VENHANCER: GENERATIVE SPACE-TIME ENHANCE MENT FOR VIDEO GENERATION

Anonymous authors

Paper under double-blind review

ABSTRACT

We present *VEnhancer*, a generative space-time enhancement method that can improve the existing AI-generated videos spatially and temporally through one video diffusion model. Given a generated low-quality video, our approach can increase its spatial and temporal resolution simultaneously with arbitrary up-sampling space and time scales by adding more details in spatial domain and synthesize detailed motion in temporal domain. Furthermore, VEnhancer is able to remove generated spatial artifacts and temporal flickering of generated videos. To achieve this, basing on a pretrained generative video prior, we train a Space-Time Controller and inject it to the prior as a condition on low-frame-rate and low-resolution videos. To effectively train this ST-Controller, we design *space-time data augmentation* to create diversified video training pairs as well as *video-aware conditioning* for realizing different augmentation parameters in both spatial and temporal dimensions. Benefiting from the above designs, VEnhancer can be end-to-end trained to enable multi-function in one single model. Extensive experiments show that VEnhancer surpasses existing state-of-the-art video super-resolution and space-time super-resolution methods in enhancing AI-generated videos. Moreover, VEnhancer is able to greatly improve the performance of open-source state-of-the-art text-to-video methods on video generation benchmark, VBench.



Figure 1: The enhanced screenshots for AI-generated videos (from Kling). Prompt: *Einstein plays guitar*. I2VGen-XL (refiner) (Zhang et al., 2023b) has successfully refined the video by removing distortions/artifacts, but suffers from severe identity change (*e.g.*, facial attributes) and blurry results. LaVie-SR (Wang et al., 2023b) could produce high-fidelity results but lacks generative ability in modifying and regenerating video content (*e.g.*, correcting the guitar strings). In contrast, our method could achieve effective refinement and output high-resolution videos with realistic texture details and good identity preservation. Zoom in for best view.

INTRODUCTION

With the advances of text-to-image generation (Rombach et al., 2022; Podell et al., 2023; Chen et al., 2023b; Gao et al., 2024) and large-scale video datasets with text description (Bain et al., 2021), there is fast development of text-to-video generative models (Guo et al., 2023; Chen et al., 2023c; Ho

et al., 2022; Blattmann et al., 2023b; Wang et al., 2023b;a; Chen et al., 2023a; 2024a; Gupta et al., 055 2023; Wu et al., 2023b). These developments enable users to generate compelling videos through 056 textual descriptions of the desired content. One common solution (Ho et al., 2022; Blattmann et al., 057 2023b; Wang et al., 2023b; Blattmann et al., 2023a; Gupta et al., 2023) to obtain high-quality videos 058 is to adopt cascaded pipelines, which stacks several video diffusion models, including text-to-video, temporal super-resolution and spatial super-resolution (S-SR) diffusion. These pipeline significantly reduce computation cost when generating high-resolution and high-frame-rate videos, but also pose 060 several issues. First, using two different models for spatial and temporal enhancement separately is 061 redundant, as they are strongly correlated. Also, the proposed diffusion-based spatial or temporal 062 super-resolution (Blