despite their design incorporating diverse signals, most of multi-task video editing model are inflexible during inference, as they can generally process only one control condition at a time. This prevents the model from leveraging complementary cues from multiple signals, thereby hindering the flexible transition between fine-grained and coarse-grained editing for the same task.

To address these challenges, we propose a novel unified control signal: the **o**bject **dis**tortion **co**ntrol (O-DisCo). This signal is generated by applying appropriate noise to the edited objects, effectively acting as a distortion signal for the reference video. As illustrated in Fig. 2, all other control signals can similarly be viewed as specific types of reference video distortion. Therefore, by controlling the noise, O-DisCo inherently unifies all these signals into a single representation. For training, randomness is introduced into O-DisCo. This significantly simplifies training dataset construction and model design, saving substantial training resources, as shown in Tab. 1. For inference, by adaptively manipulating the intensity and scope of O-DisCo's noise, our model can perform a wide range of tasks. While the above design primarily focuses on the edited regions, a "copy-form" preservation module is further designed to address the preservation of non-edited areas. Encapsulating these capabilities, we propose O-DisCo-Edit, a unified framework for versatile video editing.

Our comprehensive experiments confirm O-DisCo-Edit's effectiveness and versatility across diverse tasks, including object removal, outpainting, and transfers of motion, lighting, color, and style. Specifically, O-DisCo-Edit consistently surpasses the current state-of-the-art (SOTA) multi-task editing model, VACE [20], on the majority of tasks. Notably, for the object removal task on the OmnimatteRF [21] benchmark, our method also demonstrates superior performance compared to the specialized SOTA removal approach, MiniMax-Remover [22].

Overall, the contribution of this work is summarized as follows:

- A novel unified control signal, **o**bject **dis**tortion **co**ntrol (O-DisCo) is proposed to substantially reduce training resource demands and enable flexible, precise multi-task video editing from coarse to fine granularity.
- We propose a "copy-form" preservation module for non-edited region preservation, which enhances the model's ability to maintain unedited areas.
- Our proposed O-DisCo-Edit, achieving new SOTA performance across diverse tasks, offers a novel perspective for developing unified video editing frameworks.

## 2. Related Work

**Single-Task Video Editing and Control Signals.** Video editing tasks frequently require additional control signals (e.g., masks, poses, optical flows, tracking points) to mod-
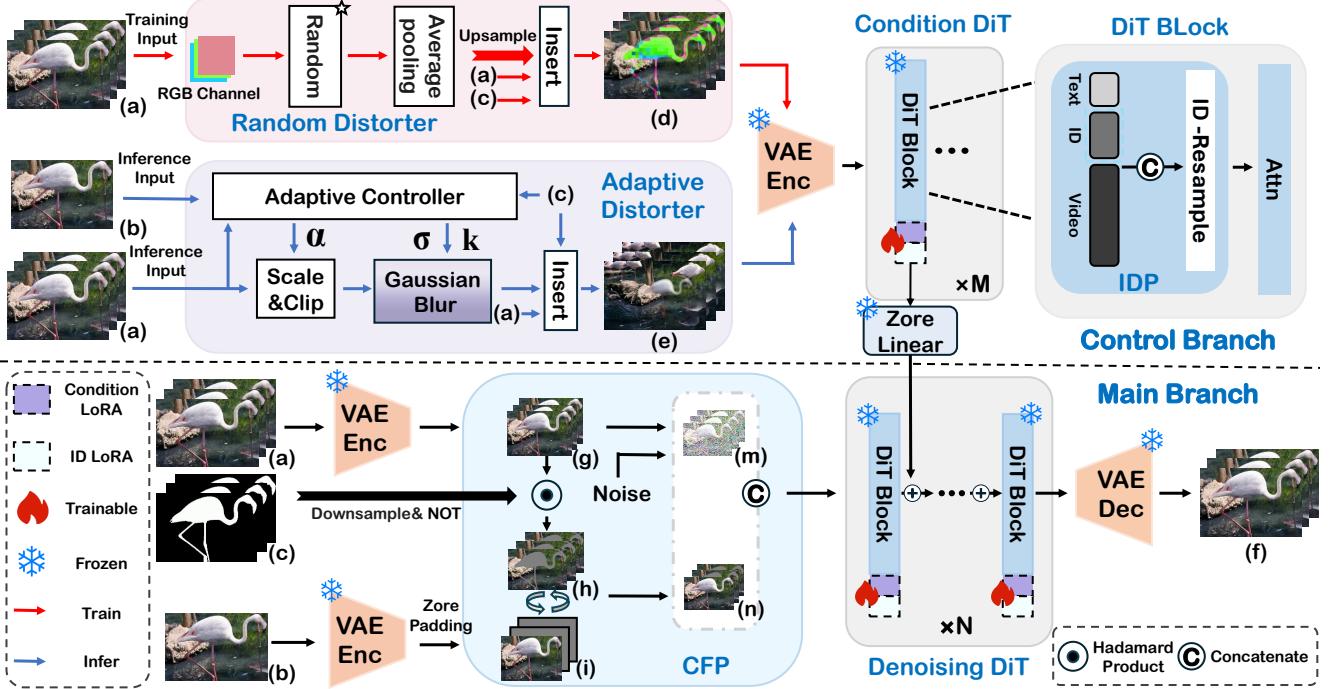
Figure 3. The framework of the proposed O-DisCo-Edit. (a) Reference video. (b) Reference image (first frame during training, edited image during inference). (c) Masks. (d) R-O-DisCo. (e) A-O-DisCo. (f) Generated video. (g) Latent of reference video. (h) Latent of the preserved region. (i) Image latent with zero-padding. (m) Noisy Latent. (n) Image Latent with the latent of preserved region. $\alpha$ represents the contrast, $\sigma$ represents the intensity of the added noise, and $k$ is the size of the gaussian blur kernel. The adaptive distorter generates A-O-DisCo for inference, and the random distorter generates R-O-DisCo for training. The CFP ensures the preservation of unedited areas. The IDP maintains object appearance consistency.

ify reference video attributes. VideoAnydoor [5] introduces masks and tracking points for object insertion, while DiffuEraser [13] leverage masks for object removal. Followyour-Canvas [23] employs 2D bounding boxes for outpainting, and ReCamMaster [24] uses camera trajectories for camera control.

**Multi-Task Video Editing.** Growing demands for creative versatility have driven the development of multi-task video editing. VACE [20] integrates sophisticated signals like optical flow and masks with a context embedder and adapter to perform tasks such as swap, animation, and outpainting. Similarly, Senorita [25] utilizes masks, canny edges, and other cues, coupled with four specialized expert models, to achieve tasks like addition, removal, and swap. Therefore, multi-task editing often demands complex training pipelines with diverse signals, specialized modules, and multi-stage training [16, 19, 20, 25]. In contrast, our proposed unified O-DisCo signal enables multi-task completion with significantly fewer training resources.

**Adaptive Inference.** Current video editing models [5, 16, 19, 20, 25] typically rely on a single control signal during inference, which limits their adaptability for multi-grained editing. Adaptive inference, conversely, allows models

to adjust outputs dynamically based on varying reference videos, images, or prompts. This is common in training-free models; for example, LMP [26] applies attention values as a loss to optimize hidden states for appearance similarity. Likewise, the authors in [27] defines motion consistency loss for gradient descent on noisy latent vectors to achieve motion transfer. Building on this principle, O-DisCo-Edit also leverages adaptive inference through the proposed O-DisCo, which dynamically adjusts injected noise based on reference videos and images for multi-grained control.

## 3. Method

As shown in Fig. 3, our approach introduces a first-frame-guided video editing model, building upon the CogVideoX-I2V architecture [1]. The object distortion control (O-DisCo), derived via the distorter, is fed into the VAE and conditional DiT for precise control over edited regions. Concurrently, a "copy-form" preservation (CFP) module processes the reference image and video, which then provides their latent output to the denoising DiT for preservation of non-edited areas. Additionally, an identity preservation (IDP) module is proposed to enhance ID fidelity within edited regions. Subsequent sections detail O-DisCo's con-

struction and the design of the CFP and IDP modules.

## 3.1. Random Object Distortion Control

During the training phase, we apply a random distoter to generate random object distortion control (R-O-DisCo). As shown in the top-left part of Fig. 3, we intentionally distort the colors of the reference video $V_{\text{ref}} \in \mathbb{R}^{F \times H \times W \times 3}$ to prevent the model from simply replicating original color information, where $F$ is the frame number, $H$ and $W$ are the height and width of the reference video. This involves applying **random** arithmetic operations to each RGB channel, and the resulted color-distorted video $V_{\text{cd}} \in \mathbb{R}^{F \times H \times W \times 3}$ is formulated as:

$$V_{\text{cd}}[:,:,:,i] = \text{clip}(V_{\text{ref}}[:,:,:,i] \star a_i), \ i \in \{0,1,2\}, \quad (1)$$

where $a_i, i \in \{0,1,2\}$ are random real numbers, with their specific ranges detailed in the Appendix A. Moreover, $\star$ denotes a randomly chosen arithmetic operation (addition, subtraction, multiplication, or division). The clip operation constrains output values to the range $[0, 255]$.

After that, $V_{\text{cd}}$ is downsampled via average pooling and then upsampled using nearest-neighbor interpolation, both by a factor of $L \in \mathbb{Z}$. $L$ is a randomly sampled integer resampling factor (more details in Appendix A). This operation intentionally disrupts fine-grained structural details, producing a mosaic-like effect video $V_{\text{cdm}} \in \mathbb{R}^{F \times H \times W \times 3}$. Consequently, the model is compelled to primarily learn video generation guided by the first frame's appearance, rather than relying on the precise visual information of $V_{\text{cd}}$.

Finally, We utilize masks $M \in \mathbb{R}^{F \times H \times W \times 1}$ to insert the edited object from $V_{\text{cdm}}$ into $V_{\text{ref}}$, which produces the R-O-DisCo $V_{\text{RODC}}$ ((d) in Fig. 3) during the training stage:

$$V_{\text{RODC}} = V_{\text{cdm}} \odot M + V_{\text{ref}} \odot (\mathbf{1} - M), \quad (2)$$

where $\odot$ represents the Hadamard product.

Overall, during the training phase, we enhance the model's robustness and task adaptability by increasing the randomness of O-DisCo.

## 3.2. Adaptive Object Distortion Control

During inference phase, our model adapts to specific tasks or instructions via adaptive object distortion control (A-O-DisCo) (highlighted by (e) in Fig. 3), which is implemented by an adaptive distorter. It is achieved through contrast modification (scaling and clipping) and dynamic noise injection within the editable regions. The process is formally represented by the following equations:

$$V_{\text{c}}(f,x,y) = \text{clip}\big(\alpha \cdot V_{\text{ref}}(f,x,y)\big),$$

$$V_{\text{cn}}(f,x,y) = \sum_{i=-b}^{b} \sum_{j=-b}^{b} V_{\text{c}}(f,x+i,y+j) \cdot G_{\text{norm}}(i,j;\sigma),$$

$$V_{\text{AODC}} = V_{\text{cn}} \odot M + V_{\text{ref}} \odot (\mathbf{1} - M),$$

$$(3)$$

where $V_{\text{ref}}(f,x,y)$ denotes the pixel value at coordinates $(x,y)$ in the $f$-th frame of the reference video. $V_{\text{c}}(x,y)$ and $V_{\text{cn}}(x,y)$ represent the $V_{\text{ref}}$ after scale&clip and gaussian blur, respectively. $G_{\text{norm}}(i,j;\sigma)$ is the normalized gaussian blur kernel. $\alpha$ represents the contrast, $\sigma$ is the noise intensity, and $k = 2b + 1$ is the gaussian kernel size.

