

# Unified Video Editing as Temporal Reasoner

Anonymous CVPR submission

Paper ID 957



**Figure 1.** VideoCoF’s video editing capabilities emerge from its **seeing, reasoning, then editing framework**. Trained on only **50k** data (33 frames), this teaser shows multi-instance editing and robust 4x+ length generalization.

## Abstract

001      Existing video editing methods often rely on task-specific  
 002      priors, such as masks, which hinders the development of a  
 003      universal, general-purpose framework. While temporal in-  
 004      context learning offers a potential path to unification, it of-

ten lacks explicit spatial cues, leading to weak instruction-to-region mapping and imprecise localization. Motivated by these limitations and inspired by Chain-of-Thought reasoning, we propose (**VideoCoF**), a novel **Chain-of-Frames** approach. VideoCoF enforces a “see → reason → edit” procedure by compelling the video diffusion model to first pre-

005  
006  
007  
008  
009  
010

dict reasoning tokens (*edit-region latents*) before generating the target video tokens. This explicit reasoning step removes the need for user-provided masks while achieving precise instruction-to-region alignment and fine-grained video editing. Furthermore, we introduce a RoPE alignment strategy that leverages these reasoning tokens to ensure motion alignment and enable length extrapolation beyond the training duration. We demonstrate that with a minimal data cost of only 50k video pairs, VideoCoF achieves state-of-the-art performance on VideoCoF-Bench, validating the efficiency and effectiveness of our approach.

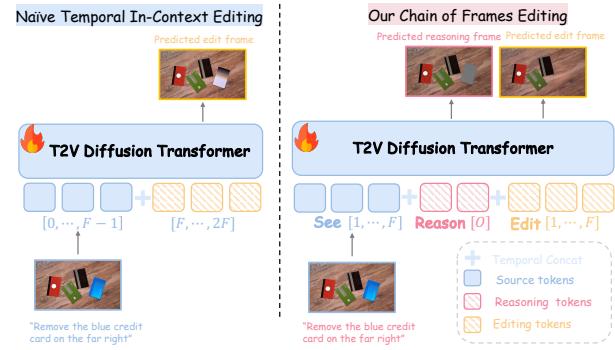
## 1. Introduction

The development of Video Diffusion Models (VDM) [12, 27, 33, 37] has enabled high-fidelity video generation across a wide range of concepts. Building on these advances, video editing methods support users in designing video by adding [26], removing [15, 43], swapping [6, 36] visual concepts, and performing global style transformation [39].

Early video editing approaches [1, 15, 26, 36, 41] typically rely on adapter-based control modules to inject spatial or appearance priors into a diffusion backbone. While accurate, these approaches require user-provided masks or references, incurring manual effort and external supervision, and limiting scalability to a unified framework.

With the scaling up of video diffusion models [27], a temporal in-context learning ability has emerged, offering a potential path [9, 16, 38] to unify these disparate editing tasks. These methods concatenate source-video tokens with noised edit tokens along the temporal dimension and use self-attention to learn from the source to guide the edit. However, this naive temporal in-context learning method (e.g., ICVE [16]), learns from video context without extra inputs, lacking explicit spatial cues for the edit region, resulting in weak instruction-to-region grounding and imprecise localization. Consequently, as shown in Fig. 2 left, the model fails to disambiguate the “blue card” on the far right, thus cannot perform the desired fine-grained editing.

Motivated by these complementary strengths and weaknesses, we seek a unified method that resolves these problems jointly. Inspired by Chain-of-Thought (CoT) multi-step reasoning [30], we compel the video diffusion model to first predict the edit region and then perform the edit, enforcing a “see → reason → edit” procedure. Accordingly, we propose **VideoCoF**, a Chain-of-Frames approach that predicts *reasoning tokens* (edit-region latents) before generating the target video tokens, thereby removing the need for user-provided masks while achieving precise instruction-to-region alignment. To explicitly model the reasoning process, we leverage visual grounding, which is naturally suited to simulating reasoning about the edit region. Empirically, we find a soft, gradually highlighted grayscale region



**Figure 2.** Illustration of the difference between previous methods and our VideoCoF. We enhance the editing accuracy by forcing the video diffusion model to first predict the editing area, and then perform the editing.

is the most effective reasoning format. Additionally, we introduce a RoPE alignment strategy. By explicitly accounting for the reasoning latent, we reset the temporal indices of the edited video’s rotary position embeddings to match those of the source segment, ensuring motion alignment and length extrapolation.

Our contributions can be summarized as follows:

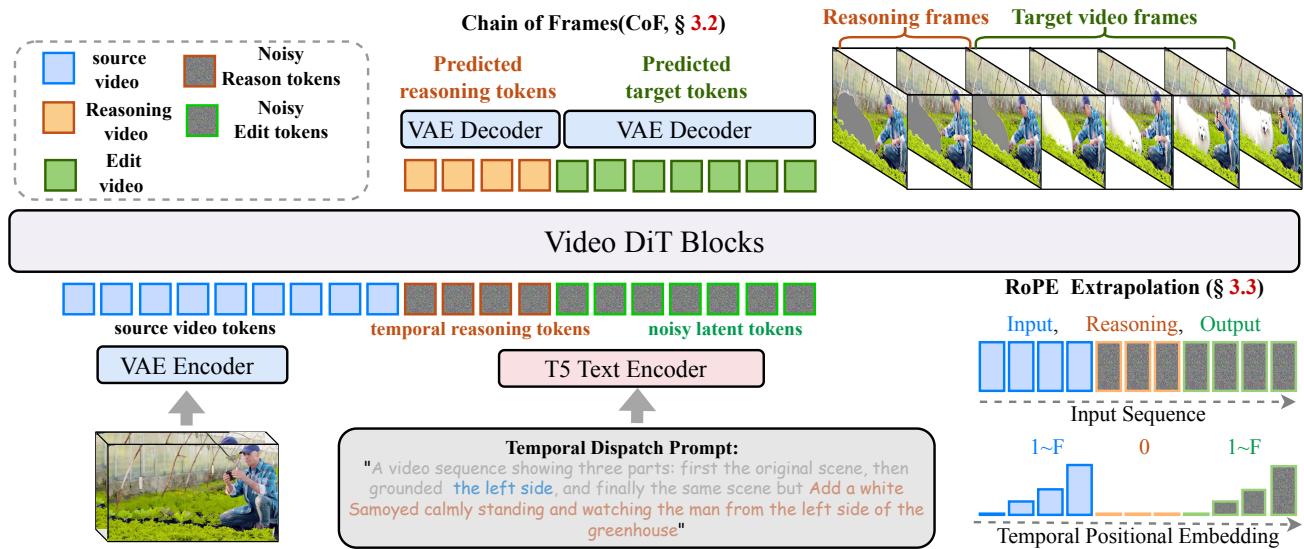
- We propose VideoCoF, the first framework to introduce a Chain of Frames approach to video editing, enabling temporal reasoning for fine-grained video editing.
- Building on VideoCoF, we explore an effective reasoning format for video diffusion models, and introduce a RoPE alignment strategy that allows generalization to longer frames exceeding the training duration.
- We demonstrate that with a minimal data cost (only 50k video pairs), we achieve state-of-the-art quantitative and qualitative performance on VideoCoF-Bench, validating the efficiency and effectiveness of our approach.

