

# ENHANCE-A-VIDEO: BETTER GENERATED VIDEO FOR FREE

Anonymous authors

Paper under double-blind review

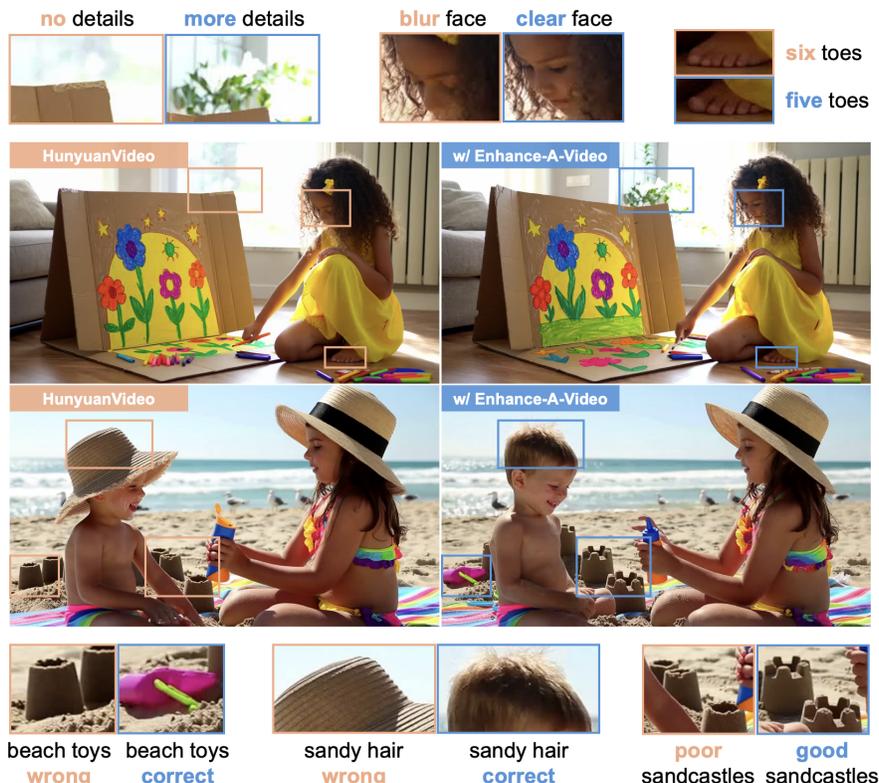


Figure 1: Enhance-A-Video boosts diffusion transformers-based video generation quality at minimal cost - no training needed, no extra learnable parameters, no memory overhead. Detailed captions are available in Appendix G.

## ABSTRACT

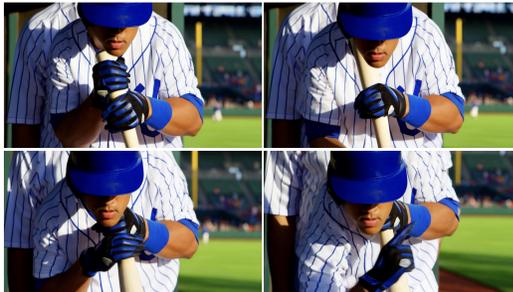
DiT-based video generation has achieved remarkable results, but research into enhancing existing models remains relatively unexplored. In this work, we introduce a training-free approach to enhance the coherence and quality of DiT-based generated videos, named **Enhance-A-Video**. The core idea is to enhance the cross-frame correlations based on non-diagonal temporal attention distributions. Thanks to its simple design, our approach can be easily applied to most DiT-based video generation frameworks without any retraining or fine-tuning. Across various DiT-based video generation models, our approach demonstrates promising improvements in both temporal consistency and visual quality. We hope this research can inspire future explorations in video generation enhancement.

## 1 INTRODUCTION

Diffusion transformer (DiT) models (Peebles & Xie, 2022) have revolutionized video generation, enabling the creation of realistic and compelling videos (Yang et al., 2024; Brooks et al., 2024; Lin

et al., 2024; Xu et al., 2024a; Kong et al., 2025; Wan et al., 2025). However, achieving temporal consistency across frames while maintaining fine-grained details remains a significant challenge. Many existing methods generate videos that suffer from unnatural transitions and degraded quality, as illustrated in Fig. 2, which limits their practical applicability in real-world scenarios and professional applications (Yan et al., 2023; Henschel et al., 2024).

Video generation enhancement (He et al., 2024; Ma et al., 2025) is designed for addressing the above limitations, where two objectives are primarily considered: (i) maintaining temporal consistency across frames, which ensures smooth and coherent transitions, and (ii) improving spatial details, which enhances the visual quality of each frame for more realistic video outputs. In UNet-based video generation (Zhang et al., 2023a; Guo et al., 2024; Xu et al., 2024b; Zhang et al., 2024; Xia et al., 2024; Bu et al., 2025; Yuan et al., 2025; Li et al., 2025b), Upscale-A-Video (Zhou et al., 2024) integrated a local-global temporal strategy for better temporal coherence, and VEnhancer (He et al., 2024) designed a video ControlNet (Zhang et al., 2023b) to enhance spatial and temporal resolution simultaneously. Nevertheless, the exploration of enhancing DiT-based video generation remains limited, particularly in addressing challenges of temporal consistency and spatial detail preservation.



A baseball player grips a bat in black gloves, wearing a blue-and-white uniform and cap, with a blurred crowd and green field highlighting his focused stance.

Figure 2: Video sample of HunyuanVideo model with unnatural head movements, repeated right hands and conflicting glove color.

In DiT-based video generation, temporal attention (Singer et al., 2022; Villegas et al., 2023; Tan et al., 2023) plays a crucial role in ensuring coherence among frames. Through careful analysis of temporal attention in DiT blocks, we made an important observation as shown in Fig. 3: In DiT-based models, temporal attention is typically concentrated along the diagonal, enabling more efficient utilization of model capacity to capture the most relevant temporal dependencies (Lu et al., 2024; Cai et al., 2024; Li et al., 2025c). However, the imbalance between cross-frame (non-diagonal elements) and intra-frame attention (diagonal elements) may cause foundation models to underutilize global cross-frame information, resulting in inconsistencies across frames – such as abrupt transitions or blurred details in the generated videos.

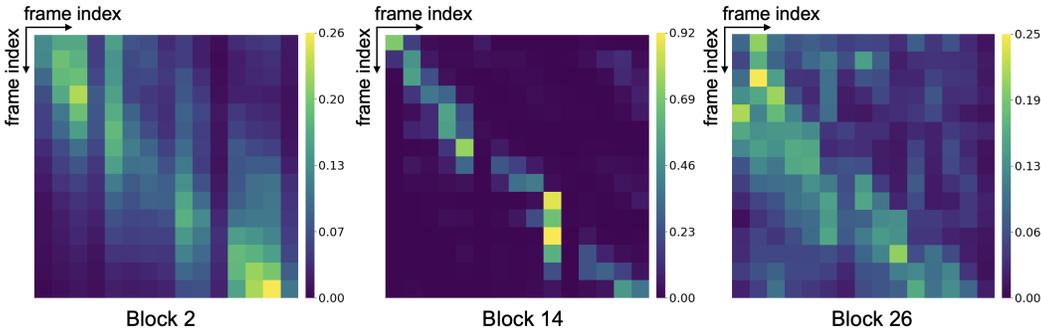


Figure 3: Visualization of temporal attention distributions in Open-Sora for blocks 2, 14, and 26 at denoising step 30, where non-diagonal elements are considerably weaker than diagonal elements.

A straightforward approach to address the highly uneven distribution between intra-frame and cross-frame attention is to incorporate the temperature mechanism (Peeperkorn et al., 2024; Renze & Guven, 2024) into the computation of temporal attention. However, this naive strategy significantly alters the original temporal attention patterns as the depth of DiT blocks and the number of denoising steps increase. Preserving the original temporal attention patterns is important, as they encode prior knowledge acquired during training for video generation. Disrupting these patterns compromises the learned knowledge and leads to reduced generation quality, as demonstrated in Appendix C.

To more effectively integrate temperature-based adjustments into temporal attention, we propose applying a cross-frame intensity-driven temperature in the residual connection (He et al., 2015).

Specifically, the cross-frame intensity is calculated as the average of the non-diagonal temporal attention weights. This value is then used to adaptively balance dependencies across frames, leading to improved video quality. Furthermore, because attention outputs contribute relatively low compared to hidden states within the residual connection (Si et al., 2024; Ma et al., 2024a), a moderate increase of the original attention output using cross-frame intensity has a limited impact on the overall attention pattern in the final representation of each DiT block.

Building on these insights, we develop a novel, training-free, and plug-and-play approach, Enhance-A-Video, to improve the temporal and spatial quality of DiT-based generated videos. The method introduces two key innovations: a cross-frame intensity to capture cross-frame information within the residual connection and an enhance temperature parameter to scale the calculated cross-frame intensity. By strengthening cross-frame correlations from the temperature perspective, our approach enhances temporal consistency and preserves fine visual details effectively. **A notable advantage is that this method can be readily integrated into prevalent DiT-based video generation frameworks with negligible computational overhead.**

We conduct a comprehensive experimental evaluation of our approach across several benchmark DiT-based video generation models, including Wan (Wan et al., 2025), HunyuanVideo (Kong et al., 2025), CogVideoX (Yang et al., 2024), LTX-Video (HaCohen et al., 2024), Open-Sora (Zheng et al., 2024) and Open-Sora-Plan (Lin et al., 2024). By incorporating Enhance-A-Video during the inference phase, these models demonstrate a significant improvement in generated video quality by reducing temporal inconsistencies and refining visual fidelity. In particular, the enhanced cross-frame attention not only mitigates temporal inconsistencies by encouraging the model to exploit richer contextual information, but also facilitates smoother transmission of prompt and spatial information across frames, thereby enabling the generation of finer visual details and more prompt-consistent videos, as illustrated in Fig. 1.

## 2 RELATED WORK

**Video Generation.** Recent advancements in video generation have been driven by powerful diffusion transformer-based models (Chen et al., 2024; Ma et al., 2024b; Gao et al., 2024; Lu et al., 2024). Sora (Brooks et al., 2024) has demonstrated exceptional capabilities in generating realistic and long-duration videos, establishing itself as a significant milestone in text-to-video generation. CogVideoX (Yang et al., 2024) introduced a 3D full attention mechanism and expert transformers to improve motion consistency and semantic alignment. HunyuanVideo (Kong et al., 2025) introduces a hybrid stream block with enhanced semantic understanding. Nevertheless, key challenges remain in video generation, including temporal inconsistency and the degradation of fine-grained spatial details.