The adaptive controller determines suitable values for $\alpha$, $\sigma$, and $k$ by calculating two similarities: $\mathbf{Sim}_i$, the edited region's edge map similarity between the reference image and the reference video's first frame; $\mathbf{Sim}_{\text{v}}$, the intra-frame similarity within the reference video's edited region edge map. Empirically, fitting these three parameters using a quadratic polynomial of two similarity yields superior results (see specific formulas in Appendix A). Finally, the A-O-DisCo $V_{\text{AODC}}$ obtained for each task is shown in the last line of Tab. 2. Specifically, during inference for object removal and outpainting (R&O), we set $V_{\text{AODC}}$ to zero, ensuring that no additional information is introduced. This allows the model to generate the video based on other condition, thereby reducing artifacts.

## 3.3. "Copy-Form" Preservation Module

Many video editing methods [16, 20, 28] typically integrate non-edited regions with the control signals via extra branch. However, such integration often leads to mutual interference between the preserved regions and control signals, thereby limiting editing flexibility. Instead, we propose the "copy-form" preservation (CFP) module illustrated Fig. 3, which enhances editing flexibility by integrating the non-edited regions directly into the main branch of the network. Detailedly, CFP replaces conventional zero-padding (denoted as (i) in Fig. 3) with the latent of the preserved region $z_{\text{p}}^{\text{v}'}$ (marked as (h) in Fig. 3), to obtain $z_{\text{images}}$ (denoted as (n) in Fig. 3). This process is expressed as:

$$z_{\text{p}}^{\text{v}'} = z_{\text{ref}}^{\text{v}} \odot (\mathbf{1} - z_{\text{mask}}[1:]),$$
$$z_{\text{images}} = [z_{\text{ref}}^{\text{i}}, z_{\text{p}}^{\text{v}'}], \quad (4)$$

where $z_{\text{ref}}^{\text{v}}$ denotes the latent of the reference video, $z_{\text{mask}}$ represents the downsampled binary mask (with the same shape as $z_{\text{ref}}^{\text{v}}$), $z_{\text{ref}}^{\text{i}}$ is reference image latent. $[1:]$ corresponds to a slicing operation, and $[,]$ signifies tensor concatenation. The $z_{\text{images}}$ for each task are shown in Tab. 2. Notably, the CFP module achieves preservation of non-edited regions with an effect similar to "first-frame copying".

## 3.4. Identity Preservation Module

To mitigate object appearance changes during complex motion or occlusion, we design the identity preservation (IDP) module, illustrated in the upper right corner of Fig. 3. Specifically, we extract position-agnostic tokens (ID tokens) from the reference image's edited regions and con-

Table 2. Adaptive inference conditions for different tasks. R&O in the first line means object removal and outpainting.

| Condition | R&O | Style Tranfer | Other Tasks |
|---|---|---|---|
| $z_{\text{images}}$ | $[z_{\text{ref}}^{\text{i}}, z_{\text{p}}^{\text{v}'}]$ | $[z_{\text{ref}}^{\text{i}}, \mathbf{0}]$ | $[z_{\text{ref}}^{\text{i}}, z_{\text{p}}^{\text{v}'}]$ |
| $V_{\text{AODC}}$ | $\mathbf{0}$ | $V_{\text{cn}} \odot M$ $+V_{\text{ref}}(1 - M)$ | $V_{\text{cn}} \odot M$ $+V_{\text{ref}}(1 - M)$ |

catenate them with text tokens. Akin to text tokens, ID tokens act as a global guide which make the model leverage ID information throughout the video generation. Further enhancing the model's focus on ID consistency, we introduce ID-Resample to extract the key (K) and value (V) vectors from the edited regions of the generated video. These are concatenated with the K and V vectors of the original generated video, and the process compels the model to reinforce ID consistency within the edited regions.

## 4. Experiment

**Implementation Details.** Training dataset: we utilize approximately 180k video-mask pairs from the Senorita-2M grounding dataset [25]. All video-mask pairs are center-cropped and resized to a $720 \times 480$ resolution with a length of 49 frames. Moreover, prompts for masked regions are generated via Qwen2.5-VL-7B [29]. Two-stage training: Our model builds on the frozen pre-trained weights of Diffusion as Shader [13]. Firstly, condition LoRA is trained with the random distorter and CFP module (2400 steps); Secondly, the IDP module's dedicated ID LoRA is trained (5150 steps). All stages employ AdamW optimization (learning rate $1 \times 10^{-4}$) on 8 A800 GPUs with gradient accumulation for a batch size of 32.

**Baseline Methods.** For the majority of tasks, we select the SOTA unified video editing methods VACE [20], VideoPainter [28], and Senorita [25] as our primary baselines. Additionally, for the object removal, we include DiffuEraser [30], MiniMax-Remover (MiniMax) [22], and Propainter [31] as extra baselines. For the style transfer, Senorita [25] is chosen.

**Benchmarks.** We curated a benchmark from DIVAS [32] and VPData [28], specifically targeting challenging scenarios with internal motion, lighting variations, and complex object movements. For prompts, Senorita utilized the dedicated prompt, as required by its inference process, while others used same prompts from Qwen2.5VL-7B [29]. Additionally, the reference image was the first frame edited using either HiDream-E1 [33] or commercial models [1]. Subsequently, the edited frame served as input for O-DisCo-Edit and multi-task baselines. For fair comparison with Senorita,

---

only the first 33 frames of generated videos are evaluated. In parallel, OmnimatteRF [21] is selected as a benchmark for the object removal task. Further details about our benchmark's construction are provided in Appendix B.

**Metrics.** The evaluation includes automatic scoring and a manual user study. Automatic scoring metrics: (1) Non-Edited Region Preservation: Fidelity in unedited regions is assessed using PSNR ($\text{PSNR}_{\text{P}}$) and SSIM ($\text{SSIM}_{\text{P}}$). (2) Alignment: CLIP Similarity [34] (CLIP-T) measures semantic consistency between the generated video and its caption. Appearance consistency [35] (CLIP-$\text{I}_{\text{E}}$) between the output video and the reference image is calculated within the edited regions. (3) Video Generation Quality: Overall video quality is assessed via FVD [36], ArtFID [37], PSNR, SSIM, and temporal consistency (TC) [35, 38]. (4) Normalized Average Score: This score is obtained by following the work of [39] using Min-Max Normalization, with all metrics (excluding CLIP-T) weighted equally. Specifically, for the style transfer task, CFSD [40] is applied to evaluate the preservation of the reference video's content. Meanwhile, for the object removal task, we measure removal ability by calculating the SSIM ($\text{SSIM}_{\text{E}}$) and PSNR ($\text{PSNR}_{\text{E}}$) between the edited regions of the output video and the corresponding background video. Manual Assessment: The mean opinion score (MOS) is adopted, focusing on editing completeness (EC) and video quality (VQ). Anonymized generated data is randomly distributed to participants for 1-5 scale scoring. More details are in the Appendix B.

### 4.1. Comparison with State-of-the-Arts

As shown in Tab. 3, we conduct comprehensive comparisons between O-DisCo-Edit and baselines across multiple tasks including object removal, outpainting, object internal motion transfer, lighting transfer, color change, swap, addition, and style transfer. Our method demonstrates superior performance across all these tasks.

**Object Removal.** Quantitatively, our method obtains optimal results in both the Remove (49) (49-frame videos) and Remove (33) (33-frame videos) settings (Tab. 3 (a)). As shown in Fig. 4, O-DisCo-Edit successfully avoids the background damage seen in Propainter and DiffuEraser, as well as the bicycle overlap present in MiniMax-Remover. In comparison with multi-task models in Fig. 1, baselines consistently exhibit prominent artifacts, which indicate unsuccessful removal.

**Outpainting.** For evaluation, we outpaint videos from $280 \times 520$ to $480 \times 720$. As demonstrated in Tab. 3 (b), O-DisCo-Edit establishes new SOTA records across all metrics. VACE generates grainy textures in edited areas, Senorita and VideoPainter exhibit noticeable box-like artifacts shown in Fig. 5. In contrast, O-DisCo-Edit creates exceptionally well-blended, natural, and continuous results.

Table 3. Comparison of different models on various tasks using our benchmark (OmnimatteRF benchmark for object removal). The evaluation includes automatic scoring and a manual user study. The best results are in **bold**, while the second best are underlined. "Preservation" means non-edited region preservation, "TC" denotes temporal consistency, "EC" represents editing completeness, and "VQ" stands for visual quality.

| Metrics | | Video Quality | | | | Removal Capability | | Normalized | User Stuty | |
|---|---|---|---|---|---|---|---|---|---|---|
| Task | Method | TC↑ | FVD↓ | PSNR↑ | SSIM↑ | $SSIM_E$ ↑ | $PSNR_E$ ↑ | Avg. Score↑ | EC↑ | VQ↑ |
| (a.1) Object Removal (49) | DiffuEraser | 0.9964 | 422.7 | 27.89 | <u>0.9207</u> | 0.9713 | 34.09 | 0.2682 | 3.122 | 2.867 |
| | MiniMax | <u>0.9973</u> | <u>373.2</u> | 27.15 | 0.8732 | **0.9737** | <u>34.87</u> | 0.4816 | <u>3.567</u> | <u>3.333</u> |
| | propainter | 0.9971 | 410.3 | **28.30** | **0.9224** | 0.9715 | 34.12 | <u>0.4844</u> | 3.044 | 2.822 |
| | O-DisCo-Edit | **0.9974** | **300.3** | <u>28.05</u> | 0.8751 | <u>0.9730</u> | **35.43** | **0.7553** | **3.967** | **3.689** |
| (a.2) Object Removal (33) | VACE 1.3B | 0.9934 | 1376 | 23.46 | <u>0.8508</u> | 0.9551 | 26.17 | 0.4233 | 2.011 | 1.911 |
| | VACE 14B | 0.9896 | 2085 | 22.29 | 0.8372 | 0.9439 | 24.28 | 0.1316 | 1.578 | 1.444 |
| | Senorita | <u>0.9962</u> | <u>662.2</u> | <u>26.24</u> | 0.8387 | <u>0.9681</u> | <u>32.77</u> | <u>0.7058</u> | <u>3.311</u> | <u>3.156</u> |
| | VideoPainter | 0.9871 | 2403 | 21.28 | 0.8303 | 0.9452 | 23.42 | 0.0072 | 1.744 | 1.578 |
| | O-DisCo-Edit | **0.9969** | **360.1** | **28.28** | **0.8719** | **0.9740** | **36.05** | **1.000** | **3.956** | **3.689** |