## 2. Related Work

**Video Editing Methods** Early training-free video editing methods [21, 33] used inversion and consistency techniques (e.g., attention manipulation [21] or optical flow [5]) but often lack precise control and struggle with complex edits. Data-driven, training-based methods [2, 4] have become the focus, offering higher quality and edit diversity. A concurrent line of research [18, 29, 40] integrates MLLMs to guide the editing process, though this adds significant training and inference cost, which our pure VDM approach avoids.

**In-Context Video Editing** Recently, in-context learning (ICL) has emerged as a promising paradigm for unified editing [10, 35, 42]. Methods like UNIC [38] and ICVE [16] concatenate video conditions along the temporal axis to perform ICL. However, these methods are often limited by mask requirements [38] or, as we identify, suffer from fundamental issues with editing accuracy and a lack of

062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096



**Figure 3. Overview of VideoCoF framework.** Our model processes source (blue), reasoning (orange), and target (green) tokens in a unified sequence to “reason” then “edit”. **Bottom right:** Our RoPE design enables length extrapolation.

length extrapolation due to their naive temporal concatenation. While EditVerse [9] also explored unified in-context learning, it was built on a LLaMA-style DiT backbone, whereas our work explores these capabilities within a standard video diffusion transformer.

**Chain of Thought in Vision** Chain-of-Thought (CoT) prompting [11, 30] elicits multi-step reasoning in LLMs by having them “think step-by-step.” This concept of emergent reasoning has also been identified in large video generative models [3, 31] that can solve visual puzzles. However, how to leverage visual reasoning for the task of unified video editing remains unexplored. In this work, we investigate whether generative video models can perform a “chain of frames” reasoning to achieve this.

### 3. Methods

#### 3.1. VideoCoF Framework

As illustrated in Figure 3, VideoCoF employs a VideoDiT [27] for unified video editing. We model editing as a reasoning-then-generation process: the model first reasons where to edit, then generates the intended content in that area. We call this process “**Chain of Frames (CoF)**” (Sec 3.2). All visual inputs (source, reasoning, and target frames) are encoded separately by a Video VAE and then concatenated temporally. The unified frame sequence is then fed into the model, performing unified in-context learning via self-attention and language control via cross-attention. To enable video alignment and variable-length inference, we revisit the design of positional encoding. We adapt the temporal RoPE for source-to-target alignment and reasoning tokens’ RoPE for explicit spatial guidance (Sec 3.3). Subsequent sections detail the insights behind our design choices

and data curation pipeline (Sec 3.5).

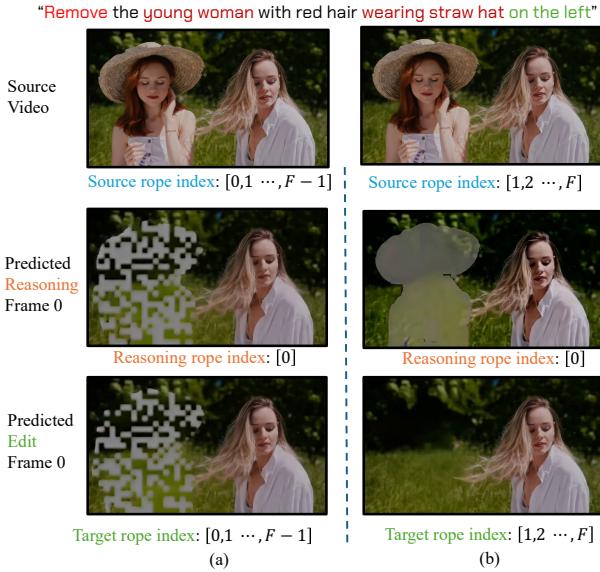
#### 3.2. Chain of Frames

**Seeing, Reasoning, then Editing** Previous video in-context editing methods, such as UNIC [38], ICVE [16], or EditVerse [9], perform in-context learning by temporally concatenating clean source video tokens with noised editing video tokens. However, this approach lacks an explicit constraint mapping the editing instruction to the specific editing region, leading to editing accuracy problems, as shown in Fig 2. Recently, VDM have been shown to possess reasoning capabilities, as demonstrated in [31]. Inspired by this, we explicitly model the reasoning tokens, forcing the model to actively learn the relationship between the editing instruction and the target edit region first. The edit is then executed *after* reasoning, following a “see, reason, then edit” process.

Inspired by Chain of Thought prompting in Large Language Models (LLMs) [30], we argue that a video generative model should also have an analogous chain-reasoning ability. Given the generative priors in video editing, the visual-chain should be progressive, moving from the original video to a visual reference of the editing region, and finally to the edited video. Visual grounding is naturally suitable for this representation. Since video diffusion models are often insensitive to grounding masks (black or white pixels). Therefore, we choose to use a gray highlight to delineate the “grounding region,” which is also evidence in [7]. Finally, the gray-highlighted area as the ground truth for the reasoning frames, teaching the diffusion model to reason about where the edit should occur.

Consequently, the entire video editing task is reformulated as a chained process: first “seeing” the original video,

128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159



**Figure 4.** How our RoPE design avoid index collision.

187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205

then “reasoning” by predicting the grounding region, and finally “editing” to generate the new video content within that specified area. We call this **Chain of Frames (CoF)**.

Let  $\mathcal{E}(\cdot)$  denote the video VAE encoder. We use  $F$  and  $L$  for frames in the source/target and reasoning latent space, respectively, and denote channel, height, and width by  $C$ ,  $H$ , and  $W$ . Given a triplet source-reasoning-target video pair  $\{\mathbf{s}, \mathbf{r}, \mathbf{e}\}$ , we first encode them into latent representations. The source  $\mathbf{s}$  and target video  $\mathbf{e}$  yield latent  $z_s = \mathcal{E}(\mathbf{s})$  and  $z_e = \mathcal{E}(\mathbf{e})$ , both with shape  $\mathbb{R}^{F \times C \times H \times W}$ . The reasoning video  $\mathbf{r}$  yields a latent  $z_r = \mathcal{E}(\mathbf{r})$  with shape  $\mathbb{R}^{L \times C \times H \times W}$ . This separate encoding ensures intra-causal relations and inter-video independence. Then, we perform temporal concatenation to get the unified latent representation:

206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218

$$\mathbf{z}_{full}^{(t)} = \underbrace{z_s^{(0)}}_{\text{seeing}} \parallel \underbrace{z_r^{(t)}}_{\text{reasoning}} \parallel \underbrace{z_e^{(t)}}_{\text{editing}} \in \mathbb{R}^{(F+L+F) \times C \times H \times W}, \quad (1)$$

Given a concatenated full latent sequence  $\mathbf{z}_{full}^{(0)} = \text{TemporalConcat}(\mathbf{z}_s^{(0)}, \mathbf{z}_r^{(0)}, \mathbf{z}_e^{(0)})$ , we treat the reasoning+editing block as the generation target during training.

Given timestep  $t \in [0, 1]$  and Gaussian noise  $\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , we only progressively noise the reasoning and editing parts,  $\mathbf{z}_{r,e}^{(t)} = (1-t)(\mathbf{z}_r^{(0)} \parallel \mathbf{z}_e^{(0)}) + t\varepsilon_{F:2F+L-1}$ , and form the model input  $\mathbf{z}^{(t)} = \mathbf{z}_s^{(0)} \parallel \mathbf{z}_{r,e}^{(t)}$ . The target velocity field is  $\mathbf{v} = \varepsilon - \mathbf{z}_{full}^{(0)}$ . Our model  $\mathbf{F}_\theta(\cdot)$  predicts this velocity field from the partially noised input, and we train it by minimizing the mean squared error between predicted and true velocities. Concretely, we only supervise the reasoning and target frames, so the training loss can be written in per-frame form as

$$\mathcal{L} = \frac{1}{L+F} \sum_{i=F}^{2F+L-1} \left\| \mathbf{v}_i - [\mathbf{F}_\theta(\mathbf{z}^{(t)}, t, \mathbf{c})]_i \right\|_2^2, \quad (2)$$

ders video length extrapolation, as the model overfits to a static  $[0, F-1] \rightarrow [F, 2F-1]$  mapping and fails to generalize to videos longer than  $F$  frames.

A better strategy is to repeat the temporal indices. For our CoF triplet (consider  $L=1$  for reasoning frame), a straightforward reset configuration is to assign temporal indices:  $[0, F-1]$  to the source, “0” to the reasoning frame, and  $[0, F-1]$  to the target.

However, as illustrated in Figure 4 (a), this naive reset leads to index collisions at temporal position 0, shared by the source, reasoning, and target frames. This overlap introduces visual artifacts that propagate from the reasoning tokens into the first target frame.

To resolve this index collision, we set the temporal indices for both the source video and the target video to the range  $[1, F]$ , while keeping the reasoning frame’s temporal index at 0. This isolates the reasoning token and prevents artifact leakage while maintaining length generalization.

### 3.4. Training and Inference Paradigm

#### Algorithm 1 Chain of Frame (CoF) Training

**Input:** Dataset  $\mathcal{D}$  with tuples  $(\mathbf{z}_s^{(0)}, \mathbf{z}_r^{(0)}, \mathbf{z}_e^{(0)}, \mathbf{c})$

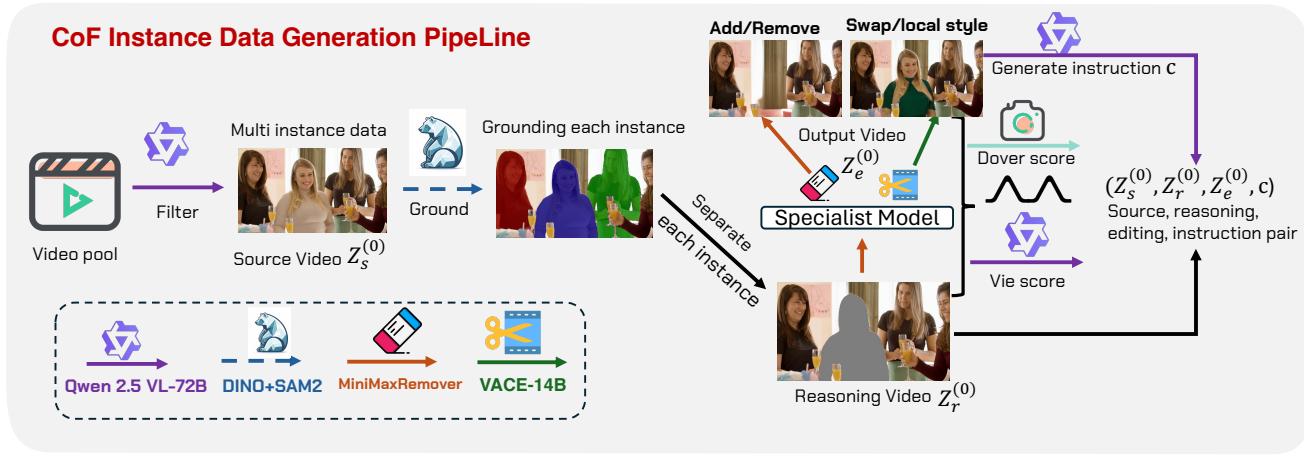
**Output:** Fine-tuned parameters  $\theta$

```
foreach minibatch  $(\mathbf{z}_s^{(0)}, \mathbf{z}_r^{(0)}, \mathbf{z}_e^{(0)}, \mathbf{c}) \sim \mathcal{D}$  do
    foreach sample in minibatch do
         $\mathbf{z}_{full}^{(0)} \leftarrow \mathbf{z}_s^{(0)} \parallel \mathbf{z}_r^{(0)} \parallel \mathbf{z}_e^{(0)}$  Sample  $t \sim \mathcal{U}[0, 1]$ 
        Sample  $\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  with the same shape  $\mathbf{z}_{full}^{(0)}$ 
         $\mathbf{v} \leftarrow (\varepsilon - \mathbf{z}_{full}^{(0)})$ 
         $\mathbf{z}_{r,e}^{(t)} \leftarrow (1-t)(\mathbf{z}_r^{(0)} \parallel \mathbf{z}_e^{(0)}) + t\varepsilon_{F:2F+L-1}$ 
         $\mathbf{z}^{(t)} \leftarrow \mathbf{z}_s^{(0)} \parallel \mathbf{z}_{r,e}^{(t)}$ 
         $\hat{\mathbf{v}} \leftarrow \mathbf{F}_\theta(\mathbf{z}^{(t)}, t, \mathbf{c})$ 
         $\mathcal{L} \leftarrow \frac{1}{L+F} \sum_{i=F}^{2F+L-1} \|\mathbf{v}_i - \hat{\mathbf{v}}_i\|_2^2$ 
    Update  $\theta$  using gradients of  $\mathcal{L}$ 
```

