

# Scaling Instruction-Based Video Editing with a High-Quality Synthetic Dataset

## Supplementary Material

### A. Overview

In this appendix, we provide additional details and results to supplement our main paper. In Sec. B and Sec. C, we present comparisons of data pipelines and the analysis of the in-context data generator. Sec. D contains additional quantitative and qualitative results of our dataset and model, and a word cloud visualizing the distribution of instructions in our dataset. Finally, we showcase the interface used for our human evaluations in Sec. E. We strongly recommend reviewing the video samples in “[index.html](#)” in Supplementary Materials for a better understanding.

### B. Comparison of Data Pipelines

We compare the data synthesis and filtering strategies of Ditto with InsViE and Señorita in Table S1. Regarding the synthesis pipeline, InsViE relies on I2V inversion which is inherently slow and costly, while Señorita necessitates the complex training and maintenance of 19 distinct expert models. In contrast, Ditto employs all-in-one models, offering a significantly more efficient and scalable generation process. In terms of data filtering, we address the limitations of Señorita’s CLIP-based approach by introducing tracking + VLM protocol. Specifically, our tracking mechanism precisely filters out samples with insufficient motion, ensuring superior dynamic range compared to Señorita, while the VLM rigorously guarantees editing quality and safety. For concrete visual comparisons demonstrating these advantages, we strongly recommend readers refer to the samples provided on our [HTML page](#).

Table S1. Comparison of Data Synthesis Pipelines. We contrast the synthesis models and filtering strategies.

Method	Synthesis Pipeline	Data Filtering
InsViE	I2V Inversion	Optical Flow + VLM
Señorita	Multiple Expert Models	CLIP
Ditto	All-in-one Models	Tracking + VLM

### C. Justifying Data Generation Pipeline Design

We provide an analysis of the videos synthesized by the in-context video generator, VACE, to justify the design of our data pipeline. As in Fig. S1, we first observe that using only depth maps to guide the generator results in a significant loss of content from the source video, leading to poor fidelity. Conversely, conditioning the generator on a keyframe from the original source video alongside the edit-

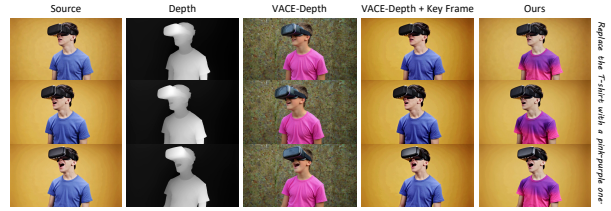


Figure S1. Results of various settings for data generation.

ing instruction fails to produce the desired edit - the output remains almost identical to the source. These findings reveal that while the base generator excels at motion transfer, its inherent instruction-following capability for editing tasks is limited. Based on this analysis, we validate our proposed approach: using a keyframe modified by an advanced image editor, in conjunction with depth guidance as the context. This method achieves the optimal balance of instruction adherence, temporal consistency, and source fidelity for our data synthesis.

### D. Additional Results of Dataset and Model

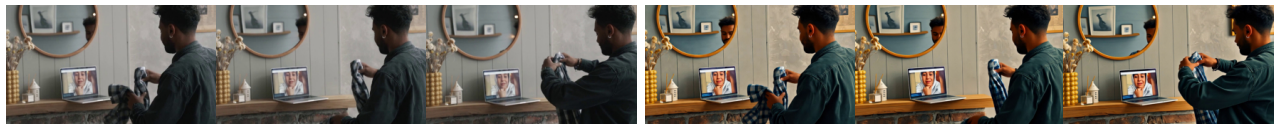
Table S2. Ablation on modality curriculum learning (MCL).

Configuration	CLIP-T $\uparrow$	CLIP-F $\uparrow$	VLM $\uparrow$
w/o MCL	24.79	98.91	7.60
w/ MCL	<b>25.54</b>	<b>99.03</b>	<b>8.10</b>

We first conduct ablation study on modality curriculum learning (MCL) in Tab. S2 to validate its effectiveness. We then include results of Vision-Language Model (VLM) data filtering in Fig. S2, as well as the local editing data on adding and removal in Fig. S3. We present additional qualitative results to further demonstrate our dataset and model’s performance across a wide range of editing instructions in Fig. S4. We also include a word cloud in Fig. S5 that illustrates the diversity and distribution of the editing prompts within the dataset, highlighting its coverage.