| Metrics | | Video Quality | | | Alignment | Preservation | | Normalized | User Stuty | |
|---|---|---|---|---|---|---|---|---|---|---|
| Task | Method | TC↑ | FVD↓ | PSNR↑ | CLIP-T↑ | $PSNR_P$ ↑ | $SSIM_P$ ↑ | Avg. Score↑ | EC↑ | VQ↑ |
| (b) Outpainting | VACE 1.3B | 0.9976 | <u>88.03</u> | 25.11 | 11.92 | 31.62 | 0.9383 | <u>0.6801</u> | <u>4.244</u> | <u>3.933</u> |
| | VACE 14B | <u>0.9977</u> | 88.75 | 23.52 | 12.03 | 30.08 | 0.9381 | 0.4972 | 4.178 | 3.911 |
| | Senorita | **0.9978** | 177.1 | <u>25.20</u> | \ | 30.90 | 0.9262 | 0.4784 | 2.889 | 2.600 |
| | VideoPainter | 0.9958 | 325.4 | 24.54 | **12.95** | <u>31.81</u> | <u>0.9442</u> | 0.3384 | 2.156 | 1.978 |
| | O-DisCo-Edit | **0.9978** | **77.03** | **26.43** | <u>12.18</u> | **33.87** | **0.9466** | **1.000** | **4.289** | **4.067** |

| Metrics | | Video Quality | | Alignment | | Preservation | | Normalized | User Stuty | |
|---|---|---|---|---|---|---|---|---|---|---|
| Task | Method | TC↑ | ArtFID↓ | CLIP-T↑ | $CLIP\text{-}I_E$ ↑ | $PSNR_P$ ↑ | $SSIM_P$ ↑ | Avg. Score↑ | EC↑ | VQ↑ |
| (c) Object Internal Motion Transfer | VACE 1.3B | **0.9946** | 7.329 | 18.75 | 93.73 | **36.41** | **0.9586** | <u>0.8515</u> | 3.011 | 2.533 |
| | VACE 14B | 0.9937 | 7.025 | <u>18.96</u> | 93.39 | 34.47 | <u>0.9582</u> | 0.7641 | <u>3.122</u> | <u>2.800</u> |
| | Senorita | <u>0.9940</u> | **6.628** | \ | <u>93.94</u> | 29.45 | 0.8477 | 0.5229 | 3.022 | 2.333 |
| | VideoPainter | 0.9908 | 8.201 | **19.47** | 91.57 | <u>36.32</u> | 0.9410 | 0.3657 | 2.300 | 1.822 |
| | O-DisCo-Edit | 0.9927 | <u>6.712</u> | 18.82 | **94.64** | 35.88 | 0.9530 | **0.8639** | **4.178** | **3.756** |
| (d) Lighting Transfer | VACE 1.3B | **0.9964** | **5.991** | 20.26 | 95.30 | 31.46 | <u>0.9292</u> | <u>0.7700</u> | 3.067 | 2.644 |
| | VACE 14B | <u>0.9958</u> | 6.187 | 20.35 | 95.32 | 30.69 | 0.9261 | 0.5489 | <u>3.411</u> | <u>3.067</u> |
| | Senorita | **0.9964** | 6.478 | \ | **96.15** | 28.71 | 0.9011 | 0.4000 | 3.033 | 2.400 |
| | VideoPainter | 0.9951 | 6.378 | **21.92** | 94.76 | <u>32.93</u> | **0.9325** | 0.4160 | 2.911 | 2.489 |
| | O-DisCo-Edit | 0.9956 | <u>6.043</u> | <u>20.86</u> | <u>96.05</u> | **33.54** | 0.9285 | **0.8157** | **3.978** | **3.689** |
| (e) Change Color | VACE 1.3B | <u>0.9955</u> | 8.150 | 11.89 | 97.16 | 30.55 | <u>0.9056</u> | <u>0.7838</u> | 3.633 | 3.244 |
| | VACE 14B | 0.9954 | 8.485 | 11.63 | 96.55 | 29.76 | 0.9036 | 0.5703 | 3.456 | 3.089 |
| | Senorita | **0.9959** | **8.002** | \ | **97.67** | 27.52 | 0.8724 | 0.6000 | **4.033** | **3.711** |
| | VideoPainter | 0.9943 | 9.388 | **12.89** | 96.41 | <u>30.99</u> | **0.9136** | 0.3996 | 3.633 | 3.267 |
| | O-DisCo-Edit | <u>0.9955</u> | <u>8.008</u> | <u>11.94</u> | <u>97.49</u> | **31.00** | 0.9049 | **0.8787** | <u>3.944</u> | <u>3.689</u> |

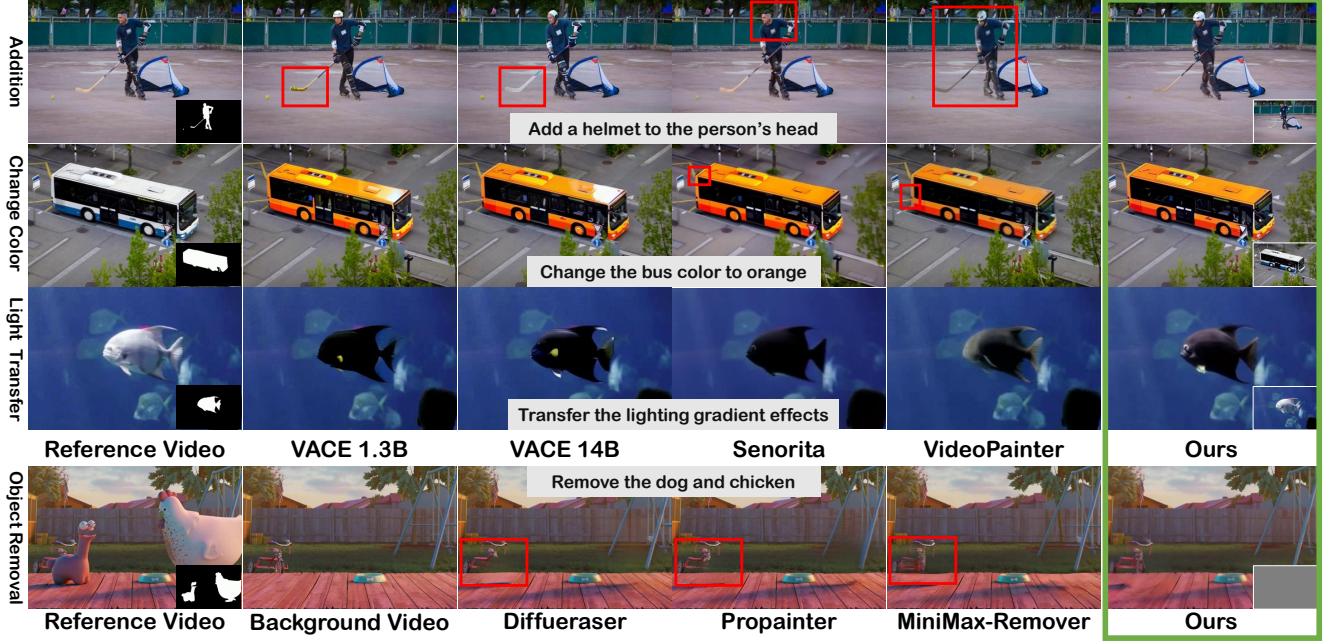| Task | Method | TC↑ | FVD↓ | CLIP-T↑ | $CLIP\text{-}I_E$ ↑ | $PSNR_P$ ↑ | $SSIM_P$ ↑ | Avg. Score↑ | EC↑ | VQ↑ |
|---|---|---|---|---|---|---|---|---|---|---|
| (f) Swap | VACE 1.3B | <u>0.9843</u> | <u>688.2</u> | 14.73 | 91.28 | <u>26.36</u> | 0.8062 | **0.7068** | 3.467 | 2.956 |
| | VACE 14B | 0.9841 | **642.3** | 14.83 | 91.03 | 25.76 | 0.8041 | <u>0.6959</u> | <u>3.556</u> | <u>3.089</u> |
| | Senorita | **0.9845** | 803.7 | \ | **92.76** | 23.66 | 0.7436 | 0.4000 | 3.456 | 2.956 |
| | VideoPainter | 0.9815 | 731.2 | **15.91** | 90.05 | **27.51** | **0.8295** | 0.4899 | 2.967 | 2.178 |
| | O-DisCo-Edit | 0.9839 | 711.8 | <u>15.27</u> | <u>91.84</u> | 26.25 | <u>0.8098</u> | 0.6950 | **4.033** | **3.689** |
| (g) Addition | VACE 1.3B | 0.9862 | 512.7 | <u>21.04</u> | 92.92 | 27.47 | <u>0.8094</u> | 0.3419 | 3.370 | 2.489 |
| | VACE 14B | <u>0.9873</u> | <u>398.9</u> | **21.41** | 92.58 | 26.86 | 0.8082 | 0.3621 | 3.200 | 2.578 |
| | Senorita | **0.9891** | **316.8** | \ | **95.39** | 27.89 | 0.7934 | **0.7375** | <u>3.496</u> | <u>3.089</u> |
| | VideoPainter | 0.9836 | 560.9 | 20.77 | 93.64 | **28.36** | **0.8246** | 0.4754 | 3.104 | 2.911 |
| | O-DisCo-Edit | 0.9871 | 448.3 | 21.03 | <u>95.37</u> | <u>28.03</u> | 0.8048 | <u>0.6470</u> | **4.037** | **3.822** |

Figure 4. Our O-DisCo-Edit method is compared against other baselines for addition, color change, light transfer, and object removal. The bottom right of the reference video displays the input masks utilized by all models, while the corresponding position in our results highlights the A-O-DisCo required by our approach.
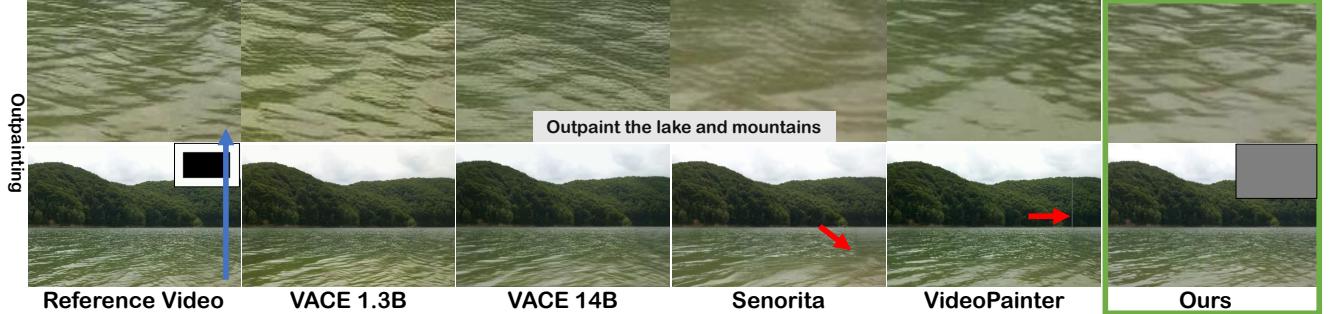


Figure 5. A comparison of O-DisCo-Edit and other baselines on the outpainting task. In the second row, the top right of the reference video displays the input masks utilized by all models, while the same position in our results highlights the A-O-DisCo required by our approach. The blue arrow indicates the source region for the magnified view presented in the first row.