Given a concatenated full latent sequence  $\mathbf{z}_{full}^{(0)} = \text{TemporalConcat}(\mathbf{z}_s^{(0)}, \mathbf{z}_r^{(0)}, \mathbf{z}_e^{(0)})$ , we treat the reasoning+editing block as the generation target during training.

Given timestep  $t \in [0, 1]$  and Gaussian noise  $\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , we only progressively noise the reasoning and editing parts,  $\mathbf{z}_{r,e}^{(t)} = (1-t)(\mathbf{z}_r^{(0)} \parallel \mathbf{z}_e^{(0)}) + t\varepsilon_{F:2F+L-1}$ , and form the model input  $\mathbf{z}^{(t)} = \mathbf{z}_s^{(0)} \parallel \mathbf{z}_{r,e}^{(t)}$ . The target velocity field is  $\mathbf{v} = \varepsilon - \mathbf{z}_{full}^{(0)}$ . Our model  $\mathbf{F}_\theta(\cdot)$  predicts this velocity field from the partially noised input, and we train it by minimizing the mean squared error between predicted and true velocities. Concretely, we only supervise the reasoning and target frames, so the training loss can be written in per-frame form as

219



**Figure 5.** Our data curation pipeline for multi-instance data.

where  $[\mathbf{F}_\theta(\mathbf{z}^{(t)}, t, \mathbf{c})]_i$  denotes the model’s prediction for frame  $i$  and  $\mathbf{c}$  is the text condition. The model parameters  $\mathbf{F}_\theta(\cdot)$  are updated via a gradient step computed from this loss. The full training procedure is summarized in Algorithm 1.

During inference we initialize the reasoning+editing block from Gaussian noise,  $\mathbf{z}_{r,e}^{(1)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and form the full latent at  $t = 1$  by temporal concatenation with the clean source  $\mathbf{z}_{\text{full}}^{(1)} = \text{TemporalConcat}(\mathbf{z}_s^{(0)}, \mathbf{z}_{r,e}^{(1)})$ . An ODE solver guided by our model  $\mathbf{F}_\theta$  evolves  $\mathbf{z}_{\text{full}}^{(1)}$  to  $\mathbf{z}_{\text{full}}^{(0)}$ . The source latents  $\mathbf{z}_s^{(0)}$  are held fixed during inference, so only the reasoning/editing parts change. We then extract the edited-target latent using the same slicing index as in training:  $\mathbf{z}_{\text{edit}}^{(0)} = (\mathbf{z}_{\text{full}}^{(0)})_{F+L:2F+L-1}$  and decode the final edited video:  $\mathbf{x}_{\text{edit}} = \mathcal{D}(\mathbf{z}_{\text{edit}}^{(0)})$ .

### 3.5. Video Data Curation

The training of our VideoCoF requires a large and diverse dataset structured as source, reasoning, and edited video triplets. However, existing video editing datasets and methods predominantly focus on single-instance-level object manipulation. This limitation is a significant barrier, as real-world videos contain complex visual cues, multiple interacting instances, and intricate spatial relationships (e.g., physical left/right, object-to-object interactions). Enabling a generative model to comprehend these complex, instance-level dynamics is a critical step toward true reasoning-based video editing. Therefore, we develop a comprehensive data curation pipeline, illustrated in Figure 5, to specifically generate and process complex, instance-level video data.

**Instance-Level Curation Pipeline.** Our pipeline begins with a large pool of diverse videos sourced from Pexels [20]. First, we employ the Qwen-VL 72B [28] to perform multi-instance identification, scanning the videos to find scenes that contain multiple, distinct objects. Once these videos are identified, we use Grounding-SAM2 [22]

to perform precise segmentation, generating distinct segmentation masks for each individual instance.

With these instance-specific masks, we generate triplets for a variety of editing tasks:

- **Object Addition/Removal:** We utilize the Minimaxremover [43] to erase a specific instance from the video. The data for object addition is then created by simply reversing this process.
- **Object Swap and Local Style Transfer:** For these tasks, we leverage the VACE 14B [8] in its inpainting mode to fill the specified masked regions. Critically, the creative prompts for these inpainting edits are generated by GPT-4o [19], as we found Qwen-VL 72B’s imaginative capabilities for this specific task to be limited.

**Filtering and Final Dataset.** All generated video pairs are rigorously evaluated to ensure quality. We use the Dover Score [32] to assess aesthetic quality and the VIE Score [13] to measure editing fidelity and coherence. A weighted combination of these scores is used to filter for high-quality, successful edits. Finally, we use this pipeline to filter from the large-scale open-source Señorita 2M [44] dataset, and distill a high-quality subset of **50k** videos to supplement our training data. This multi-pronged approach yields our final large-scale dataset, rich in the instance-level complexity required for reasoning-based video editing.