### E. Demonstration of User Study Interface

To collect human preference data, we designed a user-friendly evaluation interface, as shown in Fig. S6. For each source video and text prompt, we presented participants with the edited results from different methods in a randomized order. They were then asked to rank the videos from best (1) to worst (4) based on three criteria: Instruction Following, Temporal Consistency, and Overall preference. The final scores are calculated based on the ranking.



*Instruction: Transform the scene into a digital portrait painting style. VLM: X Instruction Adherence Issue.*



*Instruction: Transform the scene into a surreal, abstract animation with the waterfall. VLM: X Severe Spatial and Structural Misalignment.*

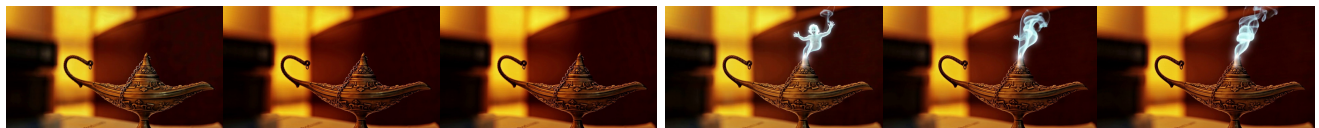


*Instruction: Recreate the video in a cinematic film noir style with heavy shadows, grainy textures, and a desaturated color palette. VLM: X Video is too dark.*

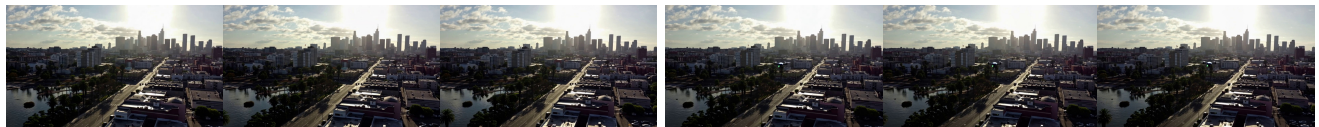


*Instruction: Make it entirely out of origami. Its form is composed of clean, sharp paper folds and crisp, geometric creases. VLM: X Too similar to source.*

Figure S2. Data filtering results from the Vision-Language Model (VLM). We design a set of rules for the VLM to guide its data filtering process. Experiments demonstrate that Large VLMs can reliably detect semantic and other types of failures.



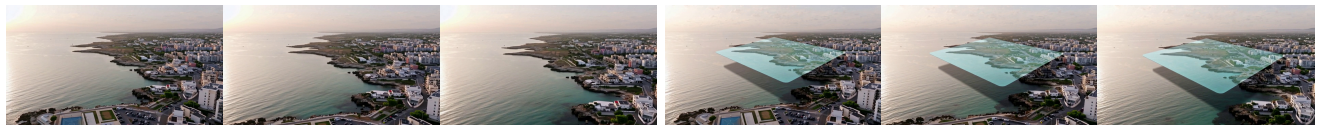
*Add / Remove the small floating genie spirit made of smoke and light emerging from the lamp's spout.*



*Add / Remove the small, glowing drone flying slowly between the buildings in the foreground.*



*Add / Remove the vintage-style pocket watch lying open on the table near her hand.*



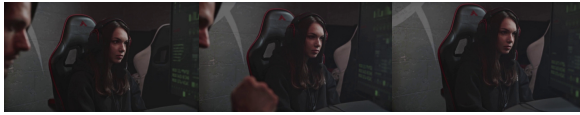
*Add / Remove the floating, translucent map overlay.*

Figure S3. Additional visualization of local editing data on adding and removing. Given the depth corresponding to the source video, our pipeline can stably synthesize data for addition-type edits, which are then *reversely* used as training data for removal-type editing.

Data Visualization



Apply an overgrown effect, as if nature is reclaiming everything. Cover all surfaces with patches of lush moss and intricate networks of fine vines. Tiny wildflowers and small mushrooms sprout from crevices. The air seems humid, with dappled light filtering through an unseen canopy.



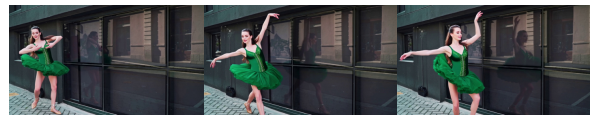
Apply a deep underwater effect. The entire scene is bathed in a moody, deep blue and cyan light.



Make it sculptured with glass.