**Object Internal Motion Transfer.** Existing control signals struggle to accurately capture intricate object internal motions, such as the milk flow depicted in the Fig. 1. Our method addresses this by leveraging the masks of the bottle and cup to generate a corresponding A-O-DisCo, thereby achieving precise transfer of internal object motion. In contrast, other baseline approaches are unable to accomplish accurate internal motion transfer. As shown in the quantitative results in Tab. 3 (c), we achieve the best overall metrics. Consequently, our method yields the most superior internal motion transfer results.

**Lighting Transfer.** O-DisCo-Edit is capable of simultaneously editing objects and transferring their lighting transfor-

mation. Quantitative results, presented in Tab. 3 (d), demonstrate that we achieve the best overall metrics. Qualitatively, as illustrated in Fig. 4, the original image exhibits obviously lighting and shadow changes. All other baseline methods fail to transfer these variations, whereas our approach successfully achieves excellent transfer performance. Consequently, our method yields the most superior lighting and shadow transfer results.

**Color Change.** O-DisCo-Edit is capable of changing color while preserving intrinsic characteristics, an advantage supported by both quantitative and qualitative results. In quantitative analysis, our approach achieves the highest average score as shown in Tab. 3 (e). Qualitatively, a visual compar-

Table 4. Ablation Studies on swap and object removal task. ① CFP module, ② A-O-DisCo, ③ IDP module. "w/o ①②③" denotes training with R-O-DisCo and inference with a fixed signal, entirely omitting IDP and CFP modules. "w/o ②③" indicates training without module IDP and inference with a fixed signal. "w/o ②" refers to using a fixed inference signal. "w/o ③" indicates training without the IDP module. "TC" denotes temporal consistency.

| Metric | Video Quality | | Alignment | | Preservation | | Video Quality | | Removal Capability | | Preservation | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Model | TC↑ | FVD↓ | CLIP-T↑ | CLIP-I$_E$↑ | PSNR$_P$↑ | SSIM$_P$↑ | TC↑ | FVD↓ | PSNR$_E$↑ | SSIM$_E$↑ | PSNR$_P$↑ | SSIM$_P$↑ |
| | **Task: Swap** | | | | | | **Task:Object Removal (33)** | | | | | |
| w/o ①②③ | 0.9837 | 887.6 | **15.92** | **92.13** | 21.57 | 0.6534 | **0.9883** | 2428 | 22.30 | 0.9164 | 25.41 | 0.8195 |
| w/o ②③ | 0.9834 | 866.3 | <u>15.75</u> | 91.03 | 25.32 | 0.8037 | <u>0.9882</u> | 2167 | 20.80 | 0.9106 | 27.17 | 0.8913 |
| w/o ② | <u>0.9839</u> | <u>776.8</u> | 15.25 | 91.69 | <u>26.09</u> | <u>0.8087</u> | <u>0.9882</u> | 2103 | 21.44 | 0.9254 | 27.69 | <u>0.8942</u> |
| w/o ③ | <u>0.9839</u> | 796.7 | 15.41 | 91.60 | 25.47 | 0.8053 | 0.9968 | <u>362.8</u> | <u>35.90</u> | <u>0.9737</u> | <u>29.10</u> | **0.9029** |
| O-DisCo-Edit | **0.9840** | **711.8** | 15.27 | <u>91.84</u> | **26.25** | **0.8096** | 0.9969 | **360.1** | **36.05** | **0.9740** | **29.11** | **0.9029** |

ison in Fig. 4 reveals that VACE produces irregular color gradients, while Senorita and VideoPainter generate subtle artifacts. Therefore, our approach avoids these issues and yields superior color transformation results. Notably, Senorita's top score in user studies comes from its first-frame propagation strategy. While this strategy creates high visual consistency, it does so at the expense of poor preservation in non-edited regions.

**Swap.** Quantitative evaluation in Tab. 3 (f) shows O-DisCo-Edit's performance is second to VACE, yet we achieve a higher CLIP-I$_E$. As shown in Fig. 1, VACE 14B struggles with ID consistency. Meanwhile, VACE 1.3B and VideoPainter overfit masks boundaries, generating anatomically incorrect outputs (e.g., polar bears with three ears). Furthermore, Senorita, VACE 14B, and VACE 1.3B exhibit motion inconsistencies (red box). Conversely, our method exhibits superior visual results, as evidenced by user study.

**Addition.** O-DisCo-Edit enables adding new objects to existing moving objects in a video. As shown in Tab. 3 (g), O-DisCo-Edit reaches competitive performance, ranking second only to Senorita. However, as Fig. 4 illustrated, Senorita fail to complete the addition task, with its high metrics solely due to "copying" the original video. Therefore, our method attains the most preferred additions results in user study.

Table 5. Comparison of different models' performance metrics on style transfer. "Preservation" means non-edited region preservation, "TC" denotes temporal consistency, "EC" represents editing completeness, and "VQ" stands for visual quality.

| Metrics | Video Quality | Preservation | User Study | |
|---|---|---|---|---|
| Method | TC↑ ArtFID↓ | CFSD↓ | EC↑ | VQ↑ |
| Senorita2m | **0.9960** 7.979 | **0.0933** | 2.989 | 2.578 |
| O-DisCo-Edit | 0.9954 **7.292** | 0.2055 | **4.322** | **4.156** |

**Style Transfer.** O-DisCo-Edit attains the highest ArtFID, as shown in Tab. 5. In contrast, Senorita exhibits a very low

CFSD, which indicates a tendency for its generated videos to align with the original reference content. As depicted in Fig. 1, such alignment is detrimental to style transfer quality. Therefore, our method consistently received higher user study evaluations.

### 4.2. Ablation Analysis

As shown in Tab. 4, we ablate on O-DisCo-Edit. (1) Comparing row 1 and row 2, a significant improvement in both PSNR$_P$ and SSIM$_P$ is observed with the inclusion of the CFP module. Thus this results validate CFP module's effectiveness in non-edited regions preservation. (2) When contrasting row 3/4 with row 2, the A-O-DisCo and IDP modules individually enhance the generated video quality and the preservation of non-edited regions for both tasks. Additionally, they improve the appearance consistency (CLIP-I$_E$) of the edited regions for the swap task and the removal capability for the object removal task. (3) A further comparison between rows 3/4 and row 5 reveals that the combination of the A-O-DisCo and IDP modules leads to an even greater improvement in model performance.

## 5. Conclusion

In this work, we introduced O-DisCo-Edit, a unified framework designed to address the key challenges in controllable video editing. Our core innovation, O-DisCo, unifies various editing signals into a single, noise-based representation. This not only dramatically simplifies the training process and reduces resource demands but also enables multi-granularity editing during inference. Paired with the designed CFP module, O-DisCo-Edit can accomplish high-fidelity editing while robustly preserving unedited regions. Comprehensive experiments on eight different tasks show that O-DisCo-Edit achieves new SOTA results, outperforming both specialized and multi-task models. This success offers a new perspective on video editing research, that a single unified control signal can be both versatile and precise without sacrificing efficiency.

# References

[1] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiao-han Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024. 1, 3

[2] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025.

[3] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022.

[4] Yoav HaCohen, Nisan Chiprut, Benny Brazowski, Daniel Shalem, Dudu Moshe, Eitan Richardson, Eran Levin, Guy Shiran, Nir Zabari, Ori Gordon, et al. Ltx-video: Realtime video latent diffusion. *arXiv preprint arXiv:2501.00103*, 2024. 1

[5] Yuanpeng Tu, Hao Luo, Xi Chen, Sihui Ji, Xiang Bai, and Hengshuang Zhao. Videoanydoor: High-fidelity video object insertion with precise motion control. *arXiv preprint arXiv:2501.01427*, 2025. 2, 3

[6] Shiyuan Yang, Liang Hou, Haibin Huang, Chongyang Ma, Pengfei Wan, Di Zhang, Xiaodong Chen, and Jing Liao. Direct-a-video: Customized video generation with user-directed camera movement and object motion. In *Proceedings of ACM SIGGRAPH 2024 Conference Papers*, pages 1–12, 2024.

[7] Quanhao Li, Zhen Xing, Rui Wang, Hui Zhang, Qi Dai, and Zuxuan Wu. Magicmotion: Controllable video generation with dense-to-sparse trajectory guidance. *arXiv preprint arXiv:2503.16421*, 2025. 2

[8] Qi Zhao, Zhan Ma, and Pan Zhou. Dreaminsert: Zero-shot image-to-video object insertion from a single image. *arXiv preprint arXiv:2503.10342*, 2025. 2

[9] Guy Yariv, Yuval Kirstain, Amit Zohar, Shelly Sheynin, Yaniv Taigman, Yossi Adi, Sagie Benaim, and Adam Polyak. Through-the-mask: Mask-based motion trajectories for image-to-video generation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 18198–18208, 2025. 2

[10] Shengming Yin, Chenfei Wu, Jian Liang, Jie Shi, Houqiang Li, Gong Ming, and Nan Duan. Dragnuwa: Fine-grained control in video generation by integrating text. *Image, and Trajectory*, 2023. 2

[11] Ge Wang, Songlin Fan, Hangxu Liu, Quanjian Song, Hewei Wang, and Jinfeng Xu. Consistent video editing as flow-driven image-to-video generation. *arXiv preprint arXiv:2506.07713*, 2025.

[12] Chang Liu, Rui Li, Kaidong Zhang, Yunwei Lan, and Dong Liu. Stablev2v: Stablizing shape consistency in video-to-video editing. *arXiv preprint arXiv:2411.11045*, 2024. 2

[13] Zekai Gu, Rui Yan, Jiahao Lu, Peng Li, Zhiyang Dou, Chenyang Si, Zhen Dong, Qifeng Liu, Cheng Lin, Ziwei Liu, et al. Diffusion as shader: 3d-aware video diffu-

[14] Xianguo Wu, Zongbao Feng, Sai Yang, Yawei Qin, Hongyu Chen, and Yang Liu. Safety risk perception and control of water inrush during tunnel excavation in karst areas: An improved uncertain information fusion method. *Automation in Construction*, page 105421, 2024.

[15] Hanlin Wang, Hao Ouyang, Qiuyu Wang, Wen Wang, Ka Leong Cheng, Qifeng Chen, Yujun Shen, and Limin Wang. Levitor: 3d trajectory oriented image-to-video synthesis. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 12490–12500, 2025. 2

[16] Sen Liang, Zhentao Yu, Zhengguang Zhou, Teng Hu, Hongmei Wang, Yi Chen, Qin Lin, Yuan Zhou, Xin Li, Qinglin Lu, et al. Omniv2v: Versatile video generation and editing via dynamic content manipulation. *arXiv preprint arXiv:2506.01801*, 2025. 2, 3, 4, 12

[17] Ruineng Li, Daitao Xing, Huiming Sun, Yuanzhou Ha, Jinglin Shen, and Chiuman Ho. Tokenmotion: Decoupled motion control via token disentanglement for human-centric video generation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1951–1961, 2025.

[18] Xu Zhang, Hao Zhou, Haoming Qin, Xiaobin Lu, Jiaxing Yan, Guanzhong Wang, Zeyu Chen, and Yi Liu. Enabling versatile controls for video diffusion models. *arXiv preprint arXiv:2503.16983*, 2025.