## 4. Experiments

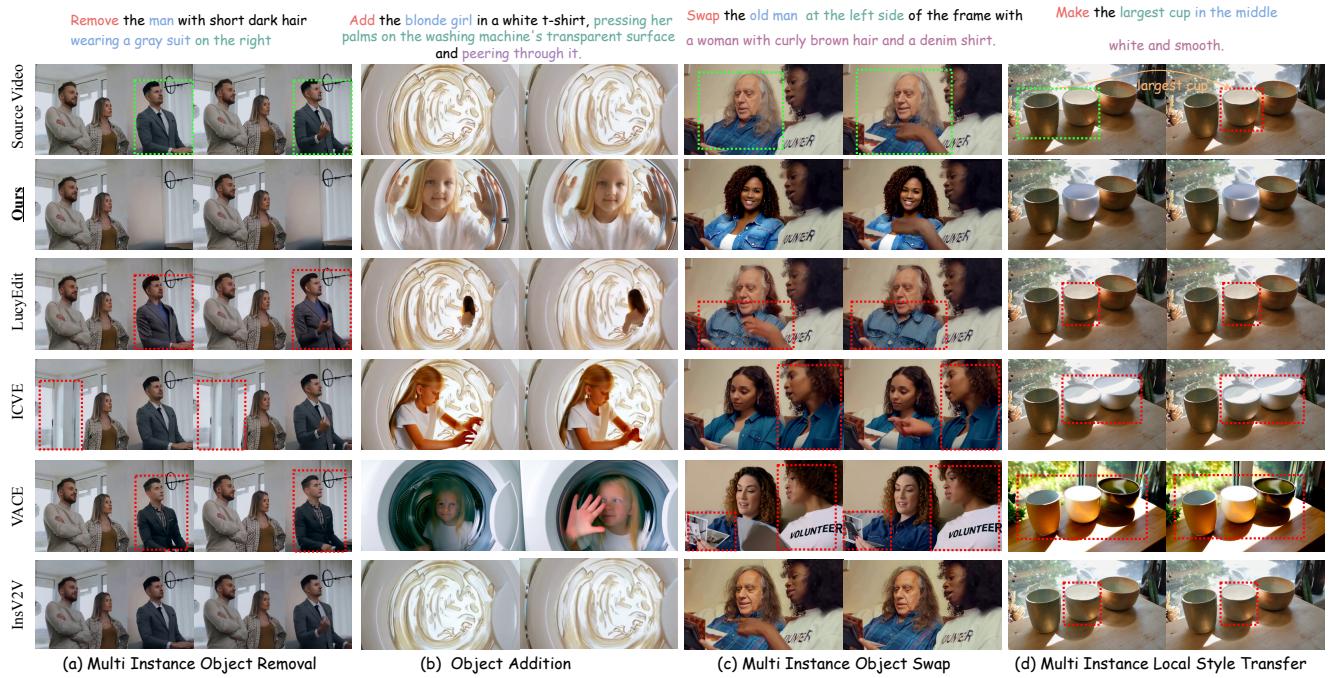
### 4.1. Implementation Details.

VideoCoF is trained on WAN-14B [27]. We employ a resolution-bucketing strategy to support multiple aspect ratios, using spatial resolutions of 336×592, 400×704, 400×752, and 400×944 (and the corresponding vertical variants, e.g., 592×336). Training videos are sourced from Señorita [44] and are 33 frames long, we only training on **50k** curated video data finally. Thanks to our RoPE alignment design, the model generalizes to longer sequences at inference (e.g., **141 frames and above**). By default we

255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279

Model	GPT-4o Score (avg.)				Perceptual Quality (avg.)		
	Instruct Follow↑	Preservation↑	Quality↑	Success Ratio↑	CLIP-T↑	CLIP-F↑	DINO↑
InsV2V [4]	3.41	6.15	5.51	6.39%	26.19	0.988	0.978
Señorita [44]	3.26	6.30	5.48	10.35%	26.04	0.994	0.988
VACE [8]	7.47	5.82	7.61	26.60%	27.02	0.994	0.990
ICVE [16]	7.79	8.06	<b>8.14</b>	57.76%	27.49	0.992	0.986
Lucy Edit [25]	5.24	6.50	6.37	29.64%	26.98	0.991	0.986
<b>VideoCoF (Ours)</b>	<b>8.97</b>	<b>8.20</b>	<u>7.77</u>	<b>76.36%</b>	<b>28.00</b>	0.992	0.991

**Table 1. Quantitative comparison on VideoCoF-Bench.** We compare VideoCoF with SOTA baselines: InsV2V [4]; Señorita [44] (an I2V model guided by an InstructPix2Pixel [2] first frame); VACE-14B [8] (using GPT-4o generated captions); the concurrent work ICVE [16] (pre-trained 1M, fine-tuned 150k); and Lucy Edit Dev [25]. Despite extensive baseline training data, our VideoCoF is fine-tuned on only 50k source-reasoning-editing triplets and shows superior instruction following and success ratio.



**Figure 6.** Visual comparison between our VideoCoF and other methods on video editing tasks.

use 33 frames source video, 33 frames edited video, and 4 frames reasoning clip. We train with a global batch size of 16 for approximately 8k iterations, optimizing with AdamW [17] and a base learning rate of  $1 \times 10^{-4}$ .

## 4.2. VideoCoF-Bench and Experimental Setting

**VideoCoF-Bench** Previous video-editing benchmarks such as V2VBench [24], TGVE [34], and FIVE-Bench [14] focus on target-prompt edits and mostly are focused on class-level object swap. They were mainly designed for training-free methods and are not suitable for instruction-guided or instance-level video editing. Real-world editing requires precise instruction understanding, including instance- and part-level control (e.g., distinguishing multiple people or

left vs. right), and complex reasoning. To address these gaps, we introduce VideoCoF-Bench. It contains 200 high-quality videos collected from Pexels [20], covering diverse scenes and both landscape and portrait aspect ratios. VideoCoF-Bench includes four tasks: Object Removal, Object Addition, Object Swap, and Local Style Transfer, each with 50 samples. Half of these samples per task are instance-level cases with instance-focused editing prompts.

**Metrics** To evaluate editing performance on VideoCoF-Bench, we employ MLLM-as-a-Judge to provide a holistic evaluation score. This is achieved by prompting **GPT-4o** [19] to assess multiple criteria given the original video, edited video, and user instruction: (1) Instruction Following (editing accuracy), (2) Preservation (unedited regions), (3)

301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317

Ablation on Chain of frames and RoPE design			
	Naive Temporal in Context	VideoCoF	
CoF	$\times$	$\times$	✓
RoPE Design	0–2F–1	0–F–1, 0–F–1	1–F, 0, 1–F
<i>GPT-4o Score</i>			
Instruct Follow↑	8.109	8.064	<b>8.973</b>
Preservation↑	7.930	7.793	<b>8.203</b>
Quality↑	7.394	7.217	<b>7.765</b>
Success Ratio↑*	72.41%	65.52%	<b>76.36%</b>
<i>Perceptual Quality</i>			
CLIP-T↑	26.880	27.088	<b>28.000</b>
CLIP-F↑	0.9907	0.9905	<b>0.9915</b>
DINO↑	0.9857	0.9826	<b>0.9913</b>

Table 2. Ablation on Chain of frames and RoPE design.