**Temperature Parameter.** The temperature parameter is a well-known concept in deep learning, primarily used to control the distribution of attention or output probabilities in generative models (Peeperkorn et al., 2024; Renze & Guven, 2024). In natural language generation tasks, the temperature is often adjusted during inference to modulate the diversity of the generated text (Holtzman et al., 2020). A higher temperature increases randomness, promoting creativity, while a lower temperature encourages deterministic and coherent outputs. Recently, the concept has been explored in vision-related tasks, such as visual question answering and multimodal learning (Chen et al., 2021), where temperature adjustments are applied to balance multimodal attention distributions. Yet, its application within DiT-based video generation, and in particular its impact on temporal attention, has not been thoroughly investigated.

## 3 METHODOLOGY

### 3.1 DIFFUSION TRANSFORMER MODELS

Diffusion Transformer models are inspired by the success of diffusion models in generating high-quality images and videos by iteratively refining noisy data (Ho et al., 2022b;a; Xu et al., 2023; Blattmann et al., 2023; Wang et al., 2023; Esser et al., 2024). These models combine the strengths of diffusion processes and transformer architectures to model temporal and spatial dependencies in video generation. The forward diffusion process adds noise to the data over  $T$  timesteps, gradually converting it into a noise distribution. Starting from clean data  $x_0$ , the noisy data at timestep  $t$  is

obtained as:

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \mathbf{z}_t, \quad \text{for } t = 1, \dots, T, \quad (1)$$

where  $\alpha_t$  controls the noise schedule and  $\mathbf{z}_t \sim \mathcal{N}(0, \mathbf{I})$  is Gaussian noise. As  $t$  increases,  $\mathbf{x}_t$  approaches a standard normal distribution  $\mathcal{N}(0, \mathbf{I})$ . To recover the original data distribution, the reverse diffusion process progressively removes noise from  $\mathbf{x}_t$  until reaching  $\mathbf{x}_0$ :

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t)), \quad (2)$$

where  $\mu_\theta$  and  $\Sigma_\theta$  are learned parameters representing the mean and covariance of the denoised distribution.

### 3.2 TEMPERATURE IN DiT-BASED VIDEO GENERATION

The temperature is a critical concept in large language model inference, controlling the randomness and coherence of the generated tokens. The probability  $P(x)$  of generating a token  $x$  is adjusted using the temperature  $\tau$  as:

$$P(x) = \frac{\exp\left(\frac{z(x)}{\tau}\right)}{\sum_{x'} \exp\left(\frac{z(x')}{\tau}\right)} \quad (3)$$

where  $z(x)$  represents the unnormalized logit for token  $x$ , and  $\tau > 0$  controls the degree of randomness: a lower  $\tau$  makes the output more deterministic, while a higher  $\tau$  increases diversity by flattening the probability distribution.

In video generation, a similar temperature principle can be considered when using DiT models, where the temporal attention mechanism controls the relationship between generated frames. Eq. 4 presents a direct usage of temperature in temporal attention of DiT models. A higher  $\tau$  yields a more uniform temporal attention, enabling broader context integration.

$$\text{Attn}(Q, K) = \text{softmax}\left(\frac{QK^\top}{\tau \cdot \sqrt{d_k}}\right) \quad (4)$$

Nevertheless, as illustrated in Appendix C, directly applying  $\tau$  to temporal attention causes increasing changes to the original attention weights as the model deepens and denoising steps accumulate, which can lead to overly smooth motion, loss of visual details, and unstable video generation.

Therefore, incorporating a temperature into the attention output within the residual connection is more practical, given the significant magnitude difference presented in Appendix D between attention outputs and hidden states. Furthermore, non-diagonal temporal attention weights can be used to quantify cross-frame correlations, enabling adaptive temperature adjustment to enhance global information aggregation across frames, thereby improving spatial diversity and temporal consistency.

### 3.3 ENHANCE BLOCK

To better adaptively adjust the temperature in the temporal attention mechanism, we propose a novel method, Enhance-A-Video, to enhance temporal consistency in video generation by utilizing the **non-diagonal temporal attention** with **enhance temperature parameter**. The cross-frame intensity is measured by the non-diagonal temporal attention, where higher values enable the model to focus on a broader temporal context, corresponding to higher temperature. By further introducing the enhance temperature parameter to scale the cross-frame intensity, we appropriately adjust the temporal attention outputs as a training-free enhancement.

As presented in Fig. 4, we design an **Enhance Block** as a parallel branch to the temporal attention mechanism. The Enhance Block operates as follows:

First, the temporal attention map  $A^{temp}$  is computed independently from the input of the 3D attention module, since direct access to the 3D attention is not possible due to the adoption of Flash Attention (Dao et al., 2022). Specifically, For video latent  $\mathbf{z} \in \mathbb{R}^{B \times (F \times H \times W) \times d}$  with batch size  $B$ ,  $F$  frames, spatial dimensions  $H \times W$  and hidden size  $d$ , we reshape features by merging spatial dimensions into the batch size, yielding  $\tilde{\mathbf{z}} \in \mathbb{R}^{(B \times H \times W) \times F \times d}$ . Self-attention (Vaswani et al., 2023) is then applied along the frame axis:

$$A^{temp} = \text{Attn}(Q(\tilde{\mathbf{z}}), K(\tilde{\mathbf{z}})) \in \mathbb{R}^{(B \times H \times W) \times F \times F} \quad (5)$$

where  $Q$  and  $K$  denote the Query and Key heads, and  $A^{temp}$  satisfies  $\sum_{j=1}^F A_{(b,i,j)}^{temp} = 1$ . The diagonal elements  $A_{ii}^{temp}$  correspond to intra-frame attention, and the non-diagonal elements  $A_{ij}^{temp}$  ( $i \neq j$ ) represent cross-frame attention.

Next, the Cross-Frame Intensity ( $CFI$ ) is calculated by averaging the non-diagonal elements of the attention map. The  $CFI$  is then multiplied by the enhance temperature parameter  $\tau$  to enhance cross-frame correlations better:

$$CFI = \frac{1}{F(F-1)} \sum_{i=1}^F \sum_{\substack{j=1 \\ j \neq i}}^F A_{ij}^{temp}. \quad (6)$$

$$CFI_{enhanced} = \text{clip}((\tau + F) \cdot CFI, 1). \quad (7)$$

Noticeably, the enhanced Cross-Frame Intensity ( $CFI_{enhanced}$ ) is clipped at a minimum value of 1, which prevents severe weakening of cross-frame correlations during enhancement.

Finally, the output of the Enhance Block ( $CFI_{enhanced}$ ) is utilized to enhance the original temporal attention block output  $\mathbf{O}_{attn}$  in the residual connection:

$$\mathbf{O}_{final} = CFI_{enhanced} \cdot \mathbf{O}_{attn} + \mathbf{H}. \quad (8)$$

where  $\mathbf{H}$  represents the hidden states that are inputs of the attention block.

When  $CFI_{enhanced}$  exceeds 1, indicating significant cross-frame information, the ratio of temporal attention block outputs is correspondingly amplified in  $\mathbf{O}_{final}$ . Otherwise, the connection defaults to a standard residual connection. Since  $\mathbf{O}_{attn}$  is relatively small compared to  $\mathbf{H}$ , modest enhancements (small  $CFI_{enhanced}$ ) to  $\mathbf{O}_{attn}$  slightly affect the  $\mathbf{O}_{final}$  distribution, enabling Enhance-A-Video to enhance cross-frame attention without substantially altering original temporal attention patterns. The complete analytical details are available in Appendix D.

The temporal attention difference map in Fig. 5 shows the difference between the temporal attention of the original CogVideoX model and w/ Enhance-A-Video, illustrating how Enhance-A-Video properly strengthens cross-frame attention without disrupting the initial temporal attention pattern. Specifically, certain non-diagonal elements (blue areas) are moderately increased (e.g.,  $0.9 \times 10^{-2}$ ), indicating enhanced cross-frame correlations. Meanwhile, the diagonal elements experience a minimal reduction ( $3.3 \times 10^{-2}$  at most), which ensures stable intra-frame attention and preserves existing fine-grained visual details. More analysis can be found in Appendix C.

## 4 EXPERIMENTS

### 4.1 SETUP

To evaluate the effectiveness of our proposed Enhance-A-Video method, we conduct experiments on video generation models incorporating two types of attention mechanisms: 3D full attention and spatial-temporal attention. Specifically, we choose several representative models for each category:

**3D Full Attention Model:** Wan (Wan et al., 2025), HunyuanVideo (Kong et al., 2025), CogVideoX (Yang et al., 2024) and LTX-Video (HaCohen et al., 2024), which employ 3D full attention to model spatial and temporal dependencies simultaneously.

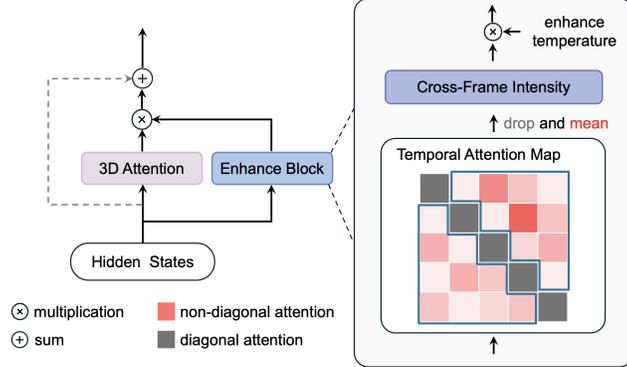


Figure 4: Overview of the Enhance Block. The block computes the average of non-diagonal elements from the temporal attention map as Cross-Frame Intensity ( $CFI$ ). The  $CFI$  is scaled by the temperature parameter and fused back to enhance the temporal attention output.

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

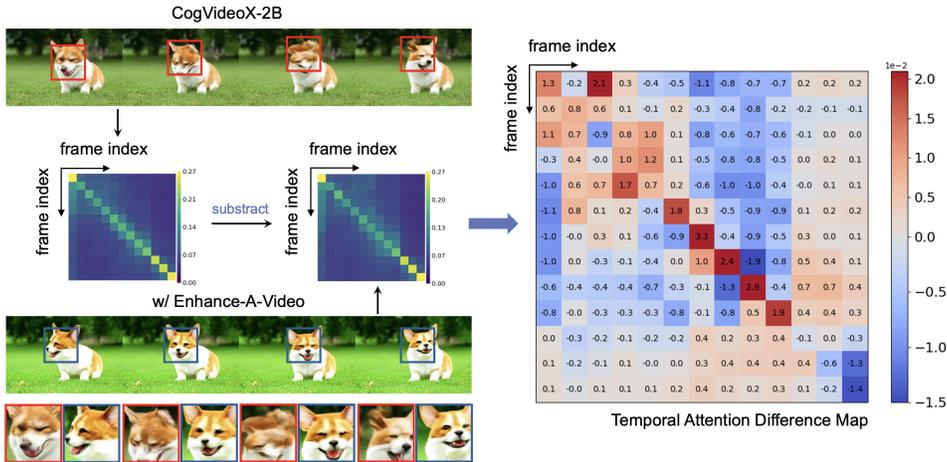


Figure 5: Temporal attention difference map between original CogVideoX model and w/ Enhance-A-Video of layer 29 at denoising step 50. Non-diagonal elements in the attention matrix of w/ Enhance-A-Video show higher values (shown in blue), while diagonal elements have reduced values (shown in red).