[19] Zixuan Ye, Xuanhua He, Quande Liu, Qiulin Wang, Xintao Wang, Pengfei Wan, Di Zhang, Kun Gai, Qifeng Chen, and Wenhan Luo. Unic: Unified in-context video editing. *arXiv preprint arXiv:2506.04216*, 2025. 3, 12

[20] Zeyinzi Jiang, Zhen Han, Chaojie Mao, Jingfeng Zhang, Yulin Pan, and Yu Liu. Vace: All-in-one video creation and editing. *arXiv preprint arXiv:2503.07598*, 2025. 2, 3, 4, 5, 12, 13

[21] Geng Lin, Chen Gao, Jia-Bin Huang, Changil Kim, Yipeng Wang, Matthias Zwicker, and Ayush Saraf. Omnimat-terf: Robust omnimatte with 3d background modeling. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 23471–23480, 2023. 2, 5, 12

[22] Bojia Zi, Weixuan Peng, Xianbiao Qi, Jianan Wang, Shihao Zhao, Rong Xiao, and Kam-Fai Wong. Minimax-remover: Taming bad noise helps video object removal. *arXiv preprint arXiv:2505.24873*, 2025. 2, 5

[23] Qihua Chen, Yue Ma, Hongfa Wang, Junkun Yuan, Wenzhe Zhao, Qi Tian, Hongmei Wang, Shaobo Min, Qifeng Chen, and Wei Liu. Follow-your-canvas: Higher-resolution video outpainting with extensive content generation. *arXiv preprint arXiv:2409.01055*, 2024. 3

[24] Jianhong Bai, Menghan Xia, Xiao Fu, Xintao Wang, Lianrui Mu, Jinwen Cao, Zuozhu Liu, Haoji Hu, Xiang Bai, Pengfei Wan, et al. Recammaster: Camera-controlled generative rendering from a single video. *arXiv preprint arXiv:2503.11647*, 2025. 3

[25] Bojia Zi, Penghui Ruan, Marco Chen, Xianbiao Qi, Shaozhe Hao, Shihao Zhao, Youze Huang, Bin Liang, Rong Xiao, and Kam-Fai Wong. Se\˜ norita-2m: A high-quality instruction-

based dataset for general video editing by video specialists. *arXiv preprint arXiv:2502.06734*, 2025. 3, 5, 13

[26] Changgu Chen, Xiaoyan Yang, Junwei Shu, Changbo Wang, and Yang Li. Lmp: Leveraging motion prior in zero-shot video generation with diffusion transformer. *arXiv preprint arXiv:2505.14167*, 2025. 3

[27] Xinyu Zhang, Zicheng Duan, Dong Gong, and Lingqiao Liu. Training-free motion-guided video generation with enhanced temporal consistency using motion consistency loss. *arXiv preprint arXiv:2501.07563*, 2025. 3

[28] Yuxuan Bian, Zhaoyang Zhang, Xuan Ju, Mingdeng Cao, Liangbin Xie, Ying Shan, and Qiang Xu. Videopainter: Any-length video inpainting and editing with plug-and-play context control. *arXiv preprint arXiv:2503.05639*, 2025. 4, 5, 12, 13

[29] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023. 5, 12

[30] Xiaowen Li, Haolan Xue, Peiran Ren, and Liefeng Bo. Diffueraser: A diffusion model for video inpainting. *arXiv preprint arXiv:2501.10018*, 2025. 5

[31] Shangchen Zhou, Chongyi Li, Kelvin CK Chan, and Chen Change Loy. Propainter: Improving propagation and transformer for video inpainting. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10477–10486, 2023. 5

[32] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alex Sorkine-Hornung, and Luc Van Gool. The 2017 davis challenge on video object segmentation. *arXiv preprint arXiv:1704.00675*, 2017. 5, 12, 15

[33] Qi Cai, Jingwen Chen, Yang Chen, Yehao Li, Fuchen Long, Yingwei Pan, Zhaofan Qiu, Yiheng Zhang, Fengbin Gao, Peihan Xu, et al. Hidream-i1: A high-efficient image generative foundation model with sparse diffusion transformer. *arXiv preprint arXiv:2505.22705*, 2025. 5, 12

[34] Chenfei Wu, Lun Huang, Qianxi Zhang, Binyang Li, Lei Ji, Fan Yang, Guillermo Sapiro, and Nan Duan. Godiva: Generating open-domain videos from natural descriptions. *arXiv preprint arXiv:2104.14806*, 2021. 5

[35] Shiyi Zhang, Junhao Zhuang, Zhaoyang Zhang, Ying Shan, and Yansong Tang. Flexiact: Towards flexible action control in heterogeneous scenarios. *arXiv preprint arXiv:2505.03730*, 2025. 5

[36] Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges. *arXiv preprint arXiv:1812.01717*, 2018. 5

[37] Matthias Wright and Björn Ommer. Artfid: Quantitative evaluation of neural style transfer. In *Proceedings of DAGM German Conference on Pattern Recognition*, pages 560–576. Springer, 2022. 5, 12

[38] Weifeng Chen, Jie Wu, Pan Xie, Hefeng Wu, Jiashi Li, Xin Xia, Xuefeng Xiao, and Liang Lin. Control-a-video: Controllable text-to-video generation with diffusion models. *CoRR*, 2023. 5

[39] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21807–21818, 2024. 5

[40] Jiwoo Chung, Sangeek Hyun, and Jae-Pil Heo. Style injection in diffusion: A training-free approach for adapting large-scale diffusion models for style transfer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8795–8805, 2024. 5, 13

[41] Shaoteng Liu, Tianyu Wang, Jui-Hsien Wang, Qing Liu, Zhifei Zhang, Joon-Young Lee, Yijun Li, Bei Yu, Zhe Lin, Soo Ye Kim, et al. Generative video propagation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 17712–17722, 2025. 12

[42] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 13

[43] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 13

[44] Yang Ye, Xianyi He, Zongjian Li, Bin Lin, Shenghai Yuan, Zhiyuan Yan, Bohan Hou, and Li Yuan. Imgedit: A unified image editing dataset and benchmark. *arXiv preprint arXiv:2505.20275*, 2025. 13

[45] Xiang Wang, Hangjie Yuan, Shiwei Zhang, Dayou Chen, Jiuniu Wang, Yingya Zhang, Yujun Shen, Deli Zhao, and Jingren Zhou. Videocomposer: Compositional video synthesis with motion controllability. *Advances in Neural Information Processing Systems*, pages 7594–7611, 2023. 15

# Appendix

## A. More Details about Methodology

### A.1. Random Object Distortion Control

During the training phase, we employ a random distorter to generate the random object distortion control (R-O-DisCo), as detailed in Algorithm 1. The algorithm takes the reference video $V_{\text{ref}}$ and mask $M$ as input, and outputs the R-O-DisCo, $V_{\text{RODC}}$. The process begins by randomly sampling several distortion parameters from uniform distributions: a scaling factor ($\theta$) from $[1.5, 3.0]$, a scaling mode ($\mu$) to determine whether to scale up or down, a target channel ($c^*$) is the scaled channel, a color offset ($\delta$) from $\{-100, -50, 50, 100\}$, and a block size ($b$) for mosaicking from $\{8, 10, 12, 15, 16, 20, 24\}$. The functions used in this process are defined as follows: ZerosLike() returns a tensor of zeros with the same shape as the input video; Clip() constrains the input values to the range $[0, 255]$; AvgPool() performs downsampling using average pooling; and Interpolate() performs upsampling using nearest-neighbor interpolation. Based on these, the distorter applies both color distortion and mosaicking to the reference video. The final $V_{\text{RODC}}$, is then obtained by using the mask $M$ to apply these distortions exclusively to the object that needs to be edited.

### A.2. Adaptive Object Distortion Control

During inference phase, our model adapts to specific tasks or instructions via adaptive object distortion control (A-O-DisCo), which is implemented by an adaptive distorter. It is achieved through contrast modification (scaling and clipping) and dynamic noise injection within the editable regions. The adaptive controller determines suitable values for contrast parameter $\alpha$, noise intensity $\sigma$, and normalized Gaussian kernel size $k$ by calculating two similarities: $\mathbf{Sim}_i$, the edited region's edge map similarity between the reference image and the reference video's first frame; $\mathbf{Sim}_v$, the intra-frame similarity within the reference video's edited region edge map. Empirically, fitting these three parameters using a quadratic polynomial of two similarity yields superior results.

In the Algorithm 2, $V_{\text{ref}}$ is the reference video, $M$ represents the mask, $I_{\text{ref}}$ is the reference image, and $T = 49$ is the number of video frames. The functions are defined as follows: Canny() computes the edge map; SSIM$(, | f_{\text{masks}})$ calculates the SSIM value on the edited region, which is extracted using the $f_{\text{masks}}$; Odd() returns the nearest odd integer value; and GaussianBlur() performs a Gaussian blur. The parameters are defined by the following quadratic functions: $f_1 = 3000 \cdot \mathbf{Sim}_i^2 + 6000 \cdot \mathbf{Sim}_i + 300$, $f_2 = 4622.64 \cdot \mathbf{Sim}_v^2 + 92453.28 \cdot \mathbf{Sim}_v + 4623.64$, and $f_3 = -36 \cdot \mathbf{Sim}_v^2 + 72 \cdot \mathbf{Sim}_v - 35$. The process involves two passes: in the first, we iterate through each video

---

**Algorithm 1** Random Distorter

**Input**: $V_{\text{ref}}$, $M$
**Output**: $V_{\text{RODC}}$

1: Initialize processing parameters:
2:     Scaling factor: $\theta \sim U(1.5, 3.0)$
3:     Target channel: $c^* \sim U(\{0, 1, 2\})$
4:     Color offset: $\delta \sim U(\{-100, -50, 50, 100\})$
5:     Block size: $b \sim U(\{8, 10, 12, 15, 16, 20, 24\})$
6:     Scaling mode: $\mu \sim U(\{0, 1\})$
7: Initialize $V_{\text{cd}} \leftarrow \text{ZerosLike}(V_{\text{ref}})$
8: Color Distortion
9: **for** each channel $c \in \{0, 1, 2\}$ **do**
10:     **if** $c = c^*$ **then**
11:       **if** $\mu = 0$ **then**
12:         $V_{\text{cd}}[:, :, :, c] \leftarrow \text{Clip}(V_{\text{ref}}[:, :, :, c] \cdot \theta, 0, 255)$
13:       **else**
14:         $V_{\text{cd}}[:, :, :, c] \leftarrow \text{Clip}(V_{\text{ref}}[:, :, :, c]/\theta, 0, 255)$
15:       **end if**
16:     **else if** $c = 0$ **then**
17:       $V_{\text{cd}}[:, :, :, c] \leftarrow \text{Clip}(V_{\text{ref}}[:, :, :, c] + \delta, 0, 255)$
18:     **else**
19:       $V_{\text{cd}}[:, :, :, c] \leftarrow \text{Clip}(V_{\text{ref}}[:, :, :, c] - \delta, 0, 255)$
20:     **end if**
21: **end for**
22: Mosaicking
23: $V_{\text{low}} \leftarrow \text{AvgPool}(V_{\text{cd}}, b)$
24: $V_{\text{cdm}} \leftarrow \text{Interpolate}(V_{\text{low}}, b)$
25: $V_{\text{RODC}} \leftarrow V_{\text{ref}} \odot (\mathbf{1} - M) + V_{\text{cdm}} \odot M$
26: **return** $V_{\text{RODC}}$

---

frame, calculate the similarity metrics (using SSIM), and determine the adaptive parameters; in the second, we apply scaling & clipping and Gaussian blurring to the reference video frames. Finally, using a mask, we obtain the A-O-DisCo video, $V_{\text{AODC}}$.