318 Video Quality. (4) Success ratio: we prompt the GPT-4o to  
 319 provide a binary Success Ratio (Yes/No) to judge the overall  
 320 success of the edit. We report three perceptual quality met-  
 321 rics quantify low- and high-level visual similarity between  
 322 source and target frames: CLIP-T for image–text alignment,  
 323 CLIP-F for temporal consistency, and DINO for structural  
 324 consistency.

### 325 4.3. Comparison on VideoCoF-Bench

326 We show qualitative and quantitative comparisons of  
 327 VideoCoF-Bench in this section. As shown in Table 1,  
 328 we evaluate VideoCoF against five baseline methods on  
 329 the VideoCoF-Bench benchmark, which spans four distinct  
 330 video editing tasks: multi-instance removal, object addition,  
 331 multi-instance swap, and multi-instance local style transfer.

332 Overall, VideoCoFdemonstrates the best performance in  
**333 Instruct Follow** and **Success Ratio** across all categories.  
 334 Compared to naive temporal in-context editing approaches  
 335 like ICVE [16], our method achieves significantly higher  
 336 success rates and better instruction adherence using only  
 337 **50k** reasoning pairs, whereas ICVE is pre-trained on 1M  
 338 samples and fine-tuned on 150k data.

339 Qualitatively (see Figure 6), our method also shows  
 340 clearer, more faithful edits at the instance level:(a) Multi-  
 341 instance removal: we precisely remove the right instance  
 342 while ICVE[16] incorrectly removes the left instance. (b)  
 343 Object addition: the added girl is correctly placed inside  
 344 the washing machine, matching the instruction. (c) Object  
 345 swap: we replace the elderly person’s face and update cloth-  
 346 ing; Lucy Edit [25] changes only clothing, ICVE fails to dis-  
 347 ambiguate instances, and VACE often alters non-target peo-  
 348 ple. (d) Local style (multi-instance): our model correctly  
 349 identifies and edits the largest cup among several similar ob-  
 350 jects; other methods either fail to edit or mistakenly edit a  
 351 bowl. These qualitative examples demonstrate VideoCoF’s  
 352 stronger instance-level reasoning and higher editing fidelity.

## 353 4.4. Ablation Study

To verify our novel Chain of Frames (CoF) design, partic-  
 354 ularly its “reasoning frames” and the RoPE design for length  
 355 exploration, we conduct an ablation study on the reasoning  
 356 frames, RoPE alignment strategy and reasoning format.

**357 Naive Temporal Incontext VS. CoF** As shown in Ta-  
 358 ble 2, we compare VideoCoF against a “Naive Temporal in-  
 359 context” baseline. This applies temporal in-context learning  
 360 by using the source video as a condition through temporal  
 361 concatenation, an approach similar to ICVE [16].

In contrast, our approach introduces **reasoning frames**  
 362 as a core component of the (CoF) design. This ensures the  
 363 video editing follows a reasoning process, i.e., forcing the  
 364 model to predict the editing region first and then execute the  
 365 versatile edit within that specific area.

The efficacy of this design is evident when comparing  
 366 the first ( $[0, 2F - 1]$ ) and third (VideoCoF) columns in Ta-  
 367 ble 2. The inclusion of CoF brings substantial gains: the  
 368 instruct follow score increases by 10.65% and the success  
 369 ratio improves by 5.46%. Furthermore, the 4.16% increase  
 370 in CLIP-T confirms that our reasoning frames effectively  
 371 enhance the model’s editing accuracy and precision.



Transform the person’s hair into realistic flames.

Figure 7. Length exploration on frames more than training

**Rope Design for length Extrapolation.** As illustrated in Fig 7, the naive approach ( $[0, 2F - 1]$ ) only learns a fixed temporal mapping (e.g., mapping frame  $0_{th}$  to frame  $33_{th}$ ). This prevents length extrapolation, causing severe degradation (blurriness, motion misalignment, and artifacts) when a 33-frame trained model is tested on 81 frames (second row).

In contrast, our RoPE alignment design ( $[1 - F, 0, 1 - F]$ ) generalizes to unseen lengths without quality degradation (third row). As demonstrated in Fig 1, our model extrapolates to 141 frames (4x training length) and beyond, supporting theoretically infinite extrapolation.

This effectiveness is also quantified in Table 2 (third vs. first column). We observe a 3.4% relative increase in the preservation score. Furthermore, the improved DINO score confirms that our RoPE design better preserves the original video’s spatio-temporal structure during editing.

**RoPE Design for Motion Alignment.** Setting the tempo-  
 391 ral index for the reasoning frame latent is a critical design  
 392 choice. A naive approach is to set its index to 0, aligning it  
 393 with the first video frame. This causes two severe issues.

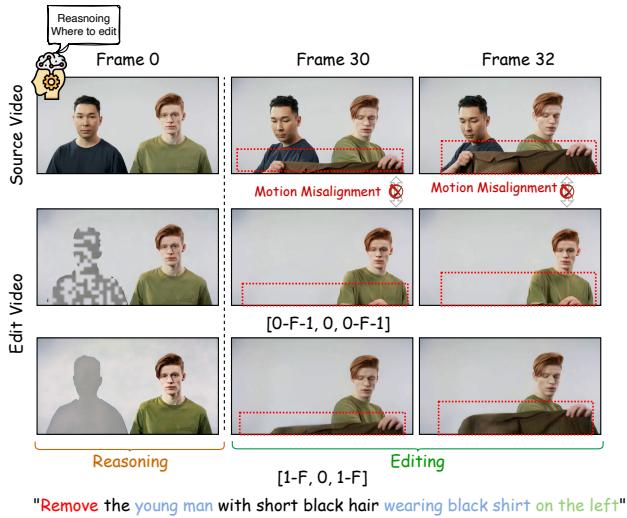


Figure 8. Motion alignment benefit by our rope design

First, it leads to significant motion misalignment (e.g., the subject fails to perform the "lifting clothes" motion in Fig 8, second row). Second, this "0-index" design causes interference with the first editing video frame (also index 0), leading to artifacts where the model incorrectly predicts the first frame as the reasoning frame (Fig 4).