**Spatial-Temporal Attention Model:** Open-Sora (Zheng et al., 2024) and Open-Sora-Plan v1.0.0 (Lin et al., 2024), which decompose the attention mechanism into separate spatial and temporal components for computational efficiency and scalability.

Detailed experimental configurations are presented in Appendix E.1. We follow the original setup of these methods exactly and incorporate the Enhance Block exclusively into the temporal attention calculation of these models during the inference phase without additional retraining or fine-tuning. To facilitate the practical use of Enhance-A-Video, we provide the recommended range for the enhance temperature parameter in Appendix E.2.

#### 4.2 3D FULL ATTENTION MODEL

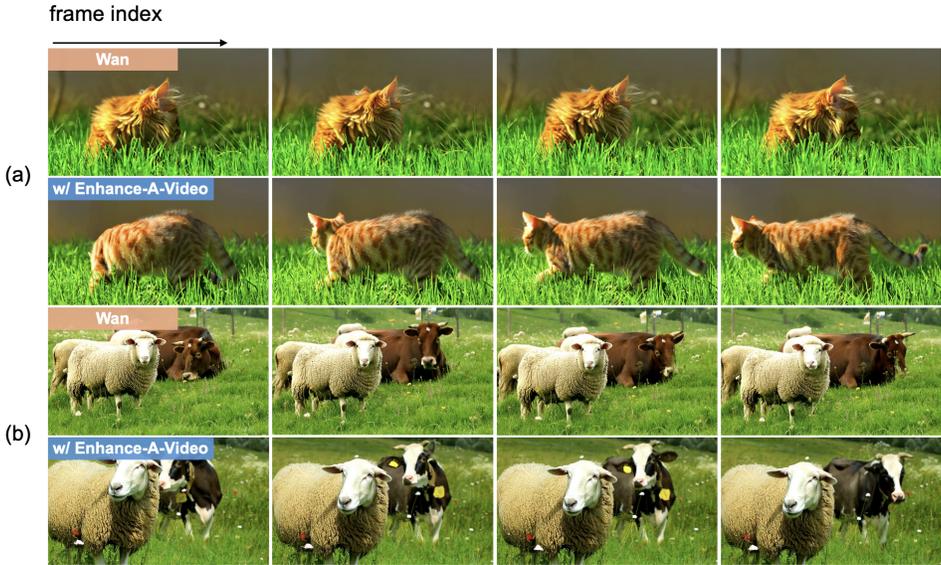


Figure 6: Qualitative results of Enhance-A-Video on Wan. Captions: (a) A cat walks on the grass, realistic. (b) A sheep and a cow.

Wan (Wan et al., 2025) is a state-of-the-art text-to-video diffusion model recognized for producing high-resolution and temporally coherent videos from textual prompts. Fig. 6 demonstrates the



Figure 7: Qualitative results of Enhance-A-Video on more 3D full attention models. Videos generated by baselines (HunyuanVideo and CogVideoX) w/ Enhance-A-Video show significantly improved detail and quality compared to the original foundation models.

effectiveness of Enhance-A-Video on improving the foundation model with more dynamic and prompt-consistent video generation.

In the first case, Wan’s output depicts a cat with unnatural fur lying motionless on the grass, whereas Enhance-A-Video generates a more realistically styled cat walking to the left, offering significantly better alignment with the prompt through improved motion dynamics and spatial content. In the second case, Enhance-A-Video generates the sheep and cow with enhanced texture sharpness, improved color fidelity, and clearer background details, resulting in frames that are more realistic and temporally consistent compared to those produced by the original model.

HunyuanVideo (Kong et al., 2025) is a large-scale text-to-video diffusion model that combines a causal 3D VAE, dual-to-single-stream Transformers, and a multimodal LLM encoder to generate high-quality videos. Our implementation of Enhance-A-Video augmentation in HunyuanVideo improves the model’s video generation capabilities effectively. The results shown in Fig. 7 demonstrate that Enhance-A-Video consistently produces more realistic images with better details.

In the first row for HunyuanVideo, HunyuanVideo introduces conflicting artifacts, such as *duplicated right hands* and *unnatural head movements*. In contrast, Enhance-A-Video captures the baseball player’s motion with greater fluidity and richer detail. In the second row, Enhance-A-Video enhances the appearance of the silver plane, making it more realistic.

By applying Enhance-A-Video to CogVideoX (Yang et al., 2024), we observe significant improvements in prompt-video consistency and visual detail. In the second row of Fig. 7 for CogVideoX, CogVideoX fails to accurately capture the prompt describing a “balloon full of water”, generating only vague water splashes without the balloon. In contrast, the enhanced model produces videos that better align with the given prompts while delivering smoother transitions and clearer visuals. In addition, the results for LTX-Video, along with additional visual examples from the evaluated models, are provided in Appendix H.

#### 4.3 SPATIAL-TEMPORAL ATTENTION MODEL

Open-Sora (Zheng et al., 2024) is an efficient text-to-video generation model that utilizes a decomposed spatial-temporal attention mechanism to balance computational efficiency and video quality. Incorporating the Enhance-A-Video augmentation into Open-Sora significantly improved temporal consistency and spatial detail preservation. As demonstrated in Fig. 8, the enhanced model produces videos with more natural motion transitions and more realistic visual details.

Open-Sora-Plan v1.0.0 (Lin et al., 2024) is a text-to-video generation model leveraging a multi-resolution latent diffusion framework for high-quality, temporally coherent videos. As shown in Fig. 8, Enhance-A-Video creates clearer leaves and sharper flower details, removing the blur seen in the baseline model. These improvements highlight Enhance-A-Video’s ability to enhance cross-frame attention and produce visually high-quality videos.



Figure 8: Qualitative results of Enhance-A-Video on spatial-temporal Attention models.

#### 4.4 QUANTITATIVE ANALYSIS

Table 1: VBench results comparing baseline models and w/ Enhance-A-Video (EAV) across 5 evaluation criteria.

Model	Subject Consistency $\uparrow$	Temporal Flickering $\uparrow$	Imaging Quality $\uparrow$	Multiple Objects $\uparrow$	Spatial Relationship $\uparrow$
Wan 2.1	91.80	98.00	64.94	91.67	98.10
w/ EAV	<b>92.18</b>	<b>98.44</b>	<b>65.18</b>	<b>92.78</b>	<b>99.59</b>
CogVideoX	92.23	90.29	57.15	52.29	52.27
w/ EAV	<b>92.93</b>	<b>91.77</b>	<b>58.56</b>	<b>54.42</b>	<b>56.82</b>
Open-Sora	94.14	98.33	60.43	49.54	60.19
w/ EAV	<b>94.27</b>	<b>98.92</b>	<b>61.04</b>	<b>53.92</b>	<b>64.58</b>

We employ VBench (Huang et al., 2024) for automatic evaluation across multiple video metrics. CogVideoX and Open-Sora use the standard VBench prompt set, while Wan is evaluated on 100 randomly sampled prompts due to computational constraints.

The results in Tab. 1 demonstrate that Enhance-A-Video consistently strengthens video generation across multiple dimensions: it improves subject consistency and temporal stability by amplifying cross-frame attention adaptively, leading to smoother and more coherent motion, while also boosting imaging quality through better preservation of fine details. Moreover, Enhance-A-Video enhances complex scene understanding and generation, as reflected in higher scores for multiple objects and spatial relationships, indicating more accurate object interactions and structural alignment.

Importantly, these benefits appear across different backbone models (Wan 2.1, CogVideoX, Open-Sora), highlighting our method’s robustness and generalizability as a lightweight, training-free plug-in that enhances temporal coherence and perceptual fidelity without retraining costs. Further results are presented in Appendix E.4.

Table 2: User study results comparing baseline models and w/ Enhance-A-Video across evaluation criteria. We compute the proportion of votes received by the baseline and our method respectively.

Model	Overall	Temporal Consistency	Prompt-Video Consistency	Visual Quality
Baseline	20.30	22.73	21.82	16.36
w/ Enhance-A-Video	<b>79.70</b>	<b>77.27</b>	<b>78.18</b>	<b>83.64</b>

Furthermore, we evaluated video quality through a blind user study of 110 participants. Each person compared two videos generated from the same text prompt and random seed - one from baseline models and one from w/ Enhance-A-Video. The videos were shown in random order to prevent bias. Participants chose which video they preferred based on three criteria: temporal consistency, prompt-video consistency, and overall visual quality. The detailed setting is given in Appendix E.3.

Tab. 2 presents the main user study results for chosen models and w/ Enhance-A-Video of each evaluation criterion. The results show that models using Enhance-A-Video received the majority of preference, demonstrating that Enhance-A-Video notably enhances the text-to-video models’ performance in all evaluated aspects:



444 Figure 9: Ablation study on the enhance temperature parameter in the Enhance Block. Moderate  
 445 values balance temporal consistency and visual diversity, while extreme values degrade performance.  
 446

447 **Temporal Consistency.** The usage of Cross-Frame Intensity (*CFI*) and the enhance temperature  
 448 parameter strengthens cross-frame connections. This results in smoother motion transitions and  
 449 improved frame-to-frame alignment, which creates a more stable and coherent visual experience in  
 450 the generated video.

451 **Prompt-Video Consistency.** In diffusion-based video generation, video frames are progressively  
 452 denoised based on the prompt. However, the lack of temporal attention in cross-frame information  
 453 transmission causes the semantic alignment between the video and the prompt to deviate gradually  
 454 during generation. Enhancing cross-frame information by Enhance-A-Video ensures that objects and  
 455 actions in the scene remain consistent with the prompt. This smooth semantic evolution avoids abrupt  
 456 or inconsistent content, improving the alignment between the generated video and the given prompt.  
 457

458 **Visual Quality.** By using *CFI* and the enhanced temperature parameter, the model makes better use  
 459 of information from contextual frames to improve details, especially in object textures and edges.  
 460 The improved cross-frame attention smooths the denoising process and reduces random changes,  
 461 allowing the model to generate more consistent motion and avoid unrealistic movements.  
 462

#### 463 4.5 ABLATION STUDY

464 **Impact of Temperature.** To better understand the impact of the temperature parameter, we conduct  
 465 an ablation study by varying the enhance temperature parameter in the Enhance Block. Results  
 466 in Fig. 9 indicate that moderate temperature values achieve the best balance between temporal  
 467 consistency and diversity, while extreme values (too low or too high) will degrade performance.  
 468

469 **Minimal Overhead.** To evaluate the inference  
 470 efficiency of the proposed Enhance-A-Video  
 471 (EAV) method, we conducted an ablation study  
 472 on two prevailing video generation models in  
 473 Tab. 3 using 1 A100 GPU. These negligible  
 474 increases in the two models indicate that the  
 475 Enhance-A-Video method is highly efficient and  
 476 scales well when integrated into large video gen-  
 477 eration models. Additional ablation studies are  
 478 provided in Appendix F.