Infuse with the visual elements of the 3D Chibi style.



Change the light-colored bodice to a deep emerald green with metallic thread details.

Model Performance



Make it Cyberpunk style.



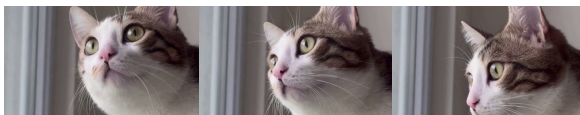
Make it a Claymation.



Overlay the video with an intense fire effect.



Make it the style of Spider-Man Into the Spider-Verse.



Turn the cat into a black cat.

Figure S4. Additional visualization of data from the proposed dataset and model outputs. Please view the site in the supplementary materials for additional video samples.

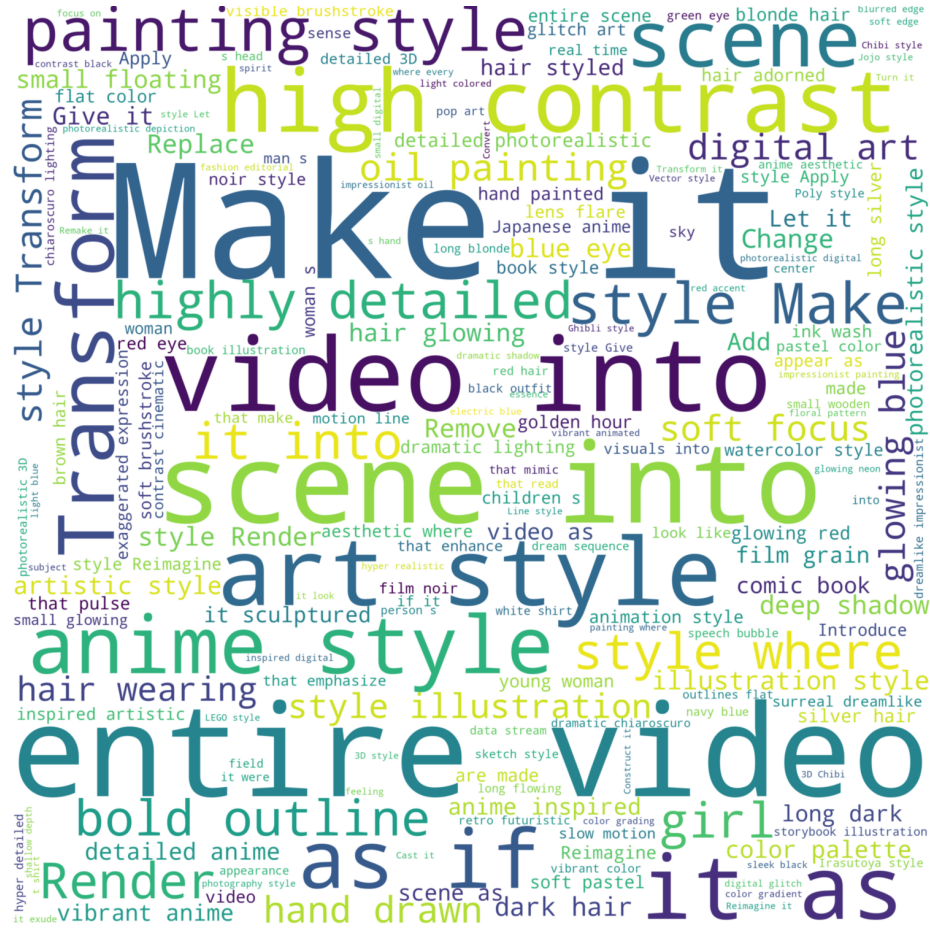


Figure S5. The word cloud of editing instructions.

### DITTO User Study

Enter Your Name

---

#### Evaluation Criteria



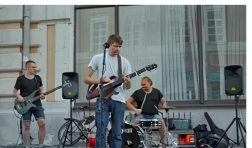


**Instruction Following:** How well the generated video follows the editing instructions.

**Temporal Consistency:** How well the generated video coherence between video frames.

**Overall:** Your overall preference for the generated video based on all factors.

---

Prompt: Let it be like the 3D Chibi style.

	Video A	Video B	Video C	Video D
Source				
Instruction Following				
Temporal Consistency				
Overall				

Ranking (Enter values 1, 2, 3, or 4) [best = 1, worse = 4]

Figure S6. The interface of the user study.