Therefore, we fix the reasoning latent's index to 0, while the source and edited video indices range from 1 to  $F$  (denoted as  $[1 - F, 0, 1 - F]$ ). This strategy allows the reasoning frame to provide clear spatial guidance on **where** to edit, without disrupting the video's temporal structure and motion alignment. The improvements across all metrics in Tab 2 (column 3 vs. column 2) validate this design.

**Reasoning Frame Format** First, we explore the most suitable color for the reasoning frame mask. As shown in Table 3, we compare three formats: (1) A black mask over the unedit region; (2) A red, 50% transparent highlight, same as veggie [40]; and (3) A pure gray mask (value 127, 0% transparency). The quantitative results show that using a gray mask (column 3) for the edit region yields the best performance.

Furthermore, we argue that the reasoning frame should act as a gradual transition from the source video to the edited video. Therefore, we test progressive gray mask. Instead of a single static mask, we interpolate gray mask reasoning frame and editing frame, with transparency is progressively increased (e.g., 0%, 25%, 50%, 75%). As shown by comparing column 4 and column 3 in Table 3, this progressive gray reasoning frame approach works best.

Qualitatively, as shown in Figure 9, the mask format is critical. The black mask fails the deletion task, while the red mask incorrectly deletes content on the right side. In contrast, our progressive gray mask accurately performs the intended deletion on the left. We conclude from these experiments that the optimal reasoning format is a gray mask

Color Transparency	Ablation on Reasoning Frame Format			
	Black (bg) (0%)	Red (50%)	Gray (0%)	Gray (0-75%)
<i>GPT-4o Score</i>				
Instruct Follow↑	7.512	7.805	8.069	<b>8.973</b>
Preservation↑	7.034	7.350	7.709	<b>8.203</b>
Quality↑	6.155	6.501	6.926	<b>7.765</b>
Success Ratio↑*	52.170%	60.330%	67.980%	<b>76.36%</b>
<i>Perceptual Quality</i>				
CLIP-T↑	26.550	26.810	27.143	<b>28.000</b>
CLIP-F↑	0.9810	0.9855	0.9890	<b>0.9915</b>
DINO↑	0.9750	0.9790	0.9826	<b>0.9913</b>

Table 3. Ablation on transparency mask settings.

with progressive transparency.

430

Predicted Reasoning Format



"Remove the woman with tattoos wearing beige bra on the left"

Figure 9. Ablation on reasoning frame format

## 5. Conclusion

In this paper, we introduced VideoCoF, a unified model for universal video editing via temporal reasoning. We identified that existing temporal in-context learning approaches often fail due to a lack of explicit spatial cues, leading to weak instruction-to-region mapping and imprecise localization. To address these issues, we proposed the innovative Chain of Frames. CoF compels the video diffusion model to follow a "see, reason, then edit" process by first predicting the editing region before executing the versatile edit. Furthermore, to solve the length generalization challenge, we developed a novel RoPE alignment paradigm that accounts for the reasoning latent. This design enables 4 times exploration in the inference. Experimental results show that VideoCoF achieves SOTA performance using a mere 50k video pairs, validating the efficiency and effectiveness of our temporal reasoning design.

## References

- [1] 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.

431

432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447

- [2] 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] Lan Chen, Yuchao Gu, and Qi Mao. Univid: Unifying vision tasks with pre-trained video generation models. *arXiv preprint arXiv:2509.21760*, 2025.
- [4] Jiaxin Cheng, Tianjun Xiao, and Tong He. Consistent video-to-video transfer using synthetic dataset. *arXiv preprint arXiv:2311.00213*, 2023.
- [5] Yuren Cong, Mengmeng Xu, Christian Simon, Shoufa Chen, Jiawei Ren, Yanping Xie, Juan-Manuel Perez-Rua, Bodo Rosenhahn, Tao Xiang, and Sen He. Flatten: Optical flow-guided attention for consistent text-to-video editing. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024.
- [6] Yuchao Gu, Yipin Zhou, Bichen Wu, Licheng Yu, Jia-Wei Liu, Rui Zhao, Jay Zhangjie Wu, David Junhao Zhang, Mike Zheng Shou, and Kevin Tang. Videoswap: Customized video subject swapping with interactive semantic point correspondence. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7621–7630, 2024.
- [7] Nicholas Guttenberg. Diffusion with offset noise. <https://www.crosslabs.org/blog/diffusion-with-offset-noise>, 2023.
- [8] 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.
- [9] Xuan Ju, Tianyu Wang, Yuqian Zhou, He Zhang, Qing Liu, Nanxuan Zhao, Zhifei Zhang, Yijun Li, Yuanhao Cai, Shaoteng Liu, et al. Editverse: Unifying image and video editing and generation with in-context learning. *arXiv preprint arXiv:2509.20360*, 2025.
- [10] Xuan Ju, Weicai Ye, Quande Liu, Qiulin Wang, Xintao Wang, Pengfei Wan, Di Zhang, Kun Gai, and Qiang Xu. Fulldit: Multi-task video generative foundation model with full attention. *arXiv preprint arXiv:2503.19907*, 2025.
- [11] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022.
- [12] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. Hunyuanyvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024.
- [13] Max Ku, Dongfu Jiang, Cong Wei, Xiang Yue, and Wenhui Chen. Viescore: Towards explainable metrics for conditional image synthesis evaluation, 2023.
- [14] Minghan Li, Chenxi Xie, Yichen Wu, Lei Zhang, and Mengyu Wang. Five: A fine-grained video editing benchmark for evaluating emerging diffusion and rectified flow models. *arXiv preprint arXiv:2503.13684*, 2025.
- [15] Xiaowen Li, Haolan Xue, Peiran Ren, and Liefeng Bo. Dif-fueraser: A diffusion model for video inpainting. *arXiv preprint arXiv:2501.10018*, 2025.
- [16] Xinyao Liao, Xianfang Zeng, Ziye Song, Zhoujie Fu, Gang Yu, and Guosheng Lin. In-context learning with unpaired clips for instruction-based video editing. *arXiv preprint arXiv:2510.14648*, 2025.
- [17] Ilya Loshchilov, Frank Hutter, et al. Fixing weight decay regularization in adam. *arXiv preprint arXiv:1711.05101*, 5 (5):5, 2017.
- [18] Chong Mou, Qichao Sun, Yanze Wu, Pengze Zhang, Xinghui Li, Fulong Ye, Songtao Zhao, and Qian He. Instructx: Towards unified visual editing with mllm guidance. *arXiv preprint arXiv:2510.08485*, 2025.
- [19] OpenAI. Hello gpt-4o. Blog post, 2024.
- [20] Pexels. Pexels: Free stock photos, royalty free stock images & videos. <https://www.pexels.com/>, 2025. Accessed: 2025-11-06.
- [21] Chenyang Qi, Xiaodong Cun, Yong Zhang, Chenyang Lei, Xintao Wang, Ying Shan, and Qifeng Chen. Fatezero: Fusing attentions for zero-shot text-based video editing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15932–15942, 2023.
- [22] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kun-chang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, Zhaoyang Zeng, Hao Zhang, Feng Li, Jie Yang, Hongyang Li, Qing Jiang, and Lei Zhang. Grounded sam: Assembling open-world models for diverse visual tasks, 2024.
- [23] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- [24] Wenhao Sun, Rong-Cheng Tu, Jingyi Liao, and Dacheng Tao. Diffusion model-based video editing: A survey. *arXiv preprint arXiv:2407.07111*, 2024.
- [25] DecartAI Team. Lucy edit: Open-weight text-guided video editing, 2025.
- [26] Yuanpeng Tu, Hao Luo, Xi Chen, Sihui Ji, Xiang Bai, and Hengshuang Zhao. Videoanydoor: High-fidelity video object insertion with precise motion control. In *Proceedings of the Special Interest Group on Computer Graphics and Interactive Techniques Conference Conference Papers*, pages 1–11, 2025.
- [27] 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.
- [28] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.
- [29] Cong Wei, Quande Liu, Zixuan Ye, Qiulin Wang, Xintao Wang, Pengfei Wan, Kun Gai, and Wenhui Chen. Univideo: Unified understanding, generation, and editing for videos. *arXiv preprint arXiv:2510.08377*, 2025.
- [30] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al.

