# Tuning-free Visual Effect Transfer across Videos

Maxwell Jones[1]    Rameen Abdal[2]    Or Patashnik[2]

Ruslan Salakhutdinov[1]    Sergey Tulyakov[2]    Jun-Yan Zhu[1]    Kuan-Chieh Jackson Wang[2]

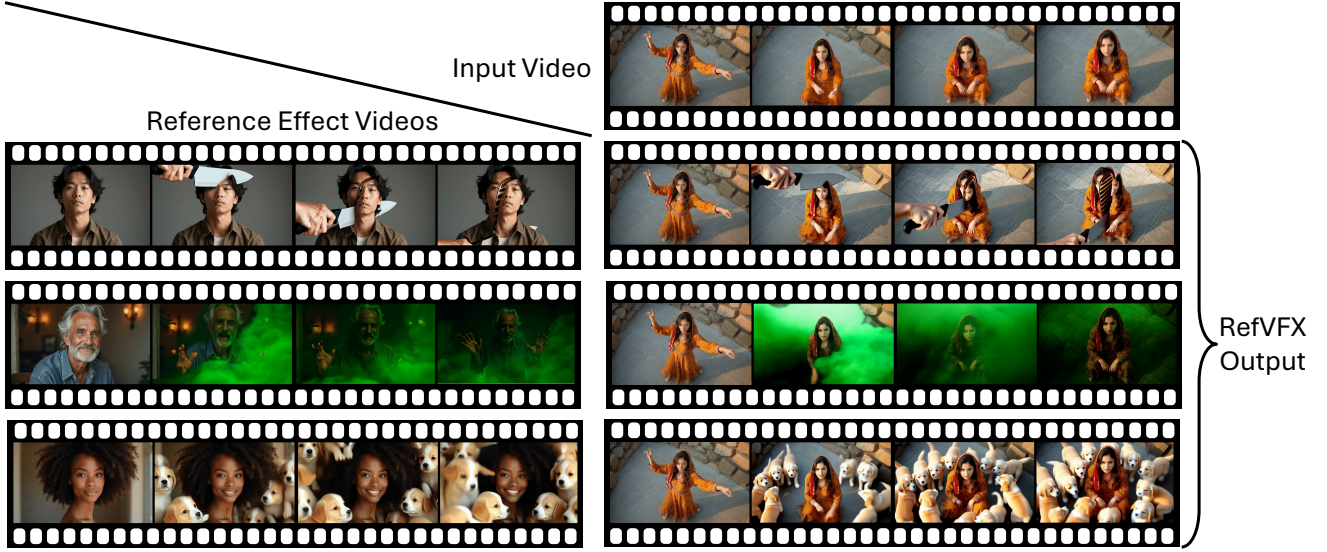[1]Carnegie Mellon University    [2]Snap Research

Figure 1. **Overview**. We present `RefVFX`, a tuning-free framework for visual effect transfer across videos. Given a reference effect video and an input video, our method produces a new output video where the reference's temporal effect is seamlessly applied to the input's content and motion. Unlike prompt-based or keyframe-conditioned approaches, `RefVFX` directly learns to interpret and reproduce complex time-varying visual effects such as material transformations, object additions, or atmospheric effects from example videos during training, enabling faithful and coherent visual effect transfer at inference time.

## Abstract

We present `RefVFX`, a new framework that transfers complex temporal effects from a reference video onto a target video or image in a feed-forward manner. While existing methods excel at prompt-based or keyframe-conditioned editing, they struggle with dynamic temporal effects such as dynamic lighting changes or character transformations, which are difficult to describe via text or static conditions. Transferring a video effect is challenging, as the model must integrate the new temporal dynamics with the input video's existing motion and appearance. To address this, we introduce a large-scale dataset of triplets, where each triplet consists of a reference effect video, an input image or video, and a corresponding output video depicting the transferred effect. Creating this data is non-trivial, espe-cially the video-to-video effect triplets, which do not exist naturally. To generate these, we propose a scalable auto-mated pipeline that creates high-quality paired videos de-signed to preserve the input's motion and structure while transforming it based on some fixed, repeatable effect. We then augment this data with image-to-video effects derived from LoRA adapters and code-based temporal effects gen-erated through programmatic composition. Building on our new dataset, we train our reference-conditioned model us-ing recent text-to-video backbones. Experimental results demonstrate that `RefVFX` produces visually consistent and temporally coherent edits, generalizes across unseen effect categories, and outperforms prompt-only baselines in both quantitative metrics and human preference. See our website at this URL.

# 1. Introduction

Generative models have enabled the automatic synthesis and manipulation of images and videos with remarkable realism and diversity [5, 17, 22, 30, 38, 43, 57, 58, 62, 74]. Recent advances in video generation have expanded user control, allowing video editing via text prompts, keyframes, or depth maps [36, 74, 88]. However, most models focus on semantic edits such as modifying objects [58, 89], scenes [58, 89], or styles [49, 83, 89], while overlooking *"temporal effects"* (see Figure 1): effects that *unfold over time* and define the emotional, stylistic, and cinematic character of a video. Effects such as dynamic lighting, intricate camera movements, or character transformations remain difficult to express through text or current visual inputs.

Reference-based conditioning offers a natural way to specify complex temporal behaviors, letting users convey the desired effect through a reference video that captures subtle cues such as motion rhythm, lighting changes, or stylistic transitions. Transferring such temporal effects, especially from one video to another, i.e., a Ref. Video + Input Video → Output Video setup, is especially challenging. To the best of our knowledge, we are the first to show results on this task.

While reference-based editing has shown promise for images [40, 65, 82], extending it to videos poses unique challenges. First, building a dataset is non-trivial: image-to-video (I2V) effects require consistent triplets in which a single static input produces multiple effect-rich outputs, and video-to-video (V2V) effects demand motion-consistent transformations that change only the temporal effect. Without such data, existing systems may fail to learn how to transfer temporal effects independently of the reference video's content or motion. Second, the model must extract temporal dynamics from the reference and seamlessly integrate them with the target's motion and appearance to produce coherent, high quality results. This requires disentangling the effect itself from the specific content and motion of the reference.

We introduce a reference-based video generation dataset and a training method dubbed `RefVFX` that transfers complex temporal effects from a reference video onto a target image or video in a feed-forward manner. Built upon recent text-to-video diffusion backbones [22, 38, 74], our model jointly conditions on three inputs: a reference video that provides the temporal effect, an input image or video that defines the scene content and motion, and a text prompt that offers high-level semantic guidance. Through this conditioning, `RefVFX` learns to integrate the reference's temporal dynamics harmoniously with the input's appearance and motion, producing coherent and visually consistent results.

A key contribution of our work is a large-scale dataset designed specifically for this task. Each sample consists of a triplet: a reference video depicting a temporal effect, an input image or video, and a corresponding output video showing the effect applied to the input. Building such a dataset at scale is challenging, as it requires coherent alignment between the reference effect and the target's content and motion. To this end, we construct a unified data curation pipeline that combines three complementary sources: (1) image-to-video effect data derived from existing LoRA-based models, (2) video-to-video transformations generated through an automated pipeline, and (3) synthetic, code-based temporal effects. Together, these sources yield over 1,700 unique effects and more than 120K triplets, providing unprecedented diversity for training reference-conditioned video editors.

We evaluate `RefVFX` across diverse unseen temporal effects, comparing it with recent prompt-based and reference-guided video editing models [26, 36, 74]. Qualitatively, our method produces coherent, visually consistent videos that successfully integrate the transferred effects while preserving the motion and appearance of the input. Quantitatively, `RefVFX` achieves higher scores in reference video embedding similarity to the reference effect video; however, existing metrics only partially reflect the aesthetic and temporal nuances of effect transfer. To better assess perceptual quality, we conduct a human preference study, where participants consistently favor our results, highlighting the importance of human judgment for evaluating dynamic visual effects. Notably, the model generalizes to unseen categories of effects and operates in a feed-forward manner without optimization at inference time, demonstrating both robustness and efficiency. Our code will be released upon publication.

Our main contributions are summarized as follows:

- We introduce `RefVFX`, a framework for reference-based video effect transfer that enables users to apply complex temporal effects from a reference video onto arbitrary videos or images in a feed-forward manner.
- We construct a large-scale, effect-aligned dataset comprising over 120K triplets of (reference, input, output) videos covering more than 1,700 distinct temporal effects along with a validation set, establishing a new benchmark for future research.
- We design a multi-source conditioning architecture built upon recent diffusion backbones that jointly encode reference video dynamics, input appearance/motion, and text prompts.
- We conduct extensive qualitative, quantitative, and human preference evaluations, demonstrating that our method achieves superior visual and temporal coherence compared to prompt-only or reference-guided baselines, and generalizes effectively to unseen effect categories.

## 2. Related Work

**Text-to-Video Generation.** Recent advances in text-to-video generation have significantly improved video quality, temporal coherence, and prompt following [5, 14, 22, 28–30, 38, 53, 58, 74]. Similar to text-to-image generation [17, 43, 57, 62], these models are largely based on diffusion or flow-matching in latent space [4, 16, 48, 52, 69]. Early work adopted U-Net backbones [20, 63, 67], while modern approaches employ diffusion transformers [56] that embed videos into latent patches and apply bidirectional attention with text [22, 28, 30, 38, 58, 74]. Beyond pure text conditioning, some models perform image-to-video or first–last-frame generation [22, 74], where a user provides keyframes and a text prompt to guide synthesis. However, these settings modify only edge frames, whereas our approach conditions on both an entire input video and a reference effect video, enabling transfer of *temporal effects* that evolve throughout the sequence while preserving input motion and appearance.