Table 3: Comparison of inference efficiency for HunyuanVideo and CogVideoX models with and without Enhance-A-Video.

Model	Time (min)		Overhead
	w/o EAV	w/ EAV	
HunyuanVideo	50.32	50.72	<b>0.8%</b>
CogVideoX	1.53	1.57	<b>2.1%</b>

## 480 5 CONCLUSION

481 This paper presents Enhance-A-Video, a simple yet effective method that improves temporal consis-  
 482 tency and visual quality in DiT-based video generation. By pioneering the exploration of cross-frame  
 483 information and the temperature concept in DiT blocks, the method offers a straightforward yet pow-  
 484 erful solution for video generation enhancement. Its robust generalization and ease of implementation  
 485 suggest promising future developments in better video generation.

## REPRODUCIBILITY STATEMENT

We employ open-source models and prompts, and provide the code necessary to reproduce our results in the supplementary material. After the blind review period, we will release the complete codebase. Detailed configurations of the evaluated models and the recommended range for the enhance temperature parameter are presented in Appendix E.1 and E.2.

## ETHICS STATEMENT

After reviewing the conference’s ethical guidelines, we believe this work poses no foreseeable ethical concerns. The proposed algorithms focus on general video generation and do not involve human subjects, harmful insights, conflicts of interest, privacy or security issues, legal compliance, or research integrity concerns.

## REFERENCES

- Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, and Dominik Lorenz. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.
- Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators. 2024.
- Jiazi Bu, Pengyang Ling, Pan Zhang, Tong Wu, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Dahua Lin, and Jiaqi Wang. Bytheway: Boost your text-to-video generation model to higher quality in a training-free way. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 12999–13008, 2025.
- Minghong Cai, Xiaodong Cun, Xiaoyu Li, Wenze Liu, Zhaoyang Zhang, Yong Zhang, Ying Shan, and Xiangyu Yue. Ditctrl: Exploring attention control in multi-modal diffusion transformer for tuning-free multi-prompt longer video generation. *arXiv:2412.18597*, 2024.
- Shoufa Chen, Mengmeng Xu, Jiawei Ren, Yuren Cong, Sen He, Yanping Xie, Animesh Sinha, Ping Luo, Tao Xiang, and Juan-Manuel Perez-Rua. Gentrion: Diffusion transformers for image and video generation, 2024. URL <https://arxiv.org/abs/2312.04557>.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. An empirical study of training self-supervised vision transformers. In *International Conference on Computer Vision (ICCV)*, 2021.
- Paul F Christiano, Jan Leike, Tom B Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- Tri Dao, Daniel Y Fu, Stefano Ermon, Atri Rudra, and Christopher Re. Flashattention: Fast and memory-efficient exact attention with IO-awareness. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=H4DqfPSibmx>.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Muller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis. *arXiv preprint arXiv:2403.03206*, 2024. doi: 10.48550/arXiv.2403.03206.
- Peng Gao, Le Zhuo, Dongyang Liu, Ruoyi Du, Xu Luo, Longtian Qiu, Yuhang Zhang, Chen Lin, Rongjie Huang, Shijie Geng, Renrui Zhang, Junlin Xi, Wenqi Shao, Zhengkai Jiang, Tianshuo Yang, Weicai Ye, He Tong, Jingwen He, Yu Qiao, and Hongsheng Li. Lumina-t2x: Transforming text into any modality, resolution, and duration via flow-based large diffusion transformers, 2024. URL <https://arxiv.org/abs/2405.05945>.

- 540 Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala,  
541 Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models  
542 without specific tuning. *International Conference on Learning Representations*, 2024.  
543
- 544 Yoav HaCohen, Nisan Chiprut, Benny Brazowski, Daniel Shalem, Dudu Moshe, Eitan Richardson,  
545 Eran Levin, Guy Shiran, Nir Zabari, Ori Gordon, Poriya Panet, Sapir Weissbuch, Victor Kulikov,  
546 Yaki Bitterman, Zeev Melumian, and Ofir Bibi. Ltx-video: Realtime video latent diffusion. *arXiv  
547 preprint arXiv:2501.00103*, 2024.
- 548 Jingwen He, Tianfan Xue, Dongyang Liu, Xinqi Lin, Peng Gao, Dahua Lin, Yu Qiao, Wanli Ouyang,  
549 and Ziwei Liu. Venhancer: Generative space-time enhancement for video generation, 2024. URL  
550 <https://arxiv.org/abs/2407.07667>.  
551
- 552 Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.  
553 *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2015.  
554 URL <https://api.semanticscholar.org/CorpusID:206594692>.
- 555 Roberto Henschel, Levon Khachatryan, Daniil Hayrapetyan, Hayk Poghosyan, Vahram Tadevosyan,  
556 Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Streamingt2v: Consistent, dynamic,  
557 and extendable long video generation from text. *arXiv preprint arXiv:2403.14773*, 2024.  
558
- 559 Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P.  
560 Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Imagen video: High  
561 definition video generation with diffusion models, 2022a. URL [https://arxiv.org/abs/  
562 2210.02303](https://arxiv.org/abs/2210.02303).
- 563 Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J.  
564 Fleet. Video diffusion models. *arXiv preprint arXiv:2204.03458*, 2022b.  
565
- 566 Ari Holtzman, Jan Buys, Liwei Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text  
567 degeneration. In *International Conference on Learning Representations (ICLR)*, 2020.  
568
- 569 Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing  
570 Wu, Qingyang Jin, Nattapol Chanpaisit, Yaohui Wang, Xinyuan Chen, Limin Wang, Dahua Lin,  
571 Yu Qiao, and Ziwei Liu. VBench: Comprehensive benchmark suite for video generative models.  
572 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- 573 Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin  
574 Li, Bo Wu, Jianwei Zhang, Kathrina Wu, Qin Lin, Junkun Yuan, Yanxin Long, Aladdin Wang,  
575 Andong Wang, Changlin Li, Duojun Huang, Fang Yang, Hao Tan, Hongmei Wang, Jacob Song,  
576 Jiawang Bai, Jianbing Wu, Jinbao Xue, Joey Wang, Kai Wang, Mengyang Liu, Pengyu Li, Shuai  
577 Li, Weiyan Wang, Wenqing Yu, Xincheng Deng, Yang Li, Yi Chen, Yutao Cui, Yuanbo Peng,  
578 Zhentao Yu, Zhiyu He, Zhiyong Xu, Zixiang Zhou, Zunnan Xu, Yangyu Tao, Qinglin Lu, Songtao  
579 Liu, Daquan Zhou, Hongfa Wang, Yong Yang, Di Wang, Yuhong Liu, Jie Jiang, and Caesar  
580 Zhong. Hunyuanvideo: A systematic framework for large video generative models, 2025. URL  
581 <https://arxiv.org/abs/2412.03603>.
- 582 Jiachen Li, Weixi Feng, Tsu-Jui Fu, Xinyi Wang, S Basu, Wenhu Chen, and William Yang Wang. T2v-  
583 turbo: Breaking the quality bottleneck of video consistency model with mixed reward feedback.  
584 In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL  
585 <https://openreview.net/forum?id=53daI9kbvf>.
- 586 Jiachen Li, Qian Long, Jian Zheng, Xiaofeng Gao, Robinson Piramuthu, Wenhu Chen, and  
587 William Yang Wang. T2v-turbo-v2: Enhancing video model post-training through data, re-  
588 ward, and conditional guidance design. In *The Thirteenth International Conference on Learning  
589 Representations*, 2025a. URL <https://openreview.net/forum?id=BZwXMqu4zG>.  
590
- 591 Xiaohui Li, Yihao Liu, Shuo Cao, Ziyang Chen, Shaobin Zhuang, Xiangyu Chen, Yinan He, Yi Wang,  
592 and Yu Qiao. Diffvsr: Enhancing real-world video super-resolution with diffusion models for  
593 advanced visual quality and temporal consistency, 2025b. URL [https://arxiv.org/abs/  
2501.10110](https://arxiv.org/abs/2501.10110).

- 594 Yumeng Li, William H. Beluch, Margret Keuper, Dan Zhang, and Anna Khoreva. VSTAR: Generative  
595 temporal nursing for longer dynamic video synthesis. In *The Thirteenth International Confer-*  
596 *ence on Learning Representations*, 2025c. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=Kb9PnkWYNT)  
597 [Kb9PnkWYNT](https://openreview.net/forum?id=Kb9PnkWYNT).
- 598 Bin Lin, Yunyang Ge, Xinhua Cheng, Zongjian Li, Bin Zhu, Shaodong Wang, Xianyi He, Yang  
599 Ye, Shenghai Yuan, Liuhan Chen, Tanghui Jia, Junwu Zhang, Zhenyu Tang, Yatian Pang, Bin  
600 She, Cen Yan, Zhiheng Hu, Xiaoyi Dong, Lin Chen, Zhang Pan, Xing Zhou, Shaoling Dong,  
601 Yonghong Tian, and Li Yuan. Open-sora plan: Open-source large video generation model, 2024.  
602 URL <https://arxiv.org/abs/2412.00131>.
- 603 Yu Lu, Yuanzhi Liang, Linchao Zhu, and Yi Yang. Freelong: Training-free long video generation with  
604 spectralblend temporal attention, 2024. URL <https://arxiv.org/abs/2407.19918>.
- 605 Guoqing Ma, Haoyang Huang, Kun Yan, Liangyu Chen, Nan Duan, Shengming Yin, Changyi  
606 Wan, Ranchen Ming, Xiaoni Song, and et al. Step-video-t2v technical report: The practice,  
607 challenges, and future of video foundation model, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2502.10248)  
608 [2502.10248](https://arxiv.org/abs/2502.10248).
- 609 Jiajun Ma, Shuchen Xue, Tianyang Hu, Wenjia Wang, Zhaoqiang Liu, Zhenguo Li, Zhi-Ming Ma,  
610 and Kenji Kawaguchi. The surprising effectiveness of skip-tuning in diffusion sampling, 2024a.  
611 URL <https://arxiv.org/abs/2402.15170>.
- 612 Xin Ma, Yaohui Wang, Gengyun Jia, Xinyuan Chen, Ziwei Liu, Yuan-Fang Li, Cunjian Chen,  
613 and Yu Qiao. Latte: Latent diffusion transformer for video generation, 2024b. URL <https://arxiv.org/abs/2401.03048>.
- 614 William S. Peebles and Saining Xie. Scalable diffusion models with transformers. *2023 IEEE/CVF*  
615 *International Conference on Computer Vision (ICCV)*, pp. 4172–4182, 2022. URL <https://api.semanticscholar.org/CorpusID:254854389>.
- 616 Max Peeperkorn, Tom Kouwenhoven, Daniel G. Brown, and Anna K. Jordanous. Is temperature the  
617 creativity parameter of large language models? *ArXiv*, 2024. doi: 10.48550/arXiv.2405.00492.
- 618 Matthew Renze and Erhan Guven. The effect of sampling temperature on problem solving in large  
619 language models. *ArXiv*, abs/2402.05201, 2024. URL <https://api.semanticscholar.org/CorpusID:267547769>.
- 620 Chenyang Si, Ziqi Huang, Yuming Jiang, and Ziwei Liu. Freeu: Free lunch in diffusion u-net. In  
621 *CVPR*, 2024.
- 622 Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry  
623 Yang, Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. Make-a-video:  
624 Text-to-video generation without text-video data, 2022. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2209.14792)  
625 [2209.14792](https://arxiv.org/abs/2209.14792).
- 626 Cheng Tan, Zhangyang Gao, Lirong Wu, Yongjie Xu, Jun Xia, Siyuan Li, and Stan Z. Li. Tem-  
627 poral attention unit: Towards efficient spatiotemporal predictive learning. In *Proceedings of the*  
628 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 18770–18782,  
629 2023.
- 630 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz  
631 Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1706.03762)  
632 [1706.03762](https://arxiv.org/abs/1706.03762).
- 633 Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kindermans, Hernan Moraldo, Han Zhang, Mo-  
634 hammad Taghi Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. Phenaki: Variable length  
635 video generation from open domain textual descriptions. In *International Conference on Learning*  
636 *Representations*, 2023. URL <https://openreview.net/forum?id=vOEXS39nOF>.
- 637 Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu,  
638 Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, and  
639 et al. Wan: Open and advanced large-scale video generative models, 2025. URL <https://arxiv.org/abs/2503.20314>.
- 640  
641  
642  
643  
644  
645  
646  
647