### A.3. "Copy-Form" Preservation Module

We propose the copy-form preservation (CFP) module, which achieves preservation of non-edited regions with an effect similar to "first-frame copying". The CFP module is expressed as:

$$
\begin{aligned}
z_p^{v'} &= z_{\text{ref}}^v \odot (\mathbf{1} - z_{\text{mask}}[1 :]), \\
z_{\text{images}} &= [z_{\text{ref}}^i, z_p^{v'}],
\end{aligned}
\tag{5}
$$

During training, we randomly dilate the masks $z_{\text{mask}}[1 \ :]$ using a max pooling operation, where the kernel size is uniformly sampled from the set $\{1, 3, 5, 7, 9, 11, 13, 15, 17, 19, 21\}$ with equal probability. During inference, users can adaptively adjust the dilation kernel size based on the target object's size in the reference image. This allows for fine-tuning the edited region in the

**Algorithm 2** Adaptive Distorter

**Input**: $V_{\text{ref}}$, $M$, $I_{\text{ref}}$
**Output**: $V_{\text{AODC}}$

1: **First Pass: Calculate SSIM Metrics**
2: **for** $t \leftarrow 0$ **to** $T-1$ **do**
3:      $f_{\text{ref}} \leftarrow V_{\text{ref}}[t]$
4:      $f_{\text{masks}} \leftarrow M[t]$
5:      **if** $t = 0$ **then**
6:          $f_{\text{ref}}^{\text{canny}} \leftarrow \text{Canny}(f_{\text{ref}})$
7:          $I_{\text{ref}}^{\text{canny}} \leftarrow \text{Canny}(I_{\text{ref}})$
8:          $f_{\text{old}} \leftarrow f_{\text{ref}}^{\text{canny}}$
9:          $\mathbf{Sim}_i \leftarrow \text{SSIM}(f_{\text{ref}}^{\text{canny}}, I_{\text{ref}}^{\text{canny}} | f_{\text{masks}})$
10:      **else**
11:          $f_{\text{ref}}^{\text{canny}} \leftarrow \text{Canny}(f_{\text{ref}})$
12:          $f_{\text{new}} \leftarrow f_{\text{ref}}^{\text{canny}}$
13:          $ssim_{\text{v}}[t] \leftarrow \text{SSIM}(f_{\text{new}}, f_{\text{old}} | f_{\text{masks}})$
14:          $f_{\text{old}} \leftarrow f_{\text{ref}}^{\text{canny}}$
15:      **end if**
16: **end for**
17: Compute average SSIM of $ssim_{\text{v}}[t]$, to obtain $\mathbf{Sim}_{\text{v}}$
18: **Calculate Adaptive Parameters**
19: $k \leftarrow 0.2 \cdot \text{Odd}(f_1(\mathbf{Sim}_i) + 1.2 \cdot f_2(\mathbf{Sim}_{\text{v}}))$
20: $\sigma \leftarrow 0.2 \cdot (f_1(\mathbf{Sim}_i) + 1.2 \cdot f_2(\mathbf{Sim}_{\text{v}}))$
21: $\alpha \leftarrow f_3(\mathbf{Sim}_{\text{v}})$
22: **Second Pass: Get A-O-DisCo**
23: Initialize $V_{\text{AODC}} \leftarrow \text{ZerosLike}(V_{\text{ref}})$
24: **for** $t \leftarrow 0$ **to** $T-1$ **do**
25:      $f_{\text{ref}} \leftarrow V_{\text{ref}}[t]$
26:      $f_{\text{masks}} \leftarrow M[t]$
27:      $f_{\text{AODC}} \leftarrow V_{\text{AODC}}[t]$
28:      Apply adaptive transformations:
29:      $f_{\text{c}} \leftarrow \text{Clip}(\alpha \cdot V_{\text{ref}}, 0, 255)$
30:      $f_{\text{cn}} \leftarrow \text{GaussianBlur}(f_{\text{ref}}, 0, 255, k, \sigma)$
31:      $f_{\text{AODC}} \leftarrow f_{\text{ref}} \odot (\mathbf{1} - M) + V_{\text{cn}} \odot M$
32: **end for**
33: **return** $V_{\text{AODC}}$

reference video to accommodate objects of different sizes for tasks such as swap or addition.

# B. More Details about Experiments

## B.1. Construction of the benchmark

Video editing is a rapidly advancing field. However, a public benchmark for first-frame-based video editing is still lacking. While the widely used DIVAS [32] is available, its lack of edited images makes it unsuitable for evaluating video editing models. The VACE benchmark [20], on the other hand, only openly provides 1-2 examples per task, which is insufficient for a comprehensive and reliable evaluation. The Genprop-Data benchmark [41] contains 40 video-image pairs, but its quality is compromised by videos with variable frame counts and dimensions. Most notably, the edited images in this benchmark exhibit significant background discrepancies compared to the original videos. Although models like OmniV2V [16] and UNIC [19] have constructed their own benchmarks, none of them are publicly available.

Given the scarcity of public video editing benchmarks, we construct a new multi-task video editing benchmark to comprehensively evaluate the performance of our method across various tasks, which includes tasks such as outpainting, object internal motion transfer, lighting transfer, color change, swap, addition, and style transfer. The data is sourced from DAVIS [32] and VPData [28]. The specific construction process is as follows:

• For outpainting: We select videos containing a variety of natural landscapes (mountains, lakes, oceans, skies) and man-made scenes (villages, buildings, roads) to ensure diversity.

• For object internal motion transfer: We choose videos with clear internal motion, such as water flowing inside a bottle, a rotating sphere, and the movement of a clock's hands. To specifically highlight the effect of internal motion transfer, we edit the first frame using relatively simple color and swap instructions.

• For lighting transfer: We select videos with significant lighting changes, including the play of light on water surfaces or metal, and lighting variations for both underwater and terrestrial object. Similarly, to focus on the lighting transfer effect, we edit the first frame with simple color and swap instructions.

• For other tasks: We randomly select several videos. We then use Qwen-VL2.5 7B [29] to randomly generate editing instructions. We use HiDream [33] and commercial models [1] to generate the edited first frame. From these, we select the best-edited first frames as reference images and their corresponding videos as reference videos. To increase the complexity of the swap and addition tasks, we further filter for videos that included obvious and complex motion.

In total, we collect a benchmark of 134 sets, each containing a reference video, a reference image, and a mask video. All of them have a resolution of $480 \times 720$ and video consist of 49 frames. Additionally, we evaluate the object removal task using the OmnimatteRF benchmark [21]. Since ground truth background videos for removal are often unavailable in real-world scenarios, we use a commercial model [1] to perform object removal on the first frame of the reference video. For a fairer comparison, this processed first frame is used as the input for multi-task baselines.

## B.2. Choice of Evaluation Metrics

For the color change, object internal motion transfer, and lighting transfer tasks, we choose ArtFID [37] as our

Table 6. A comparison between our method and specialized models on object removal task (49). The evaluation is based on several key metrics, including DBR (degree of object removal), BC (background consistency), VQ (video quality), and EC (editing completeness).

| Method | Diffuearser | Minimax | Propainter | O-DisCo-Edit |
|--------|-------------|---------|------------|--------------|
| DRO | 3.356 | 3.800 | 3.311 | **4.244** |
| BC | 2.889 | 3.333 | 2.778 | **3.689** |
| EC | 3.122 | 3.567 | 3.044 | **3.967** |
| VQ | 2.867 | 3.333 | 2.822 | **3.689** |



Figure 6. The bad cases of O-DisCo-Edit, where "T" denotes the temporal axis.

video quality metric. As the author of [40] notes, ArtFID is calculated as $(1 + \text{LPIPS}) \cdot (1 + \text{FID})$. Specifically, LPIPS [42] measures content fidelity between the generated video frames and the reference video frames, while FID [43] assesses style fidelity between the generated video frames and the corresponding edited first frame.

Just as described in the benchmark construction, the color change, object internal motion transfer, and lighting transfer tasks are analogous to local style transfer. ArtFID is a better metric for these tasks because it can tell the difference between content fidelity and style fidelity. This gives it a more detailed understanding of video quality than FVD, which only measures the overall video distribution. For this reason, ArtFID is especially well-suited for evaluating these specific tasks.

### B.3. User Study

We conduct a user study where anonymized generated videos are randomly distributed to participants for scoring on a 1-5 scale. For each task, we define specific evaluation criteria:

- For object removal: Users score the generated videos on degree of object removal (DOR), background consistency (BC) with the background video, and overall video quality (VQ).
- Outpainting: Users evaluate outpainting effectiveness (OE) and overall video quality (VQ).
- Object internal motion transfer: The criteria are the internal motion consistency (IMC) between the generated and reference videos, ID consistency (IDC) between the generated video and reference image, and overall video quality (VQ).
- Lighting transfer: Users score lighting consistency (LC) between the generated and reference videos, ID consistency (IDC) with the reference image, and overall video quality (VQ).
- Color change: User evaluate color consistency (ColorC) with the reference image, consistency of non-color content (COCC) with the reference video, and overall video quality (VQ).
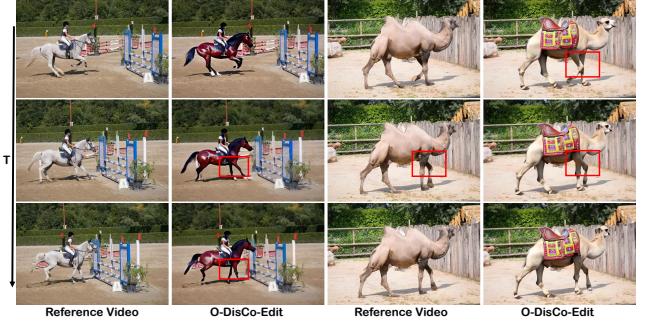- Swap: The criteria include ID consistency (IDC) with the

reference image, motion consistency (MC) with the reference video, and overall video quality (VQ).
- Addition: Users score ID consistency (IDC) with the reference image, the plausibility of the added object's motion (PAOM), the consistency of non-edited regions (CNER) with the reference video, and overall video quality (VQ).
- Style transfer: We define criteria for content consistency (CC) with the reference video, style consistency (SC) with the reference image, and overall video quality (VQ).

In our user study, we follow the methodology of [44] to ensure that editing effectiveness was prioritized. We do this by capping the video quality score at the editing effectiveness score for each submission, meaning the video's overall quality score cannot exceed its score for how well the edit is performed. For example, in the object removal task, the video quality score is capped by the degree of object removal. We then compute the overall editing completeness (EC) as the average of all metrics except for video quality. In total, we collecte 9 valid responses for the user study. The detailed results of this user study are presented in Tab. 6 and Tab. 7. Our method achieve the best user study results across all tasks except for color change. Notably, Senorita's top score in user studies comes from its first-frame propagation strategy. As discussed in the main submission, this strategy creates high visual consistency but performs poorly at preserving non-edited regions.