**Reference-Based Controllable Generation.** Reference-based conditioning uses a reference image or video that is not spatially aligned with the output. In text-to-image models, it is widely used for identity preservation [12, 21, 42, 46, 54, 55, 65, 70, 77, 79, 82] or style transfer [18, 25, 32, 64, 68, 80, 85]. These signals are injected via per-reference optimization [18, 65, 68], cross-attention [21, 46, 54, 70, 80, 82], or added reference tokens [25, 42, 64]. Video generative models extend these paradigms temporally [11, 23, 33, 47, 49, 50, 75, 83, 84], typically conditioning on one or more *static* reference images (identity or style) alongside text. Such methods maintain consistent appearance but cannot model *temporally evolving* effects. More recent finetuning based methods such as Dynamic Concepts [2, 3] are also inefficient, requiring either a new LoRA per effect instance or trained on limited data hampering scalability. In contrast, RefVFX conditions on a *reference video* that encodes dynamic, time-varying effects, allowing generalizable transfer of phenomena such as lighting changes, camera motion, stylistic transitions or identity-based conditioning.

**Video Editing.** Video editing has progressed through both zero-shot and supervised methods [7–9, 15, 24, 39, 44, 73, 76, 78, 86]. Pretrained text-to-image models can perform zero-shot edits [7, 9, 15, 24, 39, 73], while paired datasets of object or style edits [6, 76] have enabled general-purpose editing models [8, 44, 78]. Video editing methods follow similar patterns: some are zero-shot [10, 19], while others train paired text-driven editors [13, 51, 89]. However, existing methods rely solely on text and an input video, enabling semantic or appearance edits but offering no mechanism to control *how* visual properties evolve over time. Our approach introduces an additional reference video encoding a temporally varying effect, enabling control over dynamic

phenomena such as lighting transitions, motion-dependent effects, and stylistic evolution that static, per-frame, or text-only paradigms cannot capture.

## 3. Method

### 3.1. Overview

Our goal is to transfer a temporally evolving visual effect from a *reference video* onto a separate *input image or video*. The resulting output should preserve the content and motion of the input while exhibiting the temporal dynamics shown in the reference. To achieve this, we introduce RefVFX, which combines (1) a large-scale effect–aligned triplet dataset and (2) a diffusion-based video generation model conditioned jointly on the reference video, input video, and text.

First, we give a high level overview of our dataset, followed by describing how we create each individual subset of the dataset in detail. Following this, we describe the model architecture for RefVFX and implementation details for training.

### 3.2. Training Data

Our training data consists of triplets of the form (reference video, input image or video, target video). To create these triplets, we first curate $N$ effects $\{E^i\}_{i=1}^N$ (e.g., rain, motion blur, color shift), with each effect $E^i$ associated with a set $S^i$ of $K$ input output pairs: $S^i = \{(\text{input}_j^i, \text{output}_j^i)\}_{j=1}^K$. Here, $\text{output}_j^i$ is obtained by applying $E^i$ to $\text{input}_j^i$. To form a training triplet, we select an effect $E^i$ and two distinct pairs $(\text{input}_j^i, \text{output}_j^i)$ and $(\text{input}_\ell^i, \text{output}_\ell^i)$. We discard the input of the first pair, and construct triplet (reference = $\text{output}_j^i$, input = $\text{input}_\ell^i$, target = $\text{output}_\ell^i$), where the model takes in the first two videos and must predict the third. With this objective, the model learns to transfer the temporal effect from the reference onto the input while preserving the input videos motion and content. Triplets are derived from three complementary sources: (1) LoRA-based image-to-video effects, (2) paired video-to-video effects generated through our scalable pipeline, and (3) programmatically defined temporal effects, which we describe below and showcase in Fig. 2.

#### 3.2.1. Image to Video Data:

We curate a set of image-to-video effects using open-source Low-Rank Adapters (LoRAs) [31] trained on small sets of videos sharing a common visual effect atop a base image-to-video model [74]. Each adapter is treated as a distinct effect $E_i$. For every effect, we generate a large collection of videos from diverse input images synthesized via high-quality image generation models [44, 78], where each per-adapter set forms one effect category. Triplets are composed

(a) LoRA I2V Effects     (b) Custom V2V Effects     (c) Code Based Effects

Figure 2. **Dataset Examples**. We show example triplets from each of our datasets. (a) We curate a reference video + image to video dataset by collecting Low Rank Adapters (LoRAs) [31] for different Image to Video effects online. For each effect, we can apply its corresponding LoRA to two separate images to create a triplet. (b) We create a custom pipeline for generating text guided reference video + input video to video effects. For more details, see Figure 4. (c) We generate a large scale set of (ref video, input video, output video) triplets by curating specific code pipelines that apply specific, detailed effects to arbitrary videos. Armed with a specific code based effect and a fixed set of hyperparameters, we can apply the exact effect to an arbitrary number of input videos to create many triplets.
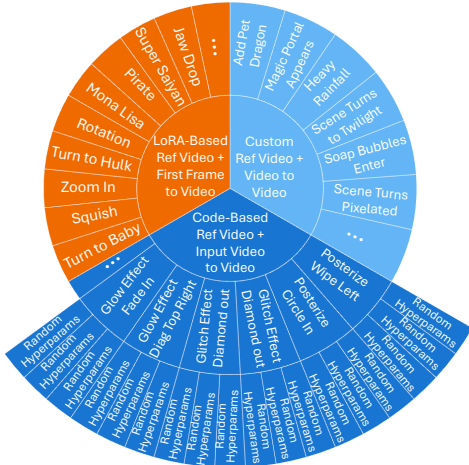


Figure 3. **Dataset Subset**. We display a summary of our dataset structure and available sample effects. These include LoRA-based Image-to-Video, our scalable V2V pipeline (see Fig. 4), and programmatic methods. For code-based effects, individual tasks are created by combining specific effects (e.g., Glow) and transitions (e.g., Fade In) with randomized hyperparameters.

of a reference video with effect $E_i$, the first frame of the target, and the corresponding output video. In total, we source 43 LoRAs with an average of 300 videos per effect, yielding over 14K video clips. See Fig. 3 for example LoRAs.

### 3.2.2. Video to Video Data:

We propose a scalable video-to-video data generation pipeline for a broad class of *motion-based effects*. These effects manifest *over time* rather than being the product of a single frame image edit followed by propagation of the edit to all frames (e.g., "turn the person into metal" as compared to "add a hat to the person throughout"). To our knowledge, this is the first scalable method designed for creating a large-scale training data pipeline for motion-based video editing.

Given an arbitrary effect prompt, our method produces paired videos $(V, V')$, consisting of an original video $V$ and its edited counterpart $V'$ exhibiting effect $E_i$ (Fig. 4, 3). **Scalable V2V Algorithm.** The process proceeds for generating these $(V, V')$ pairs is as follows: We begin by generating a subject image $I$ using a high-quality image generation model [43, 78]. Using an image editing model [44, 78], we modify $I$ to produce an image $I'$ that depicts the same subject under a new pose, camera angle, and facial expression, which will be used to anchor the last frame of our video pair. To obtain the last frame of the edited video with effect $E_i$, we apply an image editing model on $I'$ [44, 78] to apply effect $E_i$, resulting in $I''$. Given the unedited image $I$ and the edited image $I'$, we synthesize video $V$ using a first–last frame interpolation model [74]. Once we have $V$, we extract intermediate poses using an image-to-pose model [81], which will be used as conditioning for our second video generation. Finally, we generate the edited video $V'$ by leveraging a conditional video generation model [36], which can take in independent conditioning signals for each output video frame generated. We take advantage of this property, and condition the video generation on:

- Initial image $I$ for first frame conditioning
- Final edited image $I''$ for last frame conditioning
- Intermediate pose images extracted from $V$ for frames 2 through $N-1$

This enables framewise consistency between the input and output video as well as realistic temporal transitions between the unedited initial frame $I$ and stylized final frame $I''$. For the a visual, see Figure 4.

**Effect Generation.** To ensure diversity, we automatically generate a large set of motion-based effect prompts using GPT-4o [35]. These effects span multiple semantic and visual categories, which we detail in the supplement.
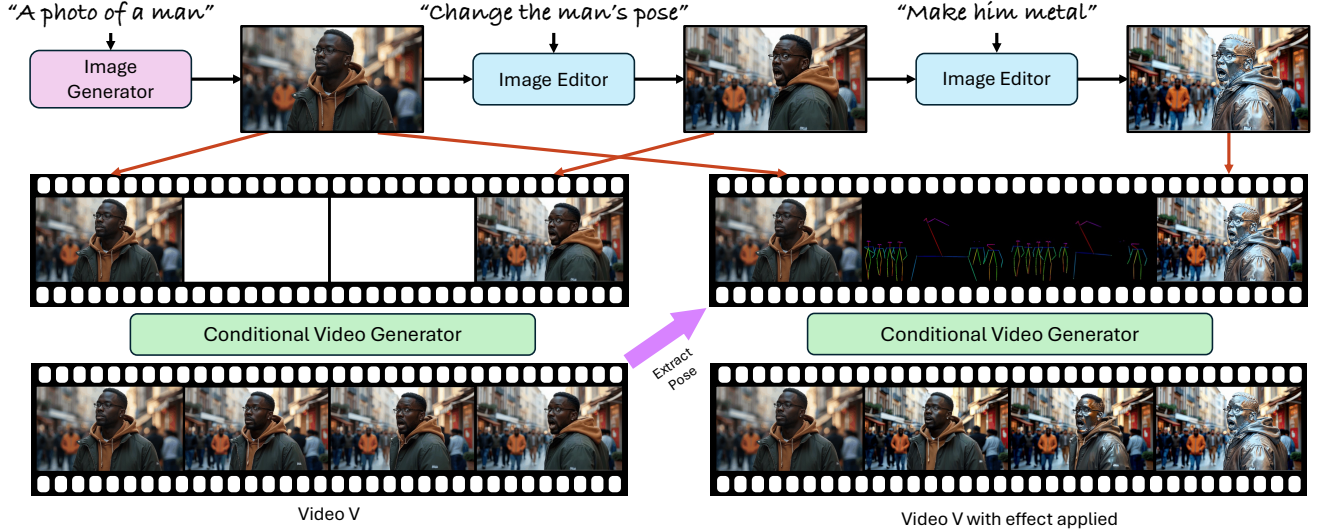
Figure 4. **Text-Guided Video Editing Pair Creation**. We present a method to generate a pair of videos $(V, V')$ from an effect prompt $E$, where $V$ is an initial video and $V'$ is the video with an effect $E$ applied. First, an image generation model is used to create an initial image. Next, an image editing model is used to change the pose, camera angle, and facial expression of the image. Finally, this image is again edited to add effect $E$. The first and second generated images are used with a first last frame model to output video $V$. Then, we use a conditional video model conditioned on the original first frame, effect edited last frame, *and intermediate poses from video $V$* as conditioning to create video $V'$.