- 648 Xiang Wang, Shiwei Zhang, Han Zhang, Yu Liu, Yingya Zhang, Changxin Gao, and Nong Sang.  
649 Videolcm: Video latent consistency model, 2023. URL [https://arxiv.org/abs/2312.](https://arxiv.org/abs/2312.09109)  
650 [09109](https://arxiv.org/abs/2312.09109).
- 651 Tian Xia, Xuweiyi Chen, and Sihan Xu. Unictrl: Improving the spatiotemporal consistency of  
652 text-to-video diffusion models via training-free unified attention control. *Transactions on Machine*  
653 *Learning Research*, 2024. ISSN 2835-8856. URL [https://openreview.net/forum?](https://openreview.net/forum?id=x2uFJ790jK)  
654 [id=x2uFJ790jK](https://openreview.net/forum?id=x2uFJ790jK).
- 655 Jiaqi Xu, Xinyi Zou, Kunzhe Huang, Yunkuo Chen, Bo Liu, MengLi Cheng, Xing Shi, and Jun  
656 Huang. Easyanimate: A high-performance long video generation method based on transformer  
657 architecture, 2024a. URL <https://arxiv.org/abs/2405.18991>.
- 658 Yiran Xu, Taesung Park, Richard Zhang, Yang Zhou, Eli Shechtman, Feng Liu, Jia-Bin Huang,  
659 and Difan Liu. Videogigagan: Towards detail-rich video super-resolution, 2024b. URL [https:](https://arxiv.org/abs/2404.12388)  
660 [//arxiv.org/abs/2404.12388](https://arxiv.org/abs/2404.12388).
- 661 Zhongcong Xu, Jianfeng Zhang, Jun Hao Liew, Hanshu Yan, Jia-Wei Liu, Chenxu Zhang, Jiashi  
662 Feng, and Mike Zheng Shou. Magicanimate: Temporally consistent human image animation using  
663 diffusion model. In *arXiv*, 2023.
- 664 Wilson Yan, Danijar Hafner, Stephen James, and Pieter Abbeel. Temporally consistent trans-  
665 formers for video generation. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara  
666 Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International*  
667 *Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp.  
668 39062–39098. PMLR, 23–29 Jul 2023. URL [https://proceedings.mlr.press/v202/](https://proceedings.mlr.press/v202/yan23b.html)  
669 [yan23b.html](https://proceedings.mlr.press/v202/yan23b.html).
- 670 Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang,  
671 Wenyi Hong, Xiaohan Zhang, Guanyu Feng, Da Yin, Xiaotao Gu, Yuxuan Zhang, Weihang Wang,  
672 Yean Cheng, Ting Liu, Bin Xu, Yuxiao Dong, and Jie Tang. Cogvideox: Text-to-video diffusion  
673 models with an expert transformer, 2024. URL <https://arxiv.org/abs/2408.06072>.
- 674 Shenghai Yuan, Jinfa Huang, Yujun Shi, Yongqi Xu, Ruijie Zhu, Bin Lin, Xinhua Cheng, Li Yuan,  
675 and Jiebo Luo. Magictime: Time-lapse video generation models as metamorphic simulators, 2025.  
676 URL <https://arxiv.org/abs/2404.05014>.
- 677 David Junhao Zhang, Jay Zhangjie Wu, Jia-Wei Liu, Rui Zhao, Lingmin Ran, Yuchao Gu, Difei  
678 Gao, and Mike Zheng Shou. Show-1: Marrying pixel and latent diffusion models for text-to-video  
679 generation. *arXiv preprint arXiv:2309.15818*, 2023a.
- 680 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image  
681 diffusion models, 2023b.
- 682 Yuang Zhang, Jiayi Gu, Li-Wen Wang, Han Wang, Junqi Cheng, Yuefeng Zhu, and Fangyuan Zou.  
683 Mimicmotion: High-quality human motion video generation with confidence-aware pose guidance.  
684 *arXiv preprint arXiv:2406.19680*, 2024.
- 685 Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou,  
686 Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all, 2024. URL  
687 <https://arxiv.org/abs/2412.20404>.
- 688 Shangchen Zhou, Peiqing Yang, Jianyi Wang, Yihang Luo, and Chen Change Loy. Upscale-A-Video:  
689 Temporal-consistent diffusion model for real-world video super-resolution. In *CVPR*, 2024.
- 690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

## 702 A USAGE OF LARGE LANGUAGE MODELS (LLMs)

703  
704 In this paper, the ideas, experiments, and manuscript writing were all carried out by the authors. LLMs  
705 were only used for minor language polishing, without contributing to the conception, experimental  
706 design, or substantive content of the work.  
707

## 708 B LIMITATIONS AND FUTURE WORK

709 In our approach, the enhancement temperature parameter needs to be manually tuned for each DiT-  
710 based text-to-video model. In future work, we plan to develop an adaptive temperature mechanism  
711 using RLHF (Christiano et al., 2017; Li et al., 2024; 2025a) to adjust this parameter based on the  
712 specific prompt context automatically. Besides, we focused solely on enhancing temporal attention  
713 without addressing spatial attention or cross-attention mechanisms, which are crucial for preserving  
714 spatial coherence and prompt alignment. Future work could explore incorporating these mechanisms  
715 to improve spatial video quality and semantic consistency.  
716  
717

## 718 C TEMPERATURE METHOD COMPARISON

719 In Fig. 10(a) and (b), where the temperature parameter  $\tau$  and Cross-Frame Intensity are directly  
720 applied in temporal attention calculation separately as presented in Equation 9 and 10, the diagonal  
721 elements (e.g., 27.4, 6.3) show a significant weakening of intra-frame attention, leading to the severe  
722 loss of spatial details and resulting in blurry and unrealistic textures. Additionally, the large negative  
723 values in the off-diagonal regions indicate an overabundant distributed enhancement of cross-frame  
724 attention, resulting in limited improvement in video quality.  
725  
726  
727

$$728 \text{Attn}(Q, K, V) = \text{softmax} \left( \frac{QK^\top}{\tau \cdot \sqrt{d_k}} \right) V \quad (9)$$

$$729 \text{Attn}(Q, K, V) = \text{softmax} \left( \frac{QK^\top}{CFI_{enhanced} \cdot \sqrt{d_k}} \right) V \quad (10)$$

730 In contrast, Fig. 10(c) using the Enhance-A-Video method shows modest changes along the diagonal,  
731 with values close to zero, preserving intra-frame attention and maintaining fine-grained details.  
732 Moreover, the negative values in the off-diagonal regions (e.g., -1.3, -0.9) reflect a targeted and  
733 moderate enhancement of cross-frame attention, significantly improving motion coherence and  
734 overall video quality.  
735  
736  
737  
738

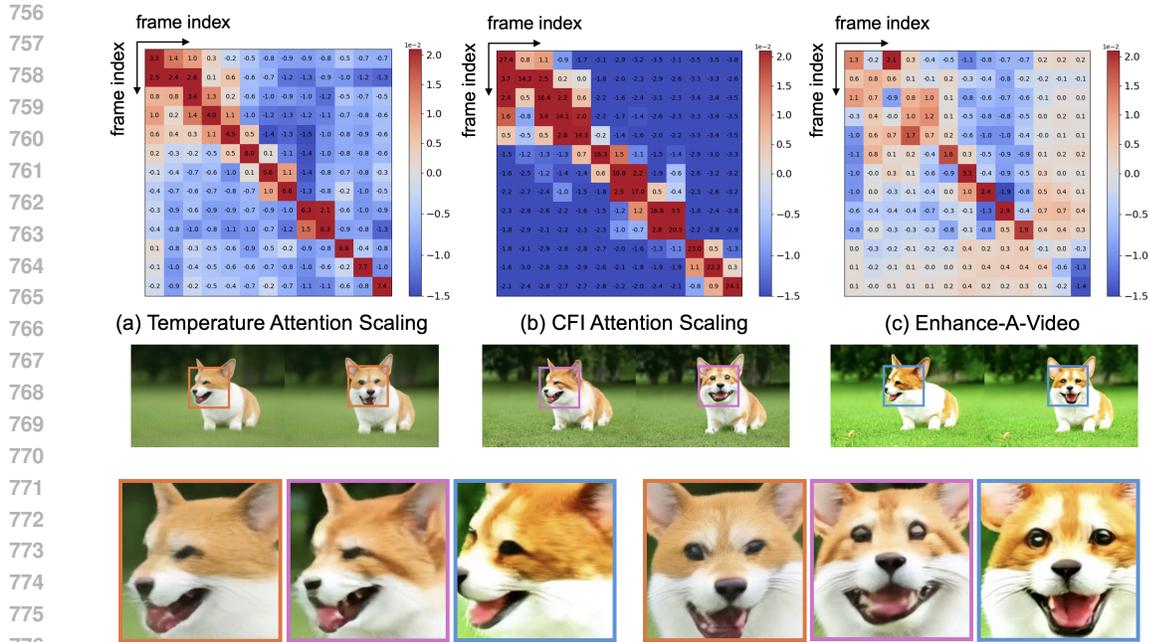
## 739 D CFI DISTRIBUTION AND L2 NORM PROPORTION IN RESIDUAL CONNECTION

740 The  $CFI_{enhanced}$  values in Fig. 11(a) range between 1.12 – 1.18, indicating a modest enhancement  
741 of keyframes containing important temporal information. Fig. 11(b) shows two low proportions  
742 calculated as follows:  
743