### B.4. Bad Case Analysis

As revealed by Tab. 3, our model receive its lowest scores on the swap task compared to its performance on other tasks. We attribute this poor performance to our O-DisCo-Edit model's difficulty in handling complex, four-limbed object motions, as illustrated by the misaligned legs of the swapped object in Fig. 6. We hypothesize this issue stems from three potential factors. First, the base model may have inherent limitations; as shown in Fig. 7, Senorita [25] and VideoPainter [28], which share our base model, exhibit similar misaligned leg problems. In contrast, while VACE 14B [20] shows some ID distortion, its legs maintain

Table 7. This table presents a comparison of user study results for different models. The metrics used for evaluation include: DBR (degree of object removal), BC (background consistency), VQ (video quality), EC (editing completeness), OE (outpainting effectiveness), IMC (internal motion consistency), IDC (ID consistency), LC (lighting consistency), ColorC (color consistency), COCC (consistency of non-color content), MC (motion consistency), PAOM (plausibility of the added object's motion), CNER (consistency of non-edited regions), CC (content consistency), and SC (style consistency).

| Task | Metric | VACE 1.3B | VACE 14B | Senorita | VideoPainter | O-DisCo-Edit |
|---|---|---|---|---|---|---|
| Object Removal (33) | DOR | 2.200 | 1.622 | 3.511 | 1.844 | **4.311** |
| | BC | 1.822 | 1.533 | 3.111 | 1.644 | **3.600** |
| | EC | 2.011 | 1.578 | 3.311 | 1.744 | **3.956** |
| | VQ | 1.911 | 1.444 | 3.156 | 1.578 | **3.689** |
| Outpainting | OE | 4.244 | 4.178 | 2.889 | 2.156 | **4.289** |
| | EC | 4.244 | 4.178 | 2.889 | 2.156 | **4.289** |
| | VQ | 3.933 | 3.911 | 2.600 | 1.978 | **4.067** |
| Object Internal Motion Transfer | IMC | 2.956 | 3.333 | 2.511 | 2.156 | **4.267** |
| | IDC | 3.067 | 2.911 | 3.533 | 2.444 | **4.089** |
| | EC | 3.011 | 3.122 | 3.022 | 2.300 | **4.178** |
| | VQ | 2.533 | 2.800 | 2.333 | 1.822 | **3.756** |
| Light Transfer | LC | 2.867 | 3.333 | 2.867 | 2.800 | **4.200** |
| | IDC | 3.267 | 3.489 | 3.200 | 3.022 | **3.756** |
| | EC | 3.067 | 3.411 | 3.033 | 2.911 | **3.978** |
| | VQ | 2.644 | 3.067 | 2.400 | 2.489 | **3.689** |
| Color Change | ColorC | 3.644 | 3.378 | **4.244** | 3.778 | 4.133 |
| | COCC | 3.622 | 3.533 | **3.822** | 3.489 | 3.756 |
| | EC | 3.633 | 3.456 | **4.033** | 3.633 | 3.944 |
| | VQ | 3.244 | 3.089 | **3.711** | 3.267 | 3.689 |
| Swap | MC | 3.711 | 3.911 | 3.533 | 3.489 | **3.956** |
| | IDC | 3.222 | 3.200 | 3.378 | 2.444 | **4.111** |
| | EC | 3.467 | 3.556 | 3.456 | 2.967 | **4.033** |
| | VQ | 2.956 | 3.089 | 2.956 | 2.178 | **3.689** |
| Addtion | PAOM | 3.289 | 3.200 | 3.222 | 3.178 | **4.111** |
| | IDC | 3.267 | 2.711 | 3.267 | 3.067 | **4.267** |
| | CNER | 3.556 | 3.689 | 4.000 | 3.067 | **3.733** |
| | EC | 3.370 | 3.200 | 3.496 | 3.104 | **4.037** |
| | VQ | 3.089 | 2.578 | 2.911 | 2.489 | **3.822** |
| Style Transfer | SC | \ | \ | 2.911 | \ | **4.356** |
| | CC | \ | \ | 3.067 | \ | **4.289** |
| | EC | \ | \ | 2.989 | \ | **4.322** |
| | VQ | \ | \ | 2.578 | \ | **4.156** |

correct depth and position. Second, the model's parameter count may be a factor, as VACE 1.3B (a smaller version of VACE 14B) also displays this failure mode. This suggests that increasing our model's parameter count could potentially resolve the issue. Finally, the training dataset we used contains a limited number of examples with complex four-limbed animal motions, which likely prevents the model from generalizing well to such challenging cases.

## B.5. More Ablation Experiment

The ablation study results for the outpainting task are shown in Tab. 8. (1) A comparison between row 1 and row 2 reveals that adding the CFP module significantly improves the model's ability to preserve non-edited regions, as evidenced by a substantial increase in both $PSNR_P$ and $SSIM_P$. (2) A comparison between rows 2 and 3 shows that adding the IDP module without A-O-DisCo results in a significant performance drop. This is because the IDP module emphasizes the position of the edited regions and the corre-
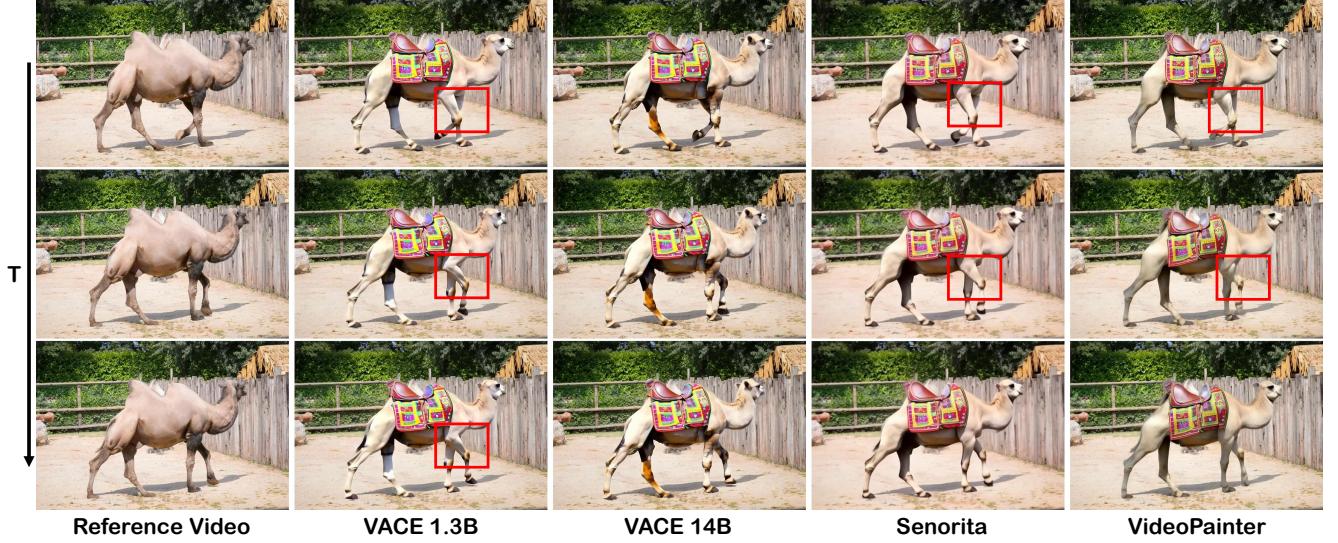
Figure 7. Qualitative results of the baseline models on swap task involving complex, four-limbed object motions. The temporal axis is denoted by 'T'.

Table 8. The results of the ablation experiment on outpainting task. ① CFP module, ② A-O-DisCo, ③ IDP module. "w/o ①②③" denotes training with R-O-DisCo and inference with a fixed signal, entirely omitting IDP and CFP modules. "w/o ②③" indicates training without module IDP and inference with a fixed signal. "w/o ②" refers to using a fixed inference signal. "w/o ③" indicates training without the IDP module. "TC" denotes temporal consistency.

| Metrics | Video Quality | | | Alignment | Preservation | |
| Model | TC↑ | FVD↓ | PSNR↑ | CLIP-T↑ | PSNR$_P$↑ | SSIM$_P$↑ |
|---|---|---|---|---|---|---|
| w/o ①②③ | **0.9983** | 621.6 | 20.11 | <u>12.50</u> | 26.66 | 0.8873 |
| w/o ②③ | 0.9976 | <u>80.89</u> | 26.03 | **12.27** | 33.75 | **0.9469** |
| w/o ② | <u>0.9979</u> | 222.7 | 25.70 | 12.20 | <u>33.85</u> | 0.9466 |
| w/o ③ | 0.9977 | 90.98 | <u>26.09</u> | 12.16 | 33.76 | **0.9469** |
| O-DisCo-Edit | 0.9978 | **77.03** | **26.43** | 12.19 | **33.87** | <u>0.9466</u> |

sponding content generation. Therefore, when A-O-DisCo is not introduced, the model is prone to retain the outpainting's boundary information, leading to residual boundary artifacts in the generated video, as illustrated in Fig. 8.

## B.6. More Examples

From Fig. 9 to Fig. 17, we present additional qualitative results for tasks including object removal, outpainting, object internal motion transfer, lighting transfer, color change, swap, addition, and style transfer. We extend the object removal task to the DAVIS dataset [32], as shown in Fig. 9 and Fig. 10. O-DisCo-Edit achieves excellent performance on the DAVIS dataset, demonstrating its effectiveness. Additionally, for the style transfer task, we incorporate Video-Composer [45] as a complementary baseline. The results
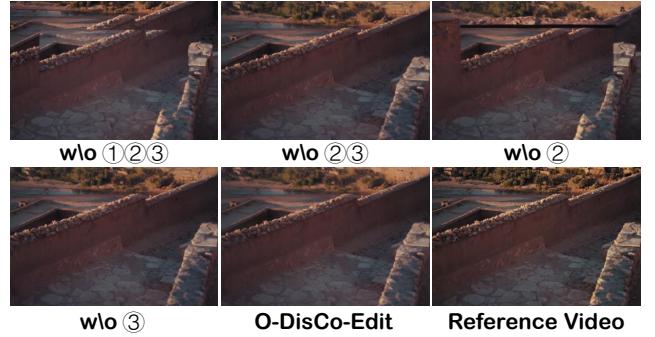


Figure 8. Ablation study results for the outpainting task. ① CFP module, ② A-O-DisCo, ③ IDP module. "w/o ①②③" denotes training with R-O-DisCo and inference with a fixed signal, entirely omitting IDP and CFP modules. "w/o ②③" indicates training without module IDP and inference with a fixed signal. "w/o ②" refers to using a fixed inference signal. "w/o ③" indicates training without the IDP module.

in Fig. 17 clearly show that O-DisCo-Edit surpasses Video-Composer in terms of style transfer performance.

## C. Limitation

Our work has several limitations. First, while we compare our model's performance with specialized models in the object removal task, we only compare it with a multi-task model in other tasks. This limited comparison prevents us from presenting a more comprehensive evaluation. Second, we conducted ablation studies for only three tasks, which does not provide a complete understanding of the contributions of different modules to other tasks. Third, the per-

formance of our model heavily depends on the quality of the first frame edit. Poor quality in the first frame edit may result in a significant drop in the model's performance.