- 567 Chain-of-thought prompting elicits reasoning in large lan- 625  
568 guage models. *Advances in neural information processing* 626  
569 *systems*, 35:24824–24837, 2022. 627
- 570 [31] Thaddäus Wiedemer, Yuxuan Li, Paul Vicol, Shixiang Shane 628  
571 Gu, Nick Matarese, Kevin Swersky, Been Kim, Priyank 629  
572 Jaini, and Robert Geirhos. Video models are zero-shot learn- 630  
573 ers and reasoners. *arXiv preprint arXiv:2509.20328*, 2025. 631
- 574 [32] Haoning Wu, Erli Zhang, Liang Liao, Chaofeng Chen, Jing- 632  
575 wen Hou Hou, Annan Wang, Wenxiu Sun Sun, Qiong Yan, 633  
576 and Weisi Lin. Exploring video quality assessment on user 634  
577 generated contents from aesthetic and technical perspectives. 635  
578 In *International Conference on Computer Vision (ICCV)*,  
579 2023.
- 580 [33] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian 631  
581 Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying Shan, Xiaohu 632  
582 Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning 633  
583 of image diffusion models for text-to-video generation. In 634  
584 *Proceedings of the IEEE/CVF international conference on* 635  
585 *computer vision*, pages 7623–7633, 2023.
- 586 [34] Jay Zhangjie Wu, Xiuyu Li, Difei Gao, Zhen Dong, Jin- 631  
587 bin Bai, Aishani Singh, Xiaoyu Xiang, Youzeng Li, Zuwei 632  
588 Huang, Yuanxi Sun, et al. Cvpr 2023 text guided video edit- 633  
589 ing competition. *arXiv preprint arXiv:2310.16003*, 2023.
- 590 [35] Shitao Xiao, Yueze Wang, Junjie Zhou, Huaying Yuan, Xin- 631  
591 grun Xing, Ruiran Yan, Chaofan Li, Shuteng Wang, Tiejun 632  
592 Huang, and Zheng Liu. Omnigen: Unified image genera- 633  
593 tion. In *Proceedings of the Computer Vision and Pattern* 634  
594 *Recognition Conference*, pages 13294–13304, 2025.
- 595 [36] Xiangpeng Yang, Linchao Zhu, Hehe Fan, and Yi Yang. 631  
596 Videograin: Modulating space-time attention for multi- 632  
597 grained video editing. In *The Thirteenth International Con- 633  
598 ference on Learning Representations*, 2025.
- 599 [37] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu 631  
600 Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiao- 632  
601 han Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video 633  
602 diffusion models with an expert transformer. *arXiv preprint* 634  
603 *arXiv:2408.06072*, 2024.
- 604 [38] Xizuan Ye, Xuanhua He, Quande Liu, Qiulin Wang, Xintao 631  
605 Wang, Pengfei Wan, Di Zhang, Kun Gai, Qifeng Chen, and 632  
606 Wenhan Luo. Unic: Unified in-context video editing. *arXiv* 633  
607 *preprint arXiv:2506.04216*, 2025.
- 608 [39] Xizuan Ye, Huijuan Huang, Xintao Wang, Pengfei Wan, Di 631  
609 Zhang, and Wenhan Luo. Stylemaster: Stylize your video 632  
610 with artistic generation and translation. In *Proceedings of* 633  
611 *the Computer Vision and Pattern Recognition Conference*, 634  
612 pages 2630–2640, 2025.
- 613 [40] Shoubin Yu, Difan Liu, Ziqiao Ma, Yicong Hong, Yang 631  
614 Zhou, Hao Tan, Joyce Chai, and Mohit Bansal. Veg- 632  
615 gie: Instructional editing and reasoning video concepts with 633  
616 grounded generation. In *Proceedings of the IEEE/CVF In- 634  
617 ternational Conference on Computer Vision*, pages 15147– 635  
618 15158, 2025.
- 619 [41] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding 631  
620 conditional control to text-to-image diffusion models. In 632  
621 *Proceedings of the IEEE/CVF international conference on* 633  
622 *computer vision*, pages 3836–3847, 2023.
- 623 [42] Zechuan Zhang, Ji Xie, Yu Lu, Zongxin Yang, and Yi Yang. 631  
624 In-context edit: Enabling instructional image editing with in- 632