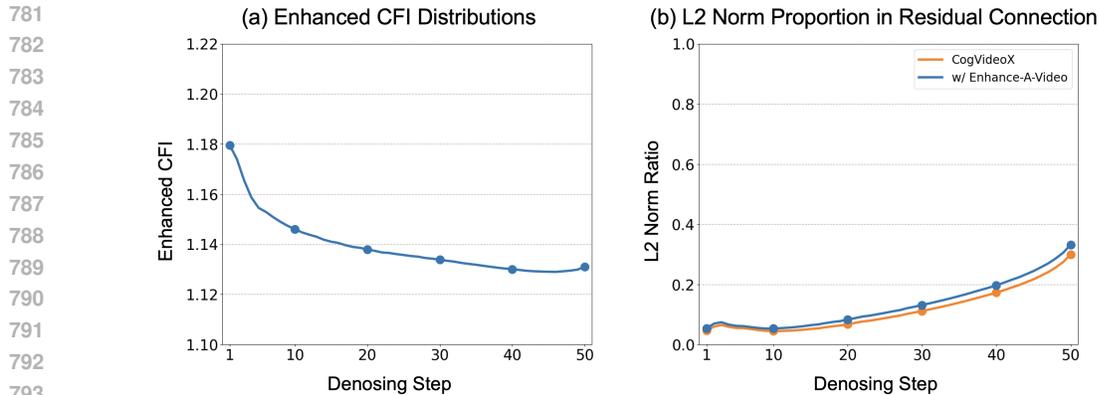
$$744 \text{prop}_{\text{CogvideoX}} = \frac{\|\mathbf{O}_{\text{attn}}\|_2}{\|\mathbf{H}\|_2} \quad (11)$$

$$745 \text{prop}_{\text{w/ Enhance-A-Video}} = \frac{\|CFI_{enhanced} \cdot \mathbf{O}_{\text{attn}}\|_2}{\|\mathbf{H}\|_2} \quad (12)$$

746 suggesting that attention outputs are relatively small compared to hidden states in the residual  
747 connection. Consequently, applying  $CFI_{enhanced}$  to attention outputs rather than attention allows for  
748 enhancing important information with minimal disruption to the original attention distribution. Thus,  
749 Enhance-A-Video improves temporal consistency while preserving existing spatial details.  
750  
751  
752  
753  
754  
755



777 Figure 10: Temporal attention difference maps and corresponding generated videos comparing three  
778 temperature enhancement methods. (a) Temperature Attention Scaling  $\tau = 1.1$ . (b) CFI Attention  
779 Scaling. (c) Enhance-A-Video Method.



794 Figure 11: (a) The distribution of  $CFI_{enhanced}$  during the inference of CogVideoX w/ Enhance-A-  
795 Video in layer 4. (b) The proportion of  $l_2$  norms between  $\mathbf{O}_{attn}$  and  $\mathbf{H}$  in residual connection in layer  
796 4.

## 799 E MORE EXPERIMENTAL INFORMATION

### 801 E.1 EVALUATED MODELS AND SETTINGS

802 We evaluated Enhance-A-Video using widely adopted open-source models with default baseline  
803 parameters. Detailed settings are given below, and all experiments were conducted on NVIDIA H100  
804 or A100-80GB GPUs.

### 807 E.2 IDEAL RANGE FOR ENHANCE TEMPERATURE PARAMETER

808 To guide users in selecting an appropriate configuration, we provide recommended ranges for the  
809 enhance temperature parameter  $\tau$  tailored to each foundation model. These ranges are determined

Table 4: Specific settings of evaluated models.

Model	Size	Video Resolution	Duration	Frame Rate (FPS)
Wan2.1	14B	480p	5s	16
HunyuanVideo	13B	720p	5s	24
CogVideoX	5B	480p	6s	8
LTX-Video	1.9B	512p	5s	24
Open-Sora 1.2	1.1B	480p	4s	16
Open-Sora-Plan v1.1.0	1.2B	512p	2s	30

based on empirical observations that balance generation quality and consistency across diverse prompts. By offering model-specific  $\tau$  intervals in Tab. 5, we facilitate more efficient tuning and ensure that users can achieve optimal results without exhaustive manual search.

Table 5: Recommended  $\tau$  ranges for each foundation model.

Model	Range
HunyuanVideo	<b>3–5</b>
CogVideoX-2B	<b>1–3</b>
LTX-Video	<b>5–7</b>
Open-Sora	<b>1–3</b>

### E.3 DETAILED SETTINGS FOR QUANTITATIVE ANALYSIS

To evaluate the effectiveness of Enhance-A-Video on different foundation models in practical scenarios, we conducted a user study based on 15 prompts sampled from the VBench benchmark (Huang et al., 2024). The number of evaluated samples per model is summarized in Tab. 6. We selected HunyuanVideo as the primary model for our quantitative analysis to ensure a more representative evaluation of Enhance-A-Video’s effectiveness on a high-performing open-source baseline.

Table 6: Number of evaluated samples per model in the user study.

Model	Number of Samples
HunyuanVideo	<b>9</b>
CogVideoX	<b>2</b>
LTX-Video	<b>2</b>
Open-Sora	<b>2</b>

### E.4 COMPARISON WITH UNET-BASED GENERATION ENHANCEMENT METHODS

We selected two representative UNet-based training-free enhancement methods for comparison with our approach. Since Flash Attention prevents direct access and modification of temporal attention in existing 3D attention models, our comparison primarily focused on Open-Sora, which employs spatial-temporal attention.

#### E.4.1 COMPARISON WITH BYTHEWAY

As shown in 7, existing UNet-based methods fail to effectively enhance video generation due to structural differences in UNet and DiT models, whereas our first DiT-oriented approach resolves this limitation without incurring additional cost.

Table 7: VBench results comparing baseline models, w/ ByTheWay and w/ Enhance-A-Video (EAV) across 5 evaluation criteria.

Model	Subject Consistency $\uparrow$	Temporal Flickering $\uparrow$	Imaging Quality $\uparrow$	Multiple Objects $\uparrow$	Spatial Relationship $\uparrow$
Open-Sora	94.14	98.33	60.43	49.54	60.19
w/ ByTheWay	94.19	98.33	60.41	47.72	60.24
w/ EAV	<b>94.27</b>	<b>98.92</b>	<b>61.04</b>	<b>53.92</b>	<b>64.58</b>

#### E.4.2 COMPARISON WITH FREEU

We experimented to assess the performance of FreeU on the Open-Sora framework using its default settings ( $b_1 = 1.1, b_2 = 1.1$ ). As illustrated in Fig. 12, applying FreeU only to the last two layers results in overly smooth videos with noticeable loss of detail, while extending it to more layers (e.g., four) causes the output to lose nearly all meaningful content. These results suggest that FreeU is not well-suited for DiT-based video generation, as it was originally developed for UNet-based image generation and does not effectively capture temporal dependencies or adapt to architectural differences. In contrast, Enhance-A-Video is specifically designed for DiT-based video models, efficiently addressing temporal attention while maintaining inference overhead comparable to FreeU’s usage in image generation.

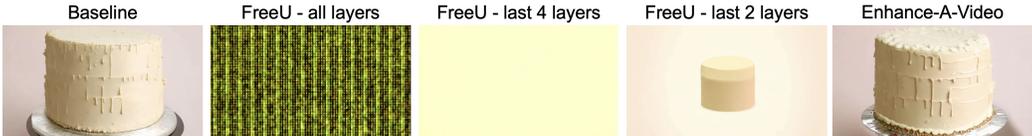


Figure 12: Qualitative results of FreeU on Open-Sora. FreeU fails to preserve structure when applied to DiT, especially with more layers. Enhance-A-Video improves visual detail and temporal consistency.

## F MORE ABLATION STUDY

### F.1 ABLATION STUDY ON CLASSIFIER-FREE GUIDANCE

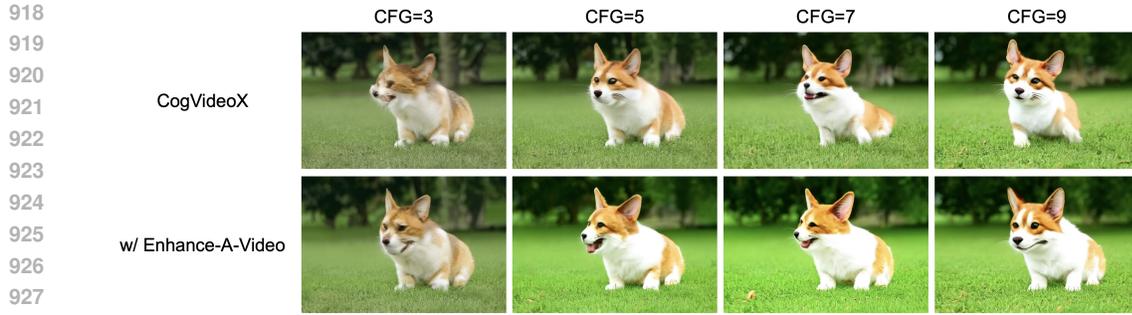
In our experiments, we adopted the default classifier-free guidance (CFG) values specified for each foundation model, as summarized in Tab. 8:

Table 8: Default classifier-free guidance (CFG) settings for each foundation model.

Model	Default CFG
HunyuanVideo	<b>6.0</b>
CogVideoX	<b>6.0</b>
LTX-Video	<b>3.0</b>
Open-Sora	<b>7.0</b>

To further investigate the influence of CFG on the effectiveness of our method, we conducted ablation studies using CogVideoX by generating a 6-second, 480p video with the prompt “A cute happy Corgi playing in park”. We varied the CFG values across 3.0, 5.0, 7.0, and 9.0, applying Enhance-A-Video in each case.

As shown in Fig. 13, the impact of Enhance-A-Video is highly dependent on the base video quality determined by the CFG setting. At higher CFG values (7.0 and 9.0), where the base outputs are of relatively high quality, Enhance-A-Video consistently improves spatial fidelity without introducing artifacts. However, at lower CFG values (3.0 and 5.0), where the base generation quality is poor, Enhance-A-Video is limited in its ability to recover detail and consistency, as it is designed to enhance—not replace—the underlying generative capacity of the model.