### 3.2.3. Synthetic Dataset Creation.

We further construct a large-scale dataset of *synthetic motion-based effects* through programmatic generation. Our approach enables scalable and reproducible creation of diverse video effects with fine-grained control over appearance and temporal dynamics.

**Base Videos.** We source input videos from the Senorita dataset [89], which provides high-quality clips with accompanying foreground–background segmentation masks obtained using SegmentAnything [61]. These masks allow effects to be applied selectively to the full frame, the foreground, or the background.

**Effect Library.** We define a library of code-based effects $c_i$ such as posterization, pixelation, and dithering. Each effect includes hyperparameters that control its visual characteristics. For example, posterization varies in the number of color bins and palette type, while pixelation varies in block size. Sampling different parameter configurations yields a wide range of distinct appearances within each effect class.

**Temporal Compositing.** To introduce temporal variation, we consider a set of temporal transition operators $t_i$ (e.g., wipe-right, diagonal wipe, circular-out, etc). Each transition is also parameterized by temporal hyperparameters controlling its start time and duration.

**Final Composition.** For each composite effect $E_i$, we randomly sample an effect mask (e.g., foreground, background, or all), a spatial effect, its parameters, a temporal transition, and its corresponding temporal parameters.

See Figure 2 (c) for an example. This final effect $E_i$ is then applied to a fixed set of source videos, generating $K$ unique clips per effect. In total, this pipeline yields approximately 100K videos spanning 1,500 distinct synthetic motion-based effects.

### 3.3. Model Architecture and Conditioning

We build upon the Wan conditioning model architecture [36, 74], which extends text-to-video diffusion models to image-to-video (I2V), first-last frame to video (FLF2V), depth to video, and other conditioning types by concatenating noisy latent channels with conditional and mask channels. We extend this mechanism to jointly condition on both a reference video effect and an input video.

The input video's latents are supplied as conditioning latents, while the reference effect video is encoded into additional latents that are concatenated across all frames. This allows for spatial self-attention between noisy latent tokens with input video conditioning information and reference effect tokens during a forward pass. A hybrid mask controls which latent conditioning frames are preserved or edited based on which are directly copied to the final output and are changeable. The architecture is shown in Figure 5. Further architectural details, including input condition dropout and cross attention configurations are provided in the supplementary material.

### 3.4. Implementation Details

We fine-tune the Wan [74] 14 Billion parameter First-Last Frame to Video model with Low Rank Adapters for 10K

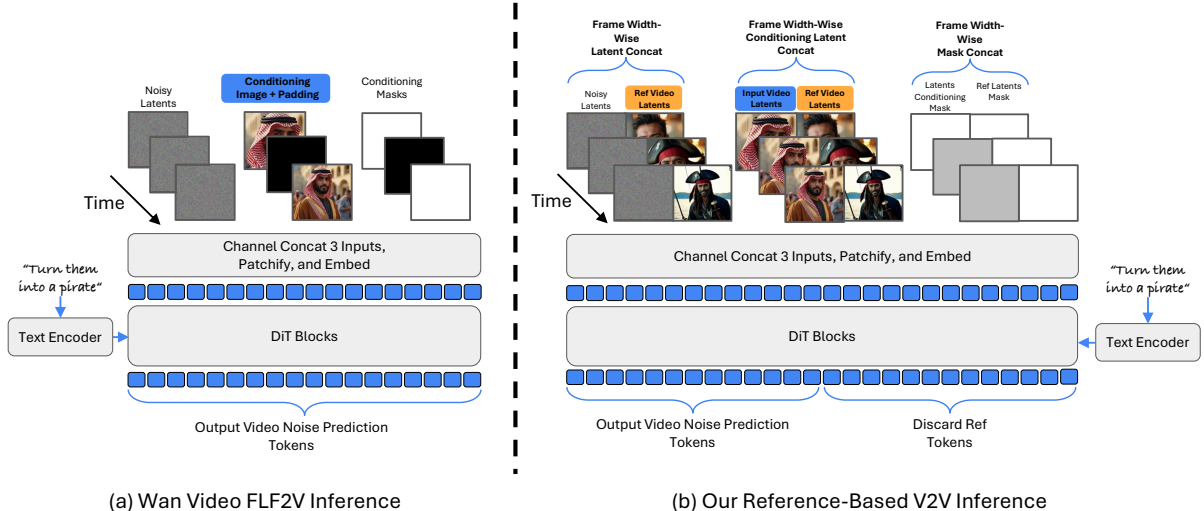(a) Wan Video FLF2V Inference      (b) Our Reference-Based V2V Inference

Figure 5. **Architecture overview**. (a) Standard Wan Video First–Last Frame to Video (FLF2V) architecture: noisy spatio-temporal latents are channel-wise concatenated with conditioning inputs and a mask, then patchified, embedded, and processed by the diffusion transformer to predict velocity. (b) In our setup, x input video latents are used as conditioning for the noisy latents, while reference video latents are concatenated width-wise to both. The latent mask is set to 1 for frames preserved exactly in the output and 0.5 for those to be modified; the reference latent mask is all ones. This design doubles the token count relative to base generation while conditioning on both the reference effect video and input video. Since all three inputs are channel-concatenated before patchification, repeated clean reference latents are merged channel-wise before embedding, ensuring no redundant reference information across tokens.

| Baseline | Neural V2V | | Code-Based V2V | |
|---|---|---|---|---|
| | RVA | IVA | RVA | IVA |
| Lucy Edit | 94.3 / +1.26 | 98.2 / +1.49 | – | – |
| VACE (Depth) | 91.7 / +1.07 | 70.0 / +0.05 | 98.3 / +1.51 | Tie / +0.00 |
| VACE (Pose) | 78.3 / +0.72 | 90.0 / +0.11 | 96.7 / +1.48 | Tie / +0.00 |
| No Ref | 92.9 / +0.18 | Tie / +0.00 | 93.6 / +0.94 | Tie / +0.00 |

Table 1. User Study Results for Neural V2V and Code-Based Edits. Each cell shows Win Rate (%) / Mean Score ($\pm 2$ scale). Our method shows large preference over base models in Reference Video Adherence (RVA), and gains in Input Video Adherence (IVA)

| Baseline | Reference Video Adherence |
|---|---|
| Wan2.1 | 74.7 / +0.57 |
| VACE (I2V) | 80.0 / +0.69 |
| No Ref | 81.5 / +0.15 |

Table 2. User Study Results for I2V. Each cell shows Win Rate (%) / Mean Score ($\pm 2$ scale). Our method shows large user preference over baselines in Reference Video Adherence.

steps with a batch size of 8 on a single node of 8 NVIDIA A100 GPUs. We sample from each of the three datasets equally during training, and drop the reference effect video conditioning, control video conditioning, and text prompts with low probability to allow for various forms of classifier-free guidance [8, 27] during inference time. We also add the true last frame from the target video to the model with low probability to allow for the first and last frame capabilities of the underlying model to be retained under our new setup.

## 4. Experiments

We evaluate our approach on both reference-based image-to-video and reference-based video-to-video generation tasks. In each setting, we provide a pair (reference video, input image or video), and assess the quality of our generated outputs against strong baselines. Our evaluation emphasizes generalization to unseen reference effects and motion patterns.

### 4.1. Validation Datasets

**Image-to-Video Testing.** We construct an evaluation set of unseen image-to-video effects sourced from publicly available LoRA-based visual effect models that were not included in training. For each LoRA effect, we generate multiple reference videos from distinct input images, resulting in 28 unique unseen LoRA effects and a total of 56 reference videos for evaluation.

**Video-to-Video Testing.** We generate a validation set of video-to-video pairs using the pipeline introduced in Section 3.2.2, with novel prompts not seen during training. Each resulting pair (reference video, input video) is manually filtered to ensure high perceptual quality. To further test cross-domain generalization, we reuse the unseen LoRA effects from the previous section but replace the input images with novel input videos, producing over 100 test pairs in total. Finally, we synthesize an additional set of temporally varying effects using the programmatic procedure from Section 3.2.3. Although these effects share the same families of operations and transitions as those used for train-

ing, their sampled hyperparameters are unseen at train time.

## 4.2. Baselines

Since existing baselines cannot directly condition on reference videos, we compare our method to state-of-the-art text-conditioned image-to-video and video-to-video models. All baselines receive the same textual descriptions of the desired effects as those implicit in our reference videos.

**Image-to-Video Baselines.** We benchmark against the Wan 2.1 image-to-video model [74] and the Wan VACE model [36], which represent current state-of-the-art text-guided video generation systems.

**Video-to-Video Baselines.** For the video-to-video task, we evaluate two configurations based upon the Wan VACE [36] conditional model and one native text-based video editing model. Pose-conditioned Wan VACE: In this setup, we condition the Wan VACE on the first frame of the input video, intermediate poses extracted from the input video [81], and the corresponding text prompt. Depth-conditioned: In this setup, we condition the Wan VACE on the first frame of the input video, intermediate depth maps [60] from the input video, and the same text prompt. Lucy Edit [26]: Lucy Edit is a text-based video editing model built on Wan 2.2 [26], which directly takes an input video and a text instruction to produce an edited output.

## 4.3. Qualitative Results

**Image-to-Video Effects.** Figure 6 shows qualitative comparisons for the image-to-video task using unseen reference effects. Baseline models, which rely solely on text prompts, struggle to capture fine-grained temporal and stylistic cues. In contrast, our method accurately transfers the reference video's camera motion, lighting, and character transformations, producing coherent temporal dynamics that align closely with the target effect.