## D. Ethical Consideration

From an ethical perspective, our model's powerful video editing capabilities provide creators with innovative tools that can inspire new ideas and enhance the artistic and creative aspects of video content. However, this also raises concerns about the potential for spreading misinformation and false content, which could undermine public trust in information. Additionally, these technologies may unintentionally reinforce existing biases and stereotypes, potentially influencing societal cultural perspectives in a negative way. These issues highlight the importance of ethical reflection and responsibility, urging policymakers, developers, and societal stakeholders to collaboratively establish appropriate regulations to ensure the healthy development of such technologies. In the future, we will make our model publicly available. This release will be accompanied by a licensing agreement that requires users to adhere to our guidelines and outlines acceptable use cases, thereby limiting potential abuse by malicious users.

Figure 9. Visual comparison of object removal capabilities, benchmarking O-DisCo-Edit against specialized baselines on the DAVIS dataset. "T" indicates the temporal axis.



Figure 10. Qualitative results for object removal, demonstrating the performance of O-DisCo-Edit versus specialized baselines. The temporal progression is shown along axis "T".
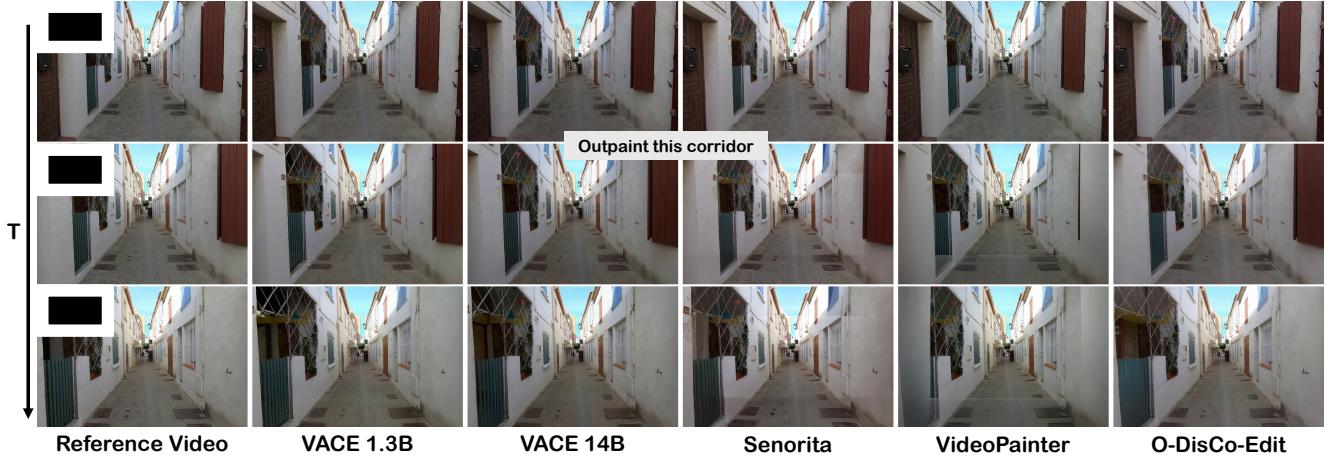
Figure 11. Performance visualization for the outpainting task, illustrating how O-DisCo-Edit compares with multi-task baselines on our benchmark. The temporal axis is denoted by "T".



Figure 12. A comparative evaluation of object internal transfer performance between O-DisCo-Edit and mutil-task baselines on the our benchmark, where "T" denotes the temporal axis.



Figure 13. A comparison of lighting transfer performance between O-DisCo-Edit and multi-task baselines on our benchmark. The temporal progression is marked by "T".

Figure 14. A comparative view of color change performance between O-DisCo-Edit and multi-task baselines on our benchmark. The time axis is indicated by 'T'.
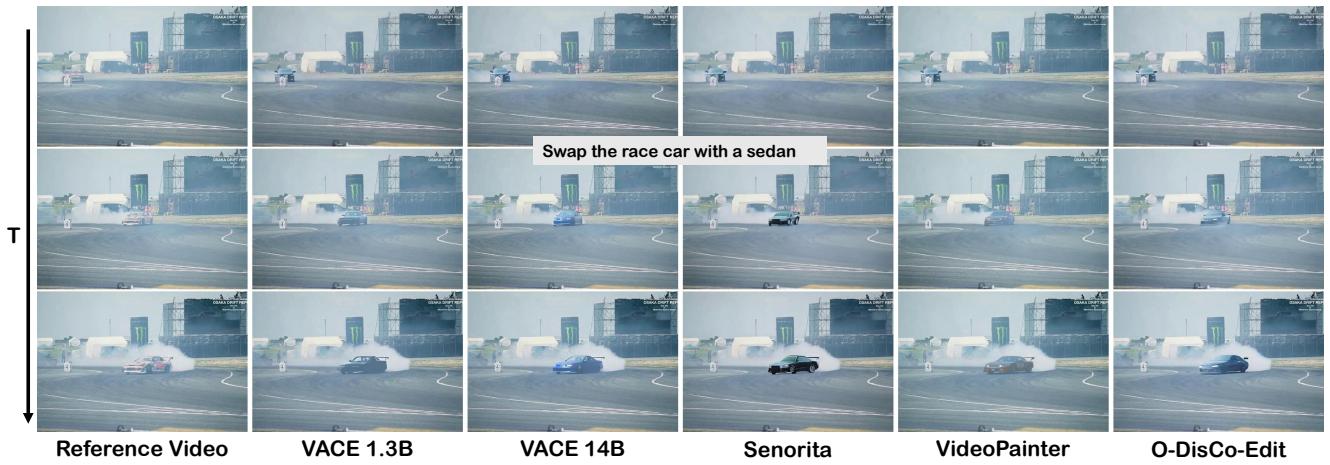


Figure 15. Qualitative results for the object swap task, comparing the output of O-DisCo-Edit against multi-task baselines on our benchmark. "T" denotes the temporal axis.
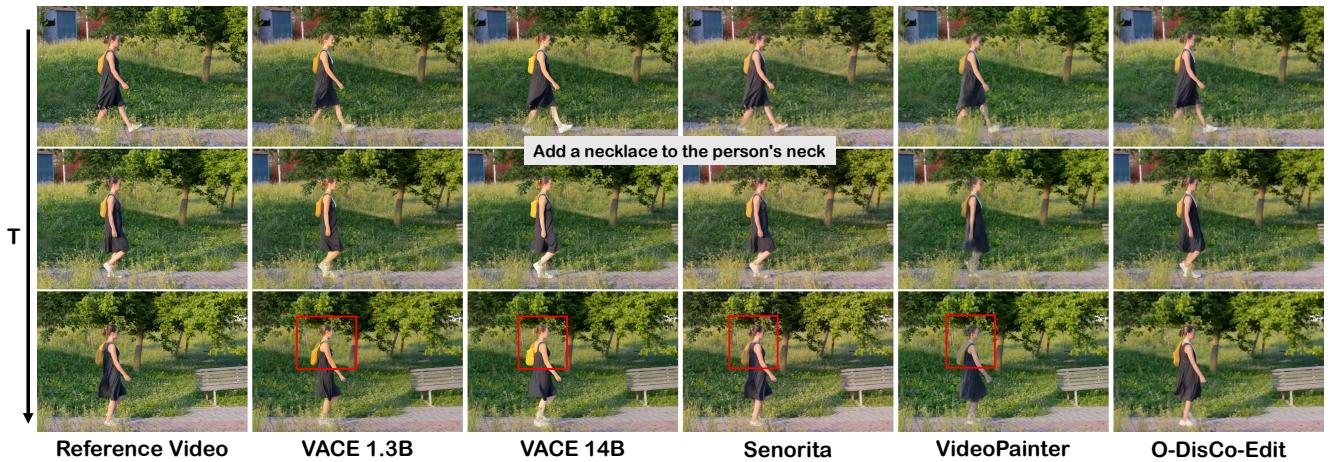


Figure 16. An evaluation of addtion performance, highlighting the differences between O-DisCo-Edit and multi-task baselines on our benchmark. The temporal axis is shown as "T".
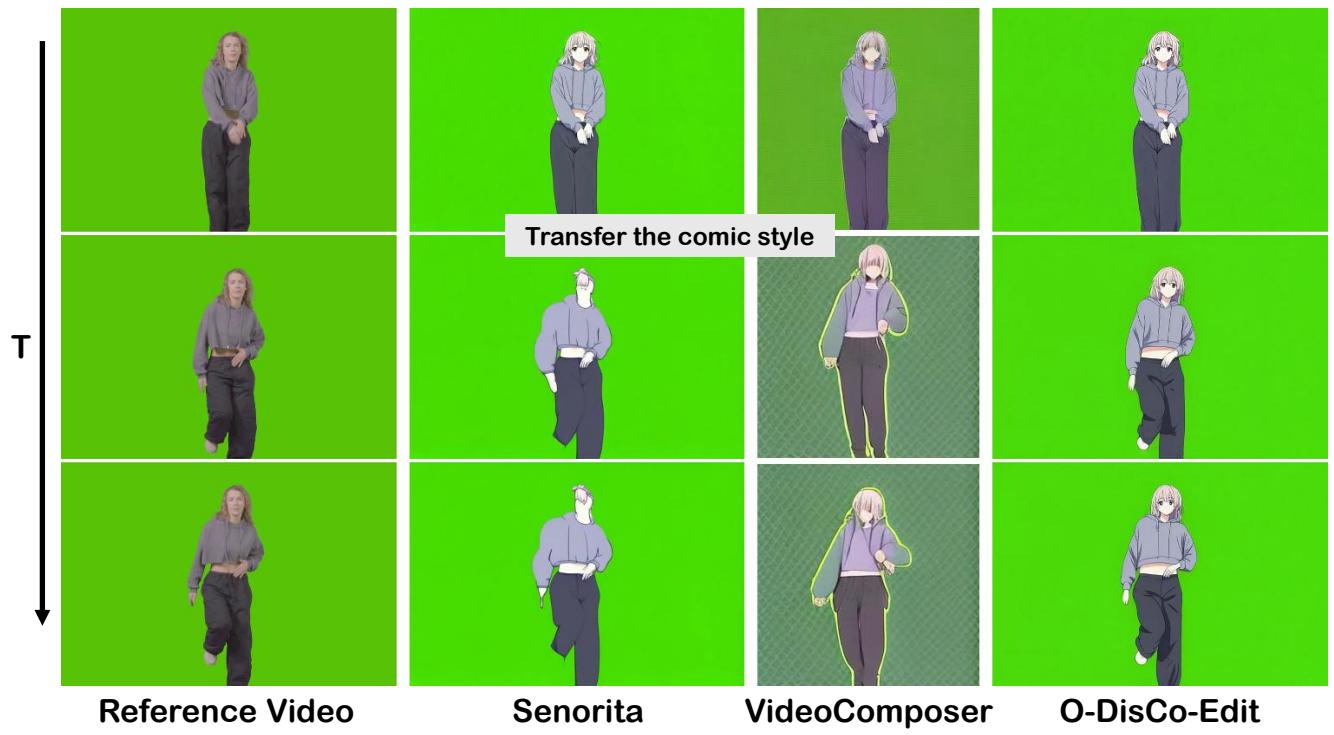
Figure 17. A side-by-side comparison of style transfer performance, showcasing O-DisCo-Edit against multi-task baselines on our benchmark. The letter "T" denotes the temporal dimension.