929 Figure 13: Qualitative results of Enhance-A-Video with different CFG settings on CogVideoX.  
 930 Caption: *A cute and happy Corgi playing in the park.*

931

932

933

934

935

936

937

938

939

940

941

942

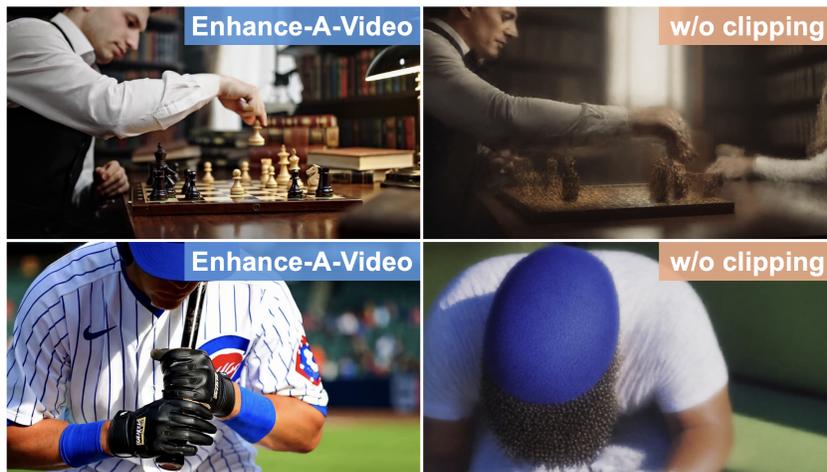
943

944

945

946

947



948 Figure 14: Visual comparison of video generation results with and without the clipping mechanism  
 949 in the Enhance Block.

## 951 F.2 EFFECTS OF CLIPPING.

952

953 Fig. 14 illustrates that applying the clipping effectively stabilizes cross-frame attention, resulting in  
 954 clearer visuals and smoother motion. Without clipping, the model produces noticeable artifacts such  
 955 as motion blur and distorted details, highlighting the necessity of clipping for maintaining temporal  
 956 consistency and preserving spatial fidelity.

## 958 F.3 DISCUSSION ON LONG VIDEO GENERATION

959

960 We employed HunyuanVideo to produce a 9-second video comprising 217 frames—well beyond its  
 961 default limit of 129 frames—and assessed the performance of Enhance-A-Video on the video.

962

963

964

965

966

967

968

969

970

971

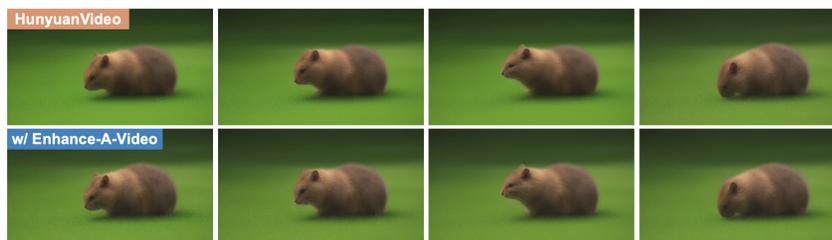


Figure 15: Visual comparison of video generation results on long video generation.

972 The video generated by HunyuanVideo in Fig. 15 shows poor visual quality, with blurred spatial  
973 details and low prompt consistency (e.g., the output resembles a mouse rather than a cat). This  
974 degradation stems from the model’s limited capacity to handle longer videos, likely due to the  
975 insufficient availability of high-quality long video data during training. Addressing this issue would  
976 require improvements in training data, model architecture, and overall scale.

977 Since Enhance-A-Video operates solely during the inference stage in a training-free way, its effective-  
978 ness depends heavily on the original video quality. While it can enhance decent videos with more  
979 details or better temporal consistency, it struggles when the input is of very low quality - as in this  
980 case - because the core limitations lie in the training process, not inference.

## 982 G CAPTIONS FOR FIGURE 1

983  
984 Caption 1 (top row): A young girl with curly hair, wearing a bright yellow dress, sits cross-legged on  
985 a wooden floor, surrounded by an array of colorful markers and crayons. She carefully colors a large  
986 piece of cardboard, her face a picture of concentration and creativity. The cardboard, propped up  
987 against a cozy living room couch, is filled with whimsical drawings of flowers, stars, and animals.  
988 Sunlight streams through a nearby window, casting a warm glow over her workspace. Her small  
989 hands move deftly, adding vibrant hues to her imaginative artwork, while her expression reflects pure  
990 joy and artistic focus.

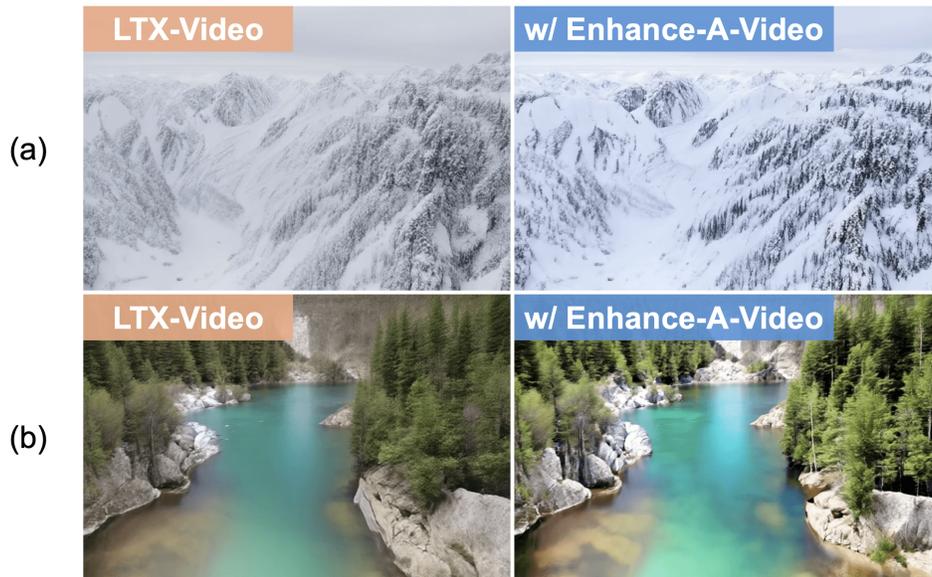
991 Caption 2 (bottom row): A young girl, wearing a wide-brimmed straw hat and a colorful swimsuit,  
992 carefully applies sunblock to her younger brother’s face on a sunlit beach. The boy, with **sandy hair**  
993 and a playful grin, sits patiently on a striped beach towel, surrounded by **sandcastles** and **beach toys**.  
994 The gentle waves of the ocean provide a soothing soundtrack as seagulls call in the distance. The  
995 girl’s hands move with care, ensuring every inch of his face is protected, while the sun casts a warm  
996 glow over the scene, highlighting the siblings’ bond and the carefree joy of a summer day by the sea.  
997

## 998 H MORE QUALITATIVE RESULTS

999  
1000 LTX-Video (HaCohen et al., 2024) is a real-time latent text-to-video diffusion model that generates  
1001 high-quality, temporally consistent videos efficiently. The integration of Enhance-A-Video into  
1002 LTX-Video further improves temporal consistency and enhances spatial details. As exhibited in Fig.  
1003 16, the enhanced model produces videos with sharper textures, more vivid colors, and smoother  
1004 transitions compared to the baseline LTX-Video.

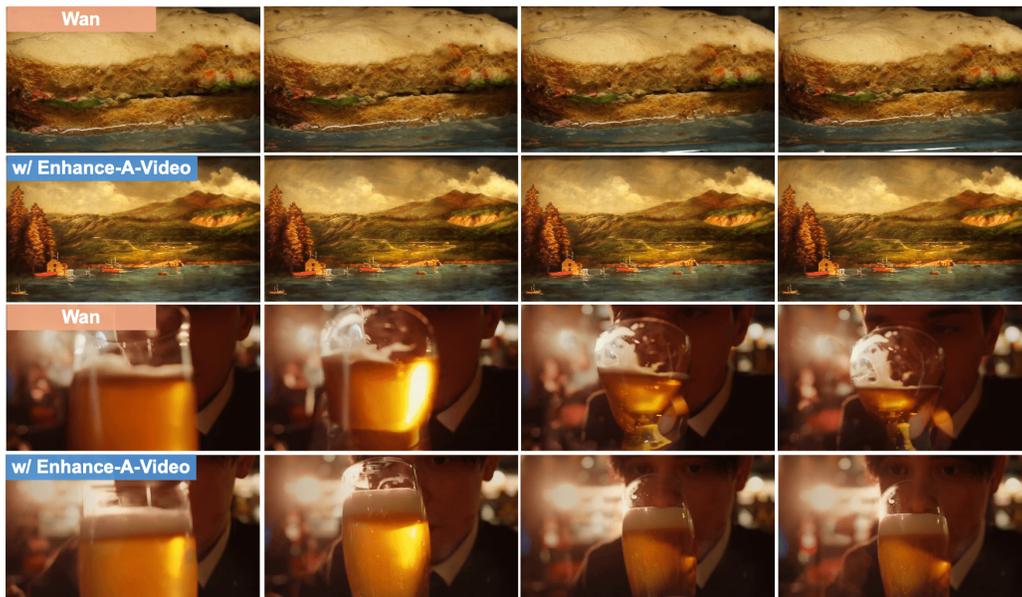
1005 The snow-covered mountains (top row) and river scene (bottom row) generated by Enhance-A-  
1006 Video display clearer structures and more natural color gradients, while the baseline results appear  
1007 less detailed and slightly blurred. This demonstrates that Enhance-A-Video effectively strengthens  
1008 cross-frame attention, leading to more realistic and visually appealing videos.  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046



1047 Figure 16: Qualitative results of Enhance-A-Video on LTX-Video. Captions: *(a) The camera pans over snow-covered mountains, revealing jagged peaks and deep, narrow valleys. (b) An emerald-green river winds through a rocky canyon, forming reflective pools amid pine trees and brown-gray rocks.*

1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074



1075 Figure 17: More qualitative results of Enhance-A-Video on Wan. Captions: *(a) A tranquil tableau of bar. (b) A person is tasting beer.* The integration of Enhance-A-Video facilitates prompt-consistent video generation while ensuring smoother and more natural transitions.