**Video-to-Video Effects.** Qualitative comparisons for the video-to-video task are shown in Figure 7. Wan VACE variants either overfit to the input video or introduce artifacts, while Lucy Edit produces static, frame-invariant edits that lack temporal evolution. By conditioning on the full reference video, our method consistently reproduces complex time-varying transformations—such as gradual color shifts or structural morphing—while maintaining spatial consistency and motion alignment with the input.

**User Study** For a subjective task such as reference based video editing or image to video, it is hard to determine any one quantitative metric that fully encompasses success. As a result, we conduct a comprehensive user study as the main way to evaluate our method across three tasks: neural video-to-video effect transfer (V2V), image-to-video effect transfer (I2V), and code-based temporal effects. Annotators performed pairwise comparisons between our method and several baselines, including Lucy Edit, VACE (Depth/-Pose/I2V variants), Wan2.1 , and ablations (Use our trained
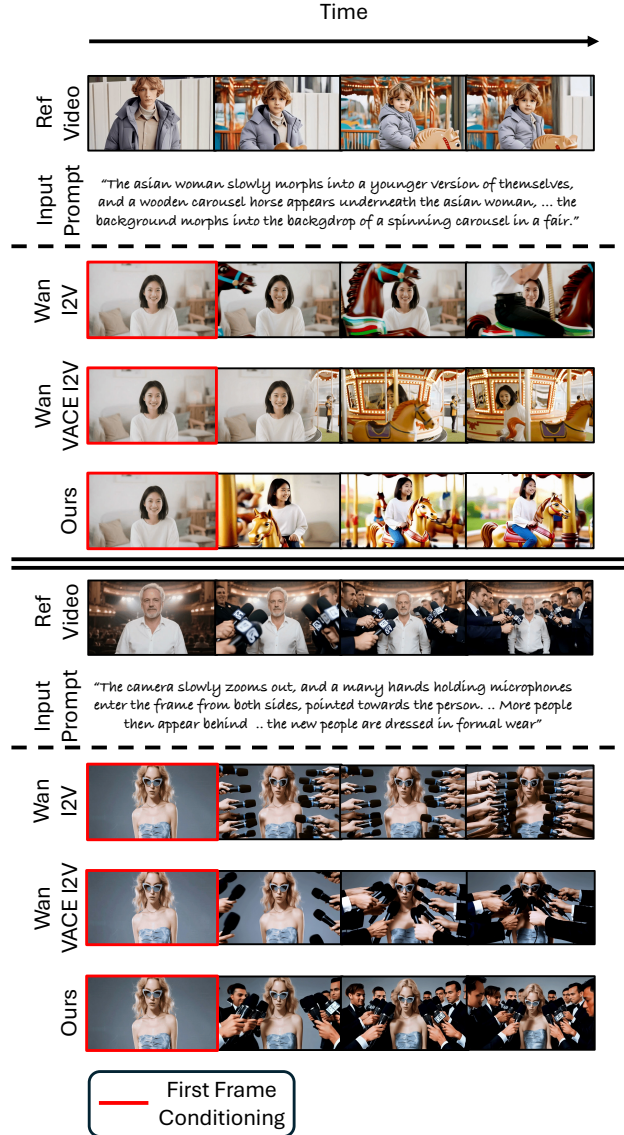


Figure 6. **Qualitative Reference Based Image to Video Results.** In the top example, baselines are unable to turn the input woman into a younger version of herself, or put her on the carousel. Our method correctly makes the subject younger, puts her on the carousel, and mimics the camera motion of the reference video. In the bottom example, baselines are unable to add the reporters into the scene, while our method correctly adds hands with microphones followed by reporters, and mimics the interactions and occlusions from the reference video.

model, but do drop the reference effect video conditioning) , rating each pair on Reference Video Adherence and Input Video Adherence dimensions using a granular 5-level scale. We collected 669 valid annotations for V2V, 455 for I2V, and 252 for code-based effects. Our method demonstrates strong performance in reference adherence, achieving win rates of 78-94% against competitive baselines in V2V, 94-98% in code-based effects, and 75-82% in I2V, as shown
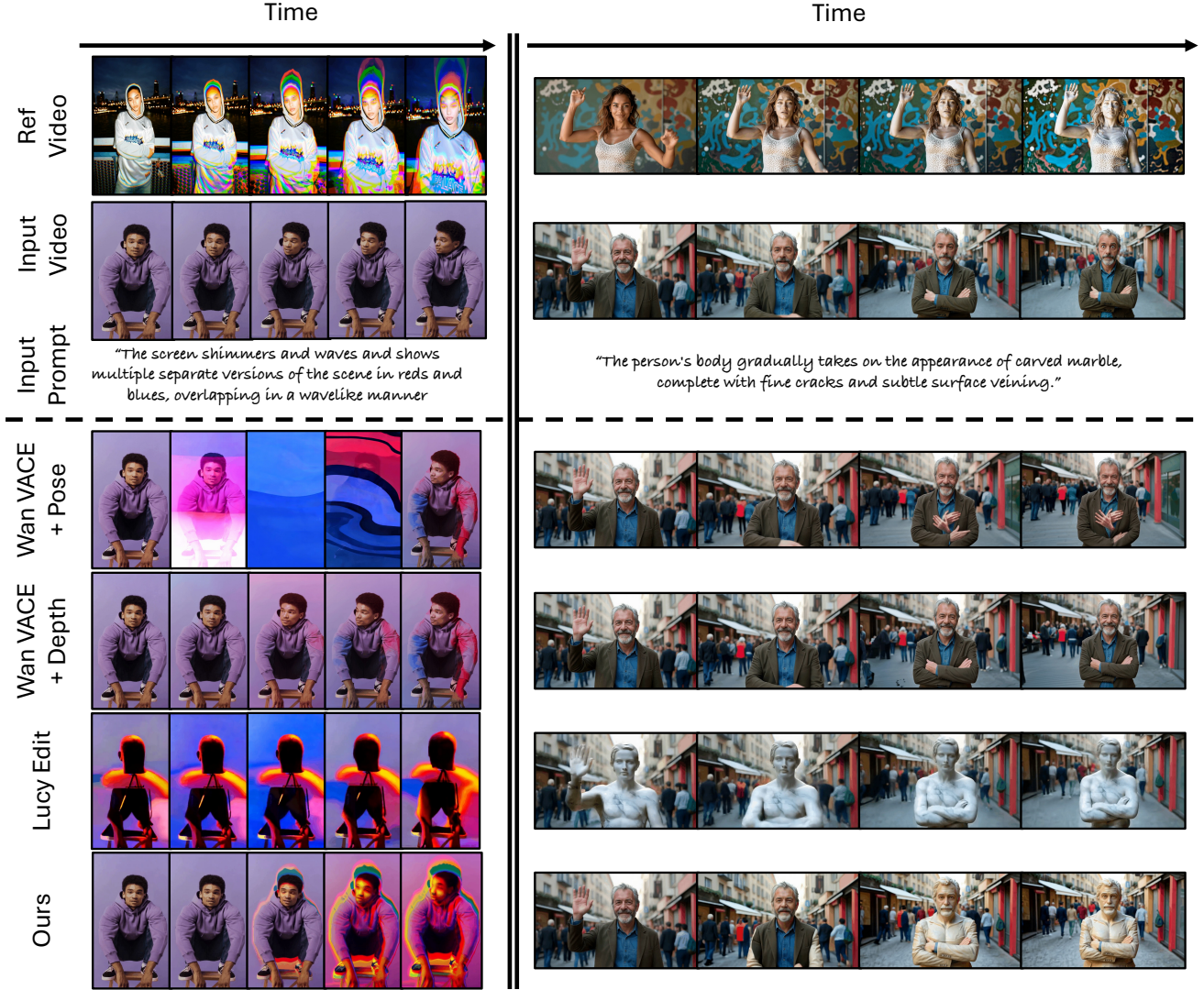
Figure 7. **Results for our reference-based video-to-video setup**. We compare three baselines: (1) Wan VACE + Pose: a conditional diffusion transformer [36] using the first input frame and extracted poses; (2) Wan VACE + Depth: using the first frame and extracted depths with Wan VACE; and (3) Lucy Edit [26], a video editing model fine-tuned from an image-to-video model. Without reference video information, all baselines fail to follow the reference from a prompt alone. Wan VACE methods overfit to the input or add irrelevant edits, while Lucy Edit applies uniform static edits across frames. Our method successfully follows the reference effect.

in Tables 2 and 1. For input video adherence in V2V, we achieve win rates of 70-98% against method-based baselines (Lucy Edit, VACE Depth/Pose), as shown in Table 1, while comparisons with our ablations show expected behavior: ties with No Style (which uses identical input processing). Mean scores across all dimensions are consistently positive (+0.15 to +1.51), indicating systematic preference for our method's outputs, and we provide full results in the supplement. These results validate our approach's effectiveness in effect transfer while preserving input content fidelity across multiple generation modalities. For more details on the user study setup, please refer to the Supplement.