1078  
1079

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

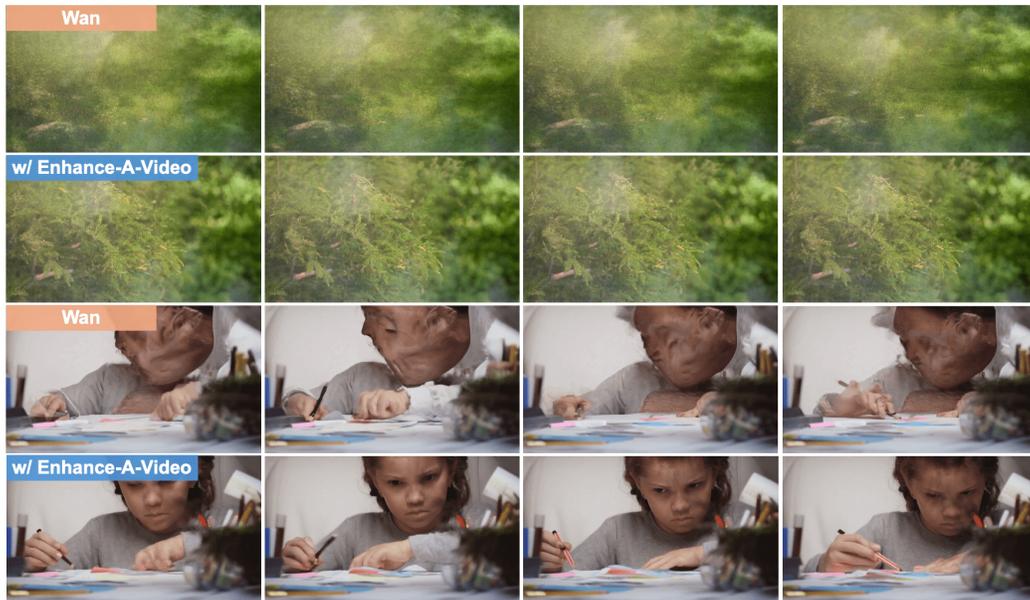


Figure 18: More qualitative results of Enhance-A-Video on Wan. Captions: (a) *A tranquil tableau of alley.* (b) *A person is drawing.* The use of Enhance-A-Video enhances image quality, allowing the model to more effectively leverage contextual information.

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

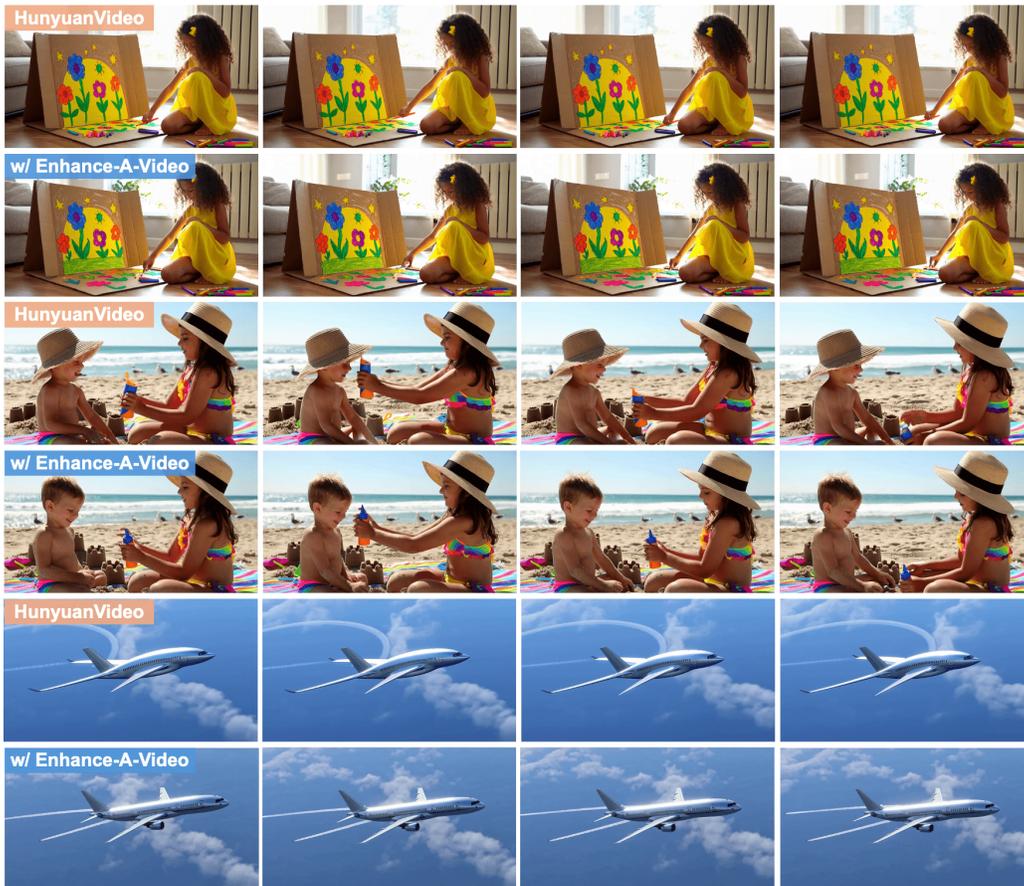
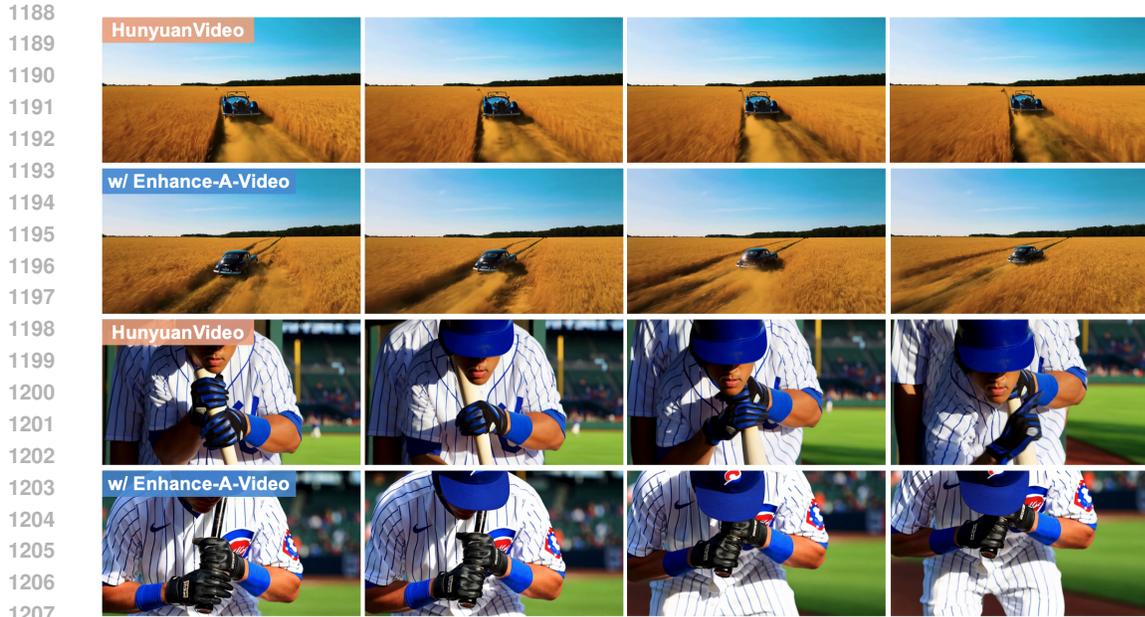
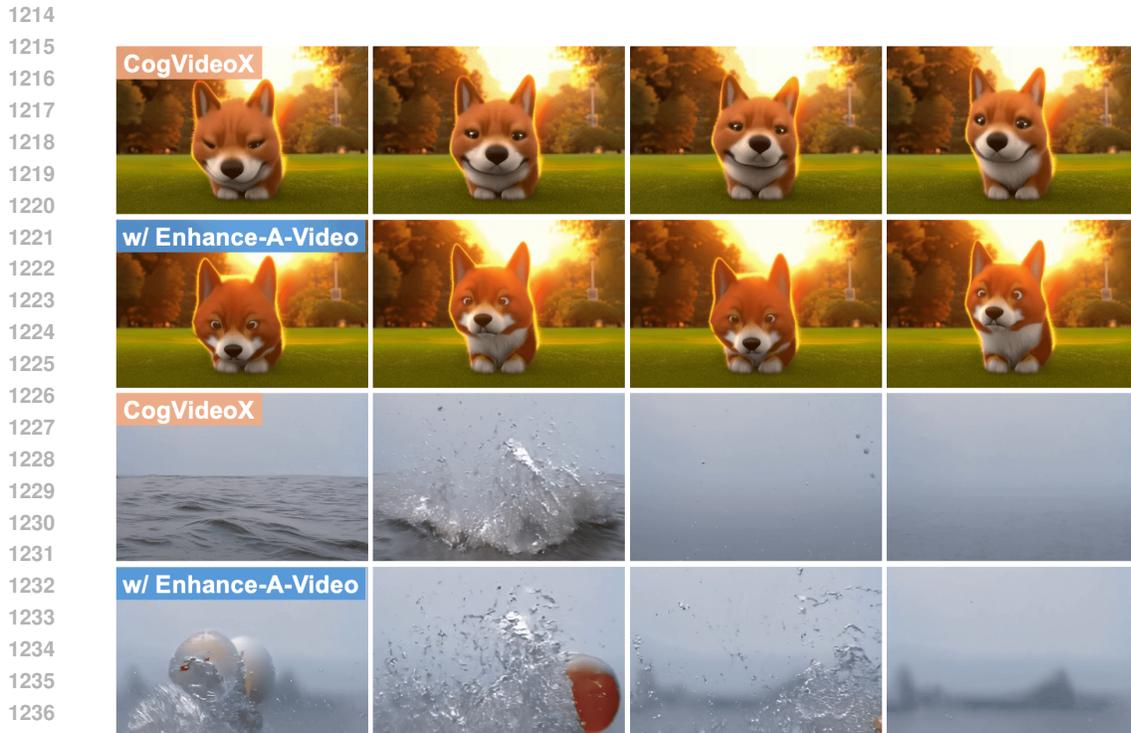


Figure 19: More qualitative results of Enhance-A-Video on HunyuanVideo. The application of Enhance-A-Video enriches visual details and ensures prompt-consistent video generation, resulting in more realistic outputs.



1208 Figure 20: More qualitative results of Enhance-A-Video on HunyuanVideo. Captions: *(a)* An antique  
1209 car drives along a dirt road through golden wheat fields. Dust rises softly as wheat brushes against  
1210 the car with distant trees meeting a blue sky. *(b)* A baseball player grips a bat in black gloves,  
1211 wearing a blue-and-white uniform and cap, with a blurred crowd and green field highlighting his  
1212 focused stance. Enhance-A-Video consistently generates more realistic frames with finer details and  
1213 more natural motion.



1238 Figure 21: More qualitative results of Enhance-A-Video on CogVideoX. Captions: *(a)* A cute and  
1239 happy Corgi playing in the park, in a surrealistic style. *(b)* Balloon full of water exploding in extreme  
1240 slow motion. The enhanced model generates videos that align more closely with the given prompts,  
1241 while providing smoother transitions and sharper visuals.