## 4.4. Quantitative Metrics

**Conditioning Input Similarity.** Since our validation sets do not include ground-truth triplets, we measure the similarity of our generated outputs to both the input and reference videos as a proxy for task success. To quantify this, we employ VideoPrism [87], a large-scale video embedding model pretrained on over 500M clips, to compute feature-space similarities between videos. For the reference + image-to-video case, we instead measure the average similarity between the generated video and its first-frame conditioning using CLIP [59] embeddings. Results are summarized in Tables 3 and 4. Our method consistently achieves higher
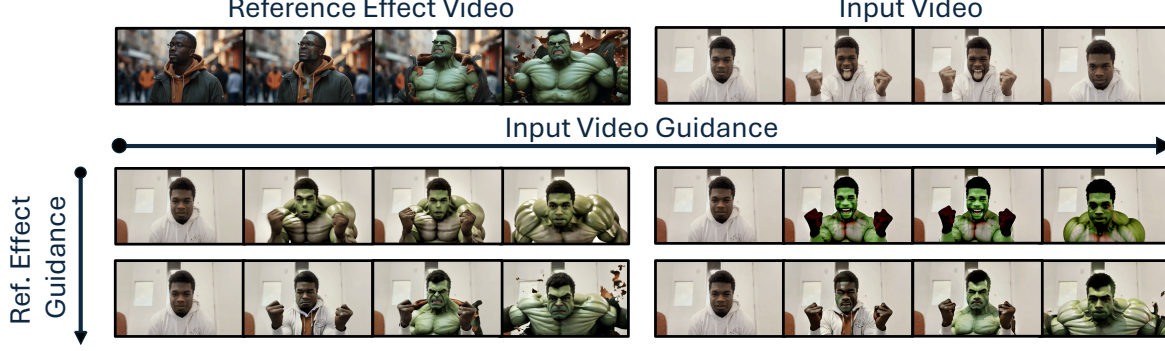
Figure 8. **Qualitative Results for Controllable Real Video Editing.** We show an example of the effect of reference effect guidance and input video guidance on model output. The first row shows the reference effect video and input video that are used for model conditioning. From left to right, we increase the input video guidance, while from top to bottom we increase the reference effect guidance. Notice that as input video guidance increases, the input video details are more faithfully kept (i.e. attire, facial shape). As the reference effect guidance increases, the reference effect video takes a larger role (i.e. reference effect muscle shape, and clothing of the man in the reference effect). Applying both gives a good mix of both input video adherence while following the effect.

similarity to the reference video than all baselines, confirming that it effectively incorporates temporal and stylistic information from the reference. Interestingly, baseline models often exhibit slightly higher similarity to the input, likely reflecting a form of under-editing. This phenomenon is also visible in Figure 7 (left) and Figure 6 (bottom), where baselines preserve input structure but fail to reproduce the desired temporal evolution.

### 4.5. Controllable Real Video Editing

When editing real videos, users may want to control the level to which the original video is maintained, as well as the amount that the reference effect video effects the final generation. To this end, we follow previous work [8, 27, 37, 41, 42, 71] and combine multiple classifier-free guidance directions during inference. Specifically, we consider applying guidance for text, input video, and reference effect video, normalizing latents as recomended in previous work [42]. This gives:

$$
\begin{aligned}
& v_\theta\left(x^t, x_{\text{ref}}, x_{\text{input}}, \varnothing\right) + \\
& \lambda_c \frac{\|g\|}{\|g_c\|} g_c + \lambda_{\text{ref}} \frac{\|g\|}{\|g_{\text{ref}}\|} g_{\text{ref}} + \lambda_{\text{in}} \frac{\|g\|}{\|g_{\text{in}}\|} g_{\text{in}},
\end{aligned} \quad (1)
$$

where

$$
g_c = v_\theta\left(x^t, x_{\text{ref}}, x_{\text{input}}, c\right) - v_\theta\left(x^t, x_{\text{ref}}, x_{\text{input}}, \varnothing\right), \quad (2)
$$
$$
g_{\text{ref}} = v_\theta\left(x^t, x_{\text{ref}}, x_{\text{input}}, c\right) - v_\theta\left(x^t, \varnothing, x_{\text{input}}, c\right), \quad (3)
$$
$$
g_{\text{in}} = v_\theta\left(x^t, x_{\text{ref}}, x_{\text{input}}, c\right) - v_\theta\left(x^t, x_{\text{ref}}, \varnothing, c\right), \quad (4)
$$
$$
\|g\| = \min(\|g_c\|, \|g_{\text{ref}}\|, \|g_{\text{in}}\|). \quad (5)
$$

and $x^t$ is the noisy video latent, $x_{\text{ref}}$ is the clean reference video effect latent, $x_{\text{input}}$ is the clean input video latent, and $c$ is the text conditioning. Interpolating $\lambda_{\text{ref}}$ and $\lambda_{\text{in}}$ results in output videos with difference amounts of adherence to

| Method | First Frame Sim. | Ref Sim. |
|---|---|---|
| Wan 2.1 | **0.7911** | 0.7230 |
| Wan VACE I2V | 0.7799 | 0.7127 |
| Ours | 0.7698 | **0.7378** |

Table 3. I2V Similarities. First Frame Sim is average clip [59] similarity to the input frame, and Ref Sim. is video embedding similarity [87] between the output video and the reference effect video

| Method | Neural V2V | | Code Based V2V | |
|---|---|---|---|---|
| | Input Sim. | Ref Sim. | Input Sim. | Ref Sim. |
| Wan VACE Pose | 0.9068 | 0.6539 | 0.9225 | 0.6002 |
| Wan VACE Depth | **0.9460** | 0.6226 | 0.9394 | 0.5998 |
| Lucy Edit | 0.7544 | 0.6852 | - | - |
| Ours | 0.8568 | **0.7014** | 0.9479 | **0.7169** |

Table 4. Neural and Code-Based V2V Similarities. Input Sim. is video embedding similarity [87] between the output video and the input video, and Ref Sim. is the video embedding similarity between the output video and the reference effect video.

the reference effect video and input video respectively, as seen in Figure 8.

### 5. Discussion

In conclusion, we presented RefVFX, a tuning-free framework for transferring complex temporal visual effects across videos using reference conditioning. By introducing a large-scale dataset of over 120K triplets encompassing 1,700 distinct effects, we enable the first systematic study of reference-based temporal effect transfer. Our scalable video-to-video generation pipeline produces diverse, motion-consistent training pairs, while our diffusion-based architecture jointly encodes reference dynamics, input appearance, and semantic guidance to produce coherent results. Perceptual and quantitative evaluations demonstrate that RefVFX surpasses prompt-only and static reference baselines in both visual fidelity and temporal coherence.

# References

[1] Wan2.1 14b 480p i2v loras. https://huggingface.co/collections/Remade-AI/wan21-14b-480p-i2v-loras. Accessed: 2025-10-21. 14

[2] Rameen Abdal, Or Patashnik, Ekaterina Deyneka, Hao Chen, Aliaksandr Siarohin, Sergey Tulyakov, Daniel Cohen-Or, and Kfir Aberman. Zero-shot dynamic concept personalization with grid-based lora, 2025. 3

[3] Rameen Abdal, Or Patashnik, Ivan Skorokhodov, Willi Menapace, Aliaksandr Siarohin, Sergey Tulyakov, Daniel Cohen-Or, and Kfir Aberman. Dynamic concepts personalization from single videos. In *Proceedings of the Special Interest Group on Computer Graphics and Interactive Techniques Conference Conference Papers*, New York, NY, USA, 2025. Association for Computing Machinery. 3

[4] Michael S Albergo, Nicholas M Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A unifying framework for flows and diffusions. *arXiv preprint arXiv:2303.08797*, 2023. 3

[5] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. 2, 3

[6] Frederic Boesel and Robin Rombach. Improving image editing models with generative data refinement. In *The Second Tiny Papers Track at ICLR 2024*, 2024. 3

[7] 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 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8861–8870, 2024. 3

[8] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18392–18402, 2023. 3, 6, 9

[9] 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 *Proceedings of the IEEE/CVF international conference on computer vision*, pages 22560–22570, 2023. 3

[10] Duygu Ceylan, Chun-Hao P Huang, and Niloy J Mitra. Pix2video: Video editing using image diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 23206–23217, 2023. 3

[11] Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, Yuwei Fang, Ivan Skorokhodov, Jun-Yan Zhu, Kfir Aberman, Ming-Hsuan Yang, and Sergey Tulyakov. Videoalchemy: Open-set personalization in video generation. 2024. 3

[12] Xi Chen, Lianghua Huang, Yu Liu, Yujun Shen, Deli Zhao, and Hengshuang Zhao. Anydoor: Zero-shot object-level image customization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6593–6602, 2024. 3

[13] Jiaxin Cheng, Tianjun Xiao, and Tong He. Consistent video-to-video transfer using synthetic dataset. *arXiv preprint arXiv:2311.00213*, 2023. 3

[14] cloneofsimo. fofr/cog-aura-flow. https://github.com/fofr/cog-aura-flow, 2024. GitHub repository. 3

[15] Gilad Deutch, Rinon Gal, Daniel Garibi, Or Patashnik, and Daniel Cohen-Or. Turboedit: Text-based image editing using few-step diffusion models. In *SIGGRAPH Asia 2024 Conference Papers*, pages 1–12, 2024. 3

[16] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 3

[17] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024. 2, 3

[18] Junyao Gao, Yanan Sun, Yanchen Liu, Yinhao Tang, Yanhong Zeng, Ding Qi, Kai Chen, and Cairong Zhao. Styleshot: A snapshot on any style. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025. 3

[19] Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Tokenflow: Consistent diffusion features for consistent video editing. *arXiv preprint arXiv:2307.10373*, 2023. 3

[20] Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023. 3

[21] Zinan Guo, Yanze Wu, Chen Zhuowei, Peng Zhang, Qian He, et al. Pulid: Pure and lightning id customization via contrastive alignment. *Advances in neural information processing systems*, 37:36777–36804, 2024. 3

[22] 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. 2, 3

[23] Xuanhua He, Quande Liu, Shengju Qian, Xin Wang, Tao Hu, Ke Cao, Keyu Yan, and Jie Zhang. Id-animator: Zero-shot identity-preserving human video generation. *arXiv preprint arXiv:2404.15275*, 2024. 3

[24] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. 3

[25] Amir Hertz, Andrey Voynov, Shlomi Fruchter, and Daniel Cohen-Or. Style aligned image generation via shared attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4775–4785, 2024. 3

[26] Sari Hleihil. Decartai/Lucy-Edit-ComfyUI. https://github.com/DecartAI/Lucy-Edit-ComfyUI, 2025. 2, 7, 8

[27] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 6, 9

[28] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 3

[29] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in neural information processing systems*, 35:8633–8646, 2022.

[30] 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. 2, 3

[31] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022. 3, 4

[32] Lianghua Huang, Di Chen, Yu Liu, Yujun Shen, Deli Zhao, and Jingren Zhou. Composer: Creative and controllable image synthesis with composable conditions. *arXiv preprint arXiv:2302.09778*, 2023. 3

[33] Yuzhou Huang, Ziyang Yuan, Quande Liu, Qiulin Wang, Xintao Wang, Ruimao Zhang, Pengfei Wan, Di Zhang, and Kun Gai. Conceptmaster: Multi-concept video customization on diffusion transformer models without test-time tuning. *arXiv preprint arXiv:2501.04698*, 2025. 3

[34] 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. 14, 15

[35] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 4

[36] 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, 4, 5, 7, 8

[37] Maxwell Jones, Sheng-Yu Wang, Nupur Kumari, David Bau, and Jun-Yan Zhu. Customizing text-to-image models with a single image pair. In *SIGGRAPH Asia 2024 Conference Papers*, pages 1–13, 2024. 9

[38] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. Hunyuanvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024. 2, 3

[39] Vladimir Kulikov, Matan Kleiner, Inbar Huberman-Spiegelglas, and Tomer Michaeli. Flowedit: Inversion-free text-based editing using pre-trained flow models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19721–19730, 2025. 3

[40] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1931–1941, 2023. 2

[41] Nupur Kumari, Grace Su, Richard Zhang, Taesung Park, Eli Shechtman, and Jun-Yan Zhu. Customizing text-to-image diffusion with object viewpoint control. 2024. 9

[42] Nupur Kumari, Xi Yin, Jun-Yan Zhu, Ishan Misra, and Samaneh Azadi. Generating multi-image synthetic data for text-to-image customization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16524–16534, 2025. 3, 9

[43] Black Forest Labs. Flux. https://github.com/black-forest-labs/flux, 2024. 2, 3, 4

[44] Black Forest Labs, Stephen Batifol, Andreas Blattmann, Frederic Boesel, Saksham Consul, Cyril Diagne, Tim Dockhorn, Jack English, Zion English, Patrick Esser, et al. Flux. 1 kontext: Flow matching for in-context image generation and editing in latent space. *arXiv preprint arXiv:2506.15742*, 2025. 3, 4

[45] LAION-AI. aesthetic-predictor. https://github.com/LAION-AI/aesthetic-predictor, 2022. GitHub repository. 14

[46] Dongxu Li, Junnan Li, and Steven Hoi. Blip-diffusion: Pre-trained subject representation for controllable text-to-image generation and editing. *Advances in Neural Information Processing Systems*, 36:30146–30166, 2023. 3

[47] Feng Liang, Haoyu Ma, Zecheng He, Tingbo Hou, Ji Hou, Kunpeng Li, Xiaoliang Dai, Felix Juefei-Xu, Samaneh Azadi, Animesh Sinha, et al. Movie weaver: Tuning-free multi-concept video personalization with anchored prompts. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 13146–13156, 2025. 3

[48] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022. 3

[49] Gongye Liu, Menghan Xia, Yong Zhang, Haoxin Chen, Jinbo Xing, Yibo Wang, Xintao Wang, Yujiu Yang, and Ying Shan. Stylecrafter: Enhancing stylized text-to-video generation with style adapter. *arXiv preprint arXiv:2312.00330*, 2023. 2, 3

[50] Lijie Liu, Tianxiang Ma, Bingchuan Li, Zhuowei Chen, Jiawei Liu, Gen Li, Siyu Zhou, Qian He, and Xinglong Wu. Phantom: Subject-consistent video generation via cross-modal alignment. *arXiv preprint arXiv:2502.11079*, 2025. 3

[51] 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. 3

[52] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022. 3

[53] Xin Ma, Yaohui Wang, Gengyun Jia, Xinyuan Chen, Ziwei Liu, Yuan-Fang Li, Cunjian Chen, and Yu Qiao. Latte: Latent diffusion transformer for video generation. *arXiv preprint arXiv:2401.03048*, 2024. 3

[54] Gaurav Parmar, Or Patashnik, Kuan-Chieh Wang, Daniil Ostashev, Srinivasa Narasimhan, Jun-Yan Zhu, Daniel Cohen-Or, and Kfir Aberman. Object-level visual prompts for compositional image generation. *arXiv preprint arXiv:2501.01424*, 2025. 3

[55] Or Patashnik, Rinon Gal, Daniil Ostashev, Sergey Tulyakov, Kfir Aberman, and Daniel Cohen-Or. Nested attention: Semantic-aware attention values for concept personalization. In *Proceedings of the Special Interest Group on Computer Graphics and Interactive Techniques Conference Conference Papers*, pages 1–12, 2025. 3

[56] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023. 3

[57] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 2, 3

[58] Adam Polyak, Amit Zohar, Andrew Brown, Andros Tjandra, Animesh Sinha, Ann Lee, Apoorv Vyas, Bowen Shi, Chih-Yao Ma, Ching-Yao Chuang, et al. Movie gen: A cast of media foundation models. *arXiv preprint arXiv:2410.13720*, 2024. 2, 3

[59] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 8, 9, 15

[60] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE transactions on pattern analysis and machine intelligence*, 44(3):1623–1637, 2020. 7

[61] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 5

[62] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2, 3

[63] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 3

[64] Litu Rout, Yujia Chen, Nataniel Ruiz, Abhishek Kumar, Constantine Caramanis, Sanjay Shakkottai, and Wen-Sheng Chu. Rb-modulation: Training-free personalization of diffusion models using stochastic optimal control. *arXiv preprint arXiv:2405.17401*, 2024. 3

[65] 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. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22500–22510, 2023. 2, 3

[66] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in neural information processing systems*, 35:25278–25294, 2022. 15

[67] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022. 3

[68] Kihyuk Sohn, Nataniel Ruiz, Kimin Lee, Daniel Castro Chin, Irina Blok, Huiwen Chang, Jarred Barber, Lu Jiang, Glenn Entis, Yuanzhen Li, et al. Styledrop: Text-to-image generation in any style. *arXiv preprint arXiv:2306.00983*, 2023. 3

[69] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 3

[70] Kunpeng Song, Yizhe Zhu, Bingchen Liu, Qing Yan, Ahmed Elgammal, and Xiao Yang. Moma: Multimodal llm adapter for fast personalized image generation. In *European Conference on Computer Vision*, pages 117–132. Springer, 2024. 3

[71] Kiwhan Song, Boyuan Chen, Max Simchowitz, Yilun Du, Russ Tedrake, and Vincent Sitzmann. History-guided video diffusion. *arXiv preprint arXiv:2502.06764*, 2025. 9

[72] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *European conference on computer vision*, pages 402–419. Springer, 2020. 15

[73] Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for text-driven image-to-image translation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1921–1930, 2023. 3

[74] 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. 2, 3, 4, 5, 7, 14

[75] 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*, 36:7594–7611, 2023. 3

[76] Cong Wei, Zheyang Xiong, Weiming Ren, Xeron Du, Ge Zhang, and Wenhu Chen. Omniedit: Building image editing generalist models through specialist supervision. In *The Thirteenth International Conference on Learning Representations*, 2024. 3

[77] Yuxiang Wei, Yabo Zhang, Zhilong Ji, Jinfeng Bai, Lei Zhang, and Wangmeng Zuo. Elite: Encoding visual concepts into textual embeddings for customized text-to-image

generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15943–15953, 2023. 3

[78] Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng-ming Yin, Shuai Bai, Xiao Xu, Yilei Chen, et al. Qwen-image technical report. *arXiv preprint arXiv:2508.02324*, 2025. 3, 4, 14

[79] Guangxuan Xiao, Tianwei Yin, William T Freeman, Frédo Durand, and Song Han. Fastcomposer: Tuning-free multi-subject image generation with localized attention. *International Journal of Computer Vision*, 133(3):1175–1194, 2025. 3

[80] Peng Xing, Haofan Wang, Yanpeng Sun, Qixun Wang, Xu Bai, Hao Ai, Renyuan Huang, and Zechao Li. Csgo: Content-style composition in text-to-image generation. *arXiv preprint arXiv:2408.16766*, 2024. 3

[81] Zhendong Yang, Ailing Zeng, Chun Yuan, and Yu Li. Effective whole-body pose estimation with two-stages distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4210–4220, 2023. 4, 7

[82] Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721*, 2023. 2, 3

[83] Zixuan Ye, Huijuan Huang, Xintao Wang, Pengfei Wan, Di Zhang, and Wenhan Luo. Stylemaster: Stylize your video with artistic generation and translation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 2630–2640, 2025. 2, 3

[84] Shenghai Yuan, Jinfa Huang, Xianyi He, Yunyang Ge, Yujun Shi, Liuhan Chen, Jiebo Luo, and Li Yuan. Identity-preserving text-to-video generation by frequency decomposition. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 12978–12988, 2025. 3

[85] Gong Zhang, Kihyuk Sohn, Meera Hahn, Humphrey Shi, and Irfan Essa. Finestyle: Fine-grained controllable style personalization for text-to-image models. *Advances in Neural Information Processing Systems*, 37:52937–52961, 2024. 3

[86] Kai Zhang, Lingbo Mo, Wenhu Chen, Huan Sun, and Yu Su. Magicbrush: A manually annotated dataset for instruction-guided image editing. *Advances in Neural Information Processing Systems*, 36:31428–31449, 2023. 3

[87] Long Zhao, Nitesh B Gundavarapu, Liangzhe Yuan, Hao Zhou, Shen Yan, Jennifer J Sun, Luke Friedman, Rui Qian, Tobias Weyand, Yue Zhao, et al. Videoprism: A foundational visual encoder for video understanding. *arXiv preprint arXiv:2402.13217*, 2024. 8, 9

[88] Min Zhao, Rongzhen Wang, Fan Bao, Chongxuan Li, and Jun Zhu. Controlvideo: conditional control for one-shot text-driven video editing and beyond. *Science China Information Sciences*, 68(3):132107, 2025. 2

[89] 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. 2, 3, 5, 14

In Section A, we present limitations of our generated model, and show a qualitative result of a current limitation. In Section B, we provide detailed descriptions of our dataset creation process, including the generation pipelines and effect categorization. Section C outlines additional implementation details such as training configurations, inference settings, and computational requirements. In Section D, we report extended quantitative comparisons between our method and baseline models across multiple metrics. **We provide extensive qualitative video results including outputs from our method, comparisons with baselines, and demonstrations from our new dataset in the accompanying website.html file located in the website folder of the supplement**.

## A. Limitations

While RefVFX achieves strong temporal coherence and visual fidelity across diverse effects, it still faces certain limitations. First, the model can struggle with accurately reproducing fine-grained occlusions or complex interactions between subjects and dynamic effects, occasionally leading to partial blending or misalignment artifacts. Second, our dataset primarily focuses on human-centric and foreground-dominant scenes, which may limit generalization to large-scale environmental effects or abstract cinematic transformations. We show an example of such limitations in Figure 9. Inference remains computationally expensive due to the dual conditioning on both input and reference videos, resulting in roughly twice the generation time compared to single-source baselines. Addressing these limitations presents promising directions for future work.

## B. Dataset Details

We source the Image-to-Video data from open-source LoRA effects trained on top of the Wan 2.1 Image-to-Video model [74], available on Hugging Face [1]. A comprehensive list of LoRA effects and their corresponding captions is provided in Table 5.

We use Qwen Image [78] and Qwen Image Edit [78], together with Wan 2.1 VACE-1.3B [74], to implement our text-guided video editing pair creation pipeline, as described in Figure 4 of the main text. We categorize our effect types and general effect categories in Table 6.

For the code-based effects, we source data from the grounding subset of the Senorita dataset [89]. Specifically, we filter for human-centric videos and include only grounding results with human-centric masks, yielding a collection of over 100K videos. The set of objects considered is listed in Table 7. In the original dataset, input videos are used to predict grounding masks; we repurpose these masks to obtain foreground or object and background segmentations, which we use to augment our effect pipelines. We manu-
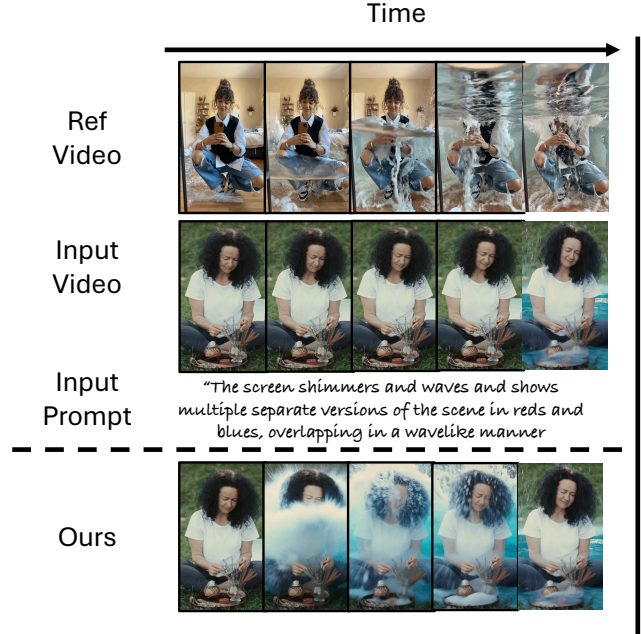


Figure 9. **Limitations.** We present a limitation of our model. Our model can struggle with imitating complete occlusion of parts of the body. In this example, the reference video has their entire head occluded by the rising water. Further, our method mistakenly adds a splashing effect that temporarily covers the subject but then subsides, as opposed to submerging them in standing water

ally curate code across a diverse set of effects and temporal transitions to sample from when creating triplets of the form (reference video, input video, output video). The full set of effect types and their corresponding hyperparameters is provided in Table 8, and the full set of temporal effect types and their corresponding hyperparameters is given in Table 9.

## C. More Implementation Details

We train our model on videos consisting of 33 frames at 15 frames per second, using a resolution of 480p for all training and inference experiments. During inference, we use a classifier-free guidance scale of 5 and perform 50 sampling steps per generation. Because our model conditions on both the input and reference effect videos, it requires approximately twice the number of latent tokens compared to the baseline model that omits reference conditioning. As a result, inference time is roughly doubled: baseline video generation takes about 3.5 minutes per video, whereas our method requires approximately 7 minutes per video on a single NVIDIA A100-SXM4 GPU.

## D. More Quantitative Results

We further assess video fidelity and perceptual quality using established metrics from VBench [34], including motion smoothness [74], aesthetic score [45], and dynamic de-

gree [34, 72].

Motion quality measures temporal consistency by computing the average cosine similarity [59] between consecutive frames:

$$\sum_{i=2}^{N} \text{cosine sim}(\text{CLIP}(f_i), \text{CLIP}(f_{i-1})). \quad (6)$$

Higher values indicate smoother and more coherent motion.

Aesthetic score follows the setup in VBench [34], using CLIP embeddings passed through a linear layer trained to predict human-rated aesthetic quality from a large-scale dataset of labeled images [66].

Dynamic degree quantifies how often videos exhibit meaningful motion. For each video, we compute optical flow using RAFT [72] and identify the top 5 percent of pixels with the highest flow magnitude per frame. Frames exceeding a motion threshold are labeled as dynamic; a video is then classified as dynamic if a sufficient proportion of its frames meet this criterion. The dynamic degree of a group of videos is defined as the percentage of videos that are labeled dynamic. We show extra qualitative results in Table 10. Across these quantitative measures, all methods perform comparably, indicating that our method can successfully transfer reference effects without sacrificing overall video quality.

| LoRA Name | Text Prompts |
| --- | --- |
| Squish | In the video, a miniature {subject} is presented. The {subject} is held in a person's hands. The person then presses on the {subject}, causing a sq41sh squish effect. The person keeps pressing down on the {subject}, further showing the sq41sh squish effect. |
| Rotate | The video shows a {subject}. The {subject} performs a r0t4tion 360 degrees rotation. |
| Inflate | (1) The video shows a {subject}, then infl4t3 inflates it, form expanding like a whimsical, cartoonish balloon. (2) The video shows a {subject}, then infl4t3 inflates it, form expanding into a perfect, inflated sphere. |
| Cakeify | The video opens on a {subject}. A knife, held by a hand, hovers over the {subject}. The knife then begins cutting into the {subject} to c4k3 cakeify it. As it slices, the inside of the {subject} is revealed to be cake with chocolate layers. |
| Deflate | The video opens with a {subject}. As the video progresses, the {subject} begins to d3d1at3 deflate it, gradually shrinking and losing shape, eventually flattening completely into a lifeless, deflated mass. |
| Crush | The video begins with a {subject}. A hydraulic press descends toward it. Upon contact, the press c5us4 crushes the {subject}, flattening and deforming it until it collapses inward. |
| Muscle | {subject} smiles slightly, then t2k1s takes off clothes revealing a lean muscular body and shows off muscles, giving a friendly smile. |
| Bride | The video begins with a {subject}, then the 8r1d3 bride effect occurs. The {subject} is now in a white wedding dress, holding a bouquet, with a sunny beige background. |
| Puppy | The video begins with a close-up portrait of a {subject}. The background changes and then the pu11y puppy effect begins. The {subject} is now surrounded by puppies and pets them. |
| Baby | The video starts with a {subject}. Then the image shifts to the 848y baby effect, with {subject} in front of a crib, surrounded by toys, then shown again in a different location as a baby version. |
| VIP | The video begins with an image of {subject}. Then the v1p red carpet transformation appears — {subject} in a black dress, gold jewelry, photographed on the red carpet. |
| Mona-Lisa | The video starts with an image of {subject}. The m0n4 Mona Lisa transformation begins, wrapping around the {subject}, who is then depicted as a Mona Lisa version seated before a landscape. |
| Princess | The video begins with a {subject}. A pr1nc355 princess transformation occurs: sparkling light appears, and {subject} is now in a silver beaded gown with tiara and gloves, seated among gifts and candles. |
| Pirate-Captain | The video begins with a {subject}. The p1r4t3 pirate captain transformation follows. The {subject} now wears a black hat, coat, and sash aboard a wooden ship with a sword. |
| Samurai | The video begins with a {subject}. The 54mur41 samurai transformation turns them into a samurai wearing traditional armor with a katana, against a misty mountain backdrop. |
| Zen | The video starts with a portrait of a {subject}. The z3n1fy zen transformation follows: pink robe, zen garden, and later a black kimono in a tranquil garden setting. |
| Assasin | The video starts with a portrait of a {subject}. The 3p1c epic transformation begins. The {subject} wears a red coat, white hair, and black gloves, holding guns in both hands. |
| Painting | The video starts with a {subject}. They appear in a gold framed mirror, then transform into a p41nt1ng painting version in red and blue attire with an old-style painted background. |
| Disney-Princess | The video starts with a {subject}. The d15n3y princess transformation occurs: the {subject} is now in a blue dress, with butterflies falling in a classic hallway setting. |
| Snow-White | The video begins with a {subject} outdoors, then cuts to the sn0w_wh1t3 transformation: classic dress, red apple, forest background. |
| Classy | The video starts with a {subject} in a suit. The c1455y transformation occurs — now in a light blue dress, smiling at the camera, seated with a flower and envelope. |

Table 5. List of LoRA models and their corresponding text prompts used in this work. Each LoRA entry spans multiple lines for readability.

| Category | Paraphrased Effects |
|---|---|
| **Object Addition** | Balance a tall stack of encyclopedias on head, wear a beekeeper suit with neon honeycomb glow, add a pizza cape, perch a tiny green dragon on shoulder, materialize toga + sandals + holographic smartwatch, hat becomes a live octopus, juggle three bright-green rubber chickens, use a banana as a phone, parade of wind-up toy robots follows behind, attach a leash to a pet cloud, grow a gnome hat and long white beard, put on a diver's helmet (library setting), don an astronaut helmet + 19th-century ball gown, tuxedo T-shirt appears and tiny penguins roam shoulders, hands become giant foam novelty hands, endless silk scarves stream from ear, open a shimmering interdimensional portal, person covered head-to-toe in rainbow alien slime, business suit catches fire (person stays calm), three personal mini-planets orbit the person. |
| **Weather + Atmospheric Effects** | Morning fog rolls in and softens background, thick industrial smog mutes colors, gentle misty drizzle begins, heavy rainfall drenches the scene, post-rain surfaces glisten like mirrors, light picturesque snowfall, heavy snowstorm with blowing flakes, full-blown blizzard white-out, hailstones streak and bounce, sandstorm tints scene orange and dusty, smoke fills air as if from nearby fire, floating embers and rising ash, ash and embers fall from sky, visible dust motes in sunbeams, aurora borealis bathes scene in eerie glow, electric storm with flashing lightning, heat haze shimmers and warps background, underwater environment with flowing caustics, moonrise casts pale long shadows, volumetric light rays through window/trees. |
| **Artistic & Stylistic Effects** | Risograph print look (overlapping magenta/cyan layers), stained-glass pane refractions, cross-hatched pencil sketch rendering, duotone (magenta + black), thermal-camera false-color map, infrared look (white foliage, dark sky), holographic scanline overlay, datamosh/pixel-sorting glitch, digital code "Matrix" rain, stop-motion claymation style, low-poly geometric rendering, blueprint schematic overlay, graphic-novel/comic-book inks, chalk drawing on dusty blackboard, aged fresco on cracked plaster, impressionistic oil painting, soft watercolor rendering, vintage 1970s film grain and fade, classic 35mm black-and-white film, sepia old-photo vignette, chrome-sphere reflection world, anamorphic lens flare, film-noir high-contrast lighting. |
| **Particle & Element Effects** | Confetti shower from above, dandelion seeds drift by, white feathers fall like snow, rainbow-hued mist swirls, fireflies blink in dark areas, countless soap bubbles float, golden coins rain down, glitter rains from halo overhead, golden sparkles twinkle throughout, cherry blossom petals whirl, glowing monarch butterflies swarm, pixie-dust trails swirl magically. |
| **Color Palette & Tonal Changes** | Deep monochrome blue grade, oversaturated vaporwave palette, selective color (B&W except vivid red), radioactive sickly-green tint, sun-bleached faded desert look, autumnal warm grade (oranges/reds/browns), teal-and-orange blockbuster grade, high-key bright airy palette, cool wintry grade (blues/whites/grays), twilight blue after sunset, warm golden-hour glow, soft morning side-light, fluorescent clinical lighting, candlelight illumination, bonfire-lit scene with flicker, scene lit by bioluminescence, cyberpunk neon glow with digital rain. |
| **Surreal & Fantasy Transformations** | Floating green-topped islands fill the sky, giant translucent spirit animals in background, overgrown with glowing bioluminescent fungi, two-sun alien planet landscape, holographic projection aura overlays scene, world fractures with spreading glass-like cracks, everything but person crystallizes, environment melts like a Dalí painting, background becomes swirling galaxy nebula, candy-themed fantasy world outside, impossible Escher-like architecture, gravity inverts and small objects float upward, person levitates into meditation pose, ground becomes still mirror-water, colors invert to negative film, ethereal mid-ground mist drifts through, shadow detaches and dances playfully, sticky notes stack up with bad cat doodles, psychedelic rainbow color-shifting sweep, psychedelic swirling liquid-paint background, scene inside a giant snow globe, neon city-night transition (cyberpunk), gritty graphic stylization with bold outlines. |

Table 6. Grouped catalog of video effects. We provide a set of general categories for our video to video effects as described in Figure 4 of the main text. We provide a paraphrased version of each effect, as the actual prompts are long. .

| Grounding Objects |
| --- |
| a a woman, a person, a woman, baby, bookshelf, person, child, children, delivery person, eye, eye mask, face, feet, foot, glasses, hair, hand, hands, head, man, nose ring, old man, old woman, pathway person, pathway, person, people, person, person man, person person, person woman, pregnant woman, skateboard, person, the person, watch, person, woman, women, young man, young people, young woman. |

Table 7. List of grounding objects kept from Senorita dataset grounding subset.

| Effect | Hyperparameters |
| --- | --- |
| Posterize Frames | color palette (30 options; 4/6/8 colors each), buckets = palette size |
| Pixelate Frames | pixel length: {4, 8, 16, 32} |
| Invert Frames | (no hyperparameters) |
| Wave Warp | amplitude: [10, 50], frequency: [0.01, 0.075] |
| Update Saturation Brightness | brightness: [-255, 254], saturation: {0.0, 0.5, 1.0, 1.5, 2.0} |
| Gaussian Blur | kernel size: {(11,11), (21,21), (51,51), (101,101), (151,151)} |
| Add Grain | amount: [10, 50], grain size: {1, 2, 3, 4, 5} |
| Black And White | (no hyperparameters) |
| Color Overlay | percent overlay: {0.25, 0.5, 0.75, 1.0}, color (RGB 0–255 each) |
| CC Ball Action | grid spacing: {5, 10, 20, 30}, ball color (RGB 0–255 each) |
| Sticker Effect | border size: {10, 20, 30, 40, 50}, sticker color (RGB 0–255 each) |
| Glow Effect | glow size: {10, 20, 30, 40, 50}, glow color (RGB 0–255 each), glow brightness: {0.5, 1.0, 1.5, 2.0}, object brightness: {0, 1, 2, 3, 4, 5} |
| Radial Blur | center: (0.5, 0.5) or uniform [0.10, 0.90] (2dp), strength: [5, 60], blur border (% of corner distance): [0, 30] |
| Rotate Pixels | center: uniform [0.000, 1.000] (3dp), max angle: {10, 20, 30, 40, 50}, radius (px): {100, 200, 300, 400, 500} |
| Glitch Effect | angle: [0, 359], red displacement: [-40, 40], green displacement: [-40, 40], blue displacement: [-5, 5] |
| Dither | size: {5, 8, 10, 12, 16, 20}, color steps per channel: {2, 3, 4, 5, 7, 9} |
| Motion Blur | angle: [0, 360], strength: [0.1, 0.6] |
| Stutter | hold duration (frames): {1, 2}, stutter frequency (every Nth frame): {3, 4, 5, 6} |
| Ghosting | ghost intensity: [0.1, 0.9] |
| Strobe | flash frequency (frames): {4, 6, 8, 10}, flash duration (frames): {1, 2}, flash color: {white, black} |

Table 8. Code-based video effects and their hyperparameters. Ranges indicate the values sampled or supported by the provided config generators.

| Temporal Effect | Hyperparameters |
|---|---|
| Alpha Blend | softness: [0.01, 0.05], center: (0.5, 0.5) or uniform [0.00, 1.00] (2dp), transition window (frames): start and end in [0, 33] |
| Wipe Left To Right | softness: [0.01, 0.05], center: (ignored), transition window (frames): start and end in [0, 33] |
| Wipe Right To Left | softness: [0.01, 0.05], center: (ignored), transition window (frames): start and end in [0, 33] |
| Wipe Top To Bottom | softness: [0.01, 0.05], center: (ignored), transition window (frames): start and end in [0, 33] |
| Wipe Bottom To Top | softness: [0.01, 0.05], center: (ignored), transition window (frames): start and end in [0, 33] |
| Diag Top Left Bottom Right | softness: [0.01, 0.05], center: (ignored), transition window (frames): start and end in [0, 33] |
| Diag Bottom Right Top Left | softness: [0.01, 0.05], center: (ignored), transition window (frames): start and end in [0, 33] |
| Diag Top Right Bottom Left | softness: [0.01, 0.05], center: (ignored), transition window (frames): start and end in [0, 33] |
| Diag Bottom Left Top Right | softness: [0.01, 0.05], center: (ignored), transition window (frames): start and end in [0, 33] |
| Circle Out | softness: [0.01, 0.05], center: (0.5, 0.5) or uniform [0.00, 1.00] (2dp), transition window (frames): start and end in [0, 33] |
| Circle In | softness: [0.01, 0.05], center: (0.5, 0.5) or uniform [0.00, 1.00] (2dp), transition window (frames): start and end in [0, 33] |
| Rect Out | softness: [0.01, 0.05], center: (0.5, 0.5) or uniform [0.00, 1.00] (2dp), transition window (frames): start and end in [0, 33] |
| Rect In | softness: [0.01, 0.05], center: (0.5, 0.5) or uniform [0.00, 1.00] (2dp), transition window (frames): start and end in [0, 33] |
| Diamond Out | softness: [0.01, 0.05], center: (0.5, 0.5) or uniform [0.00, 1.00] (2dp), transition window (frames): start and end in [0, 33] |
| Diamond In | softness: [0.01, 0.05], center: (0.5, 0.5) or uniform [0.00, 1.00] (2dp), transition window (frames): start and end in [0, 33] |

Table 9. Temporal transition effects and their hyperparameters. Note that start frame < end frame always when deciding transition window

| | I2V | | | Neural V2V | | | Code Based V2V | | |
|---|---|---|---|---|---|---|---|---|---|
| Method | AES | Motion | Dynamic | AES | Motion | Dynamic | AES | Motion | Dynamic |
| Ours | 0.5607 | 0.975 | 0.711 | 0.5649 | 0.969 | 0.712 | 0.4802 | 0.984 | 0.143 |
| Wan 2.1 | 0.5559 | 0.977 | 0.737 | - | - | - | - | - | - |
| Wan VACE I2V | 0.5479 | 0.974 | 0.816 | - | - | - | - | - | - |
| Wan VACE Pose | - | - | - | 0.5549 | 0.968 | 0.817 | 0.5270 | 0.992 | 0.190 |
| Wan VACE Depth | - | - | - | 0.5977 | 0.975 | 0.421 | 0.5532 | 0.994 | 0.048 |
| Lucy Edit | - | - | - | 0.5075 | 0.987 | 0.150 | - | - | - |
| Validation | 0.5557 | 0.972 | 0.765 | 0.6435 | 0.976 | 0.750 | 0.4922 | 0.985 | 0.190 |

Table 10. Individual video statistics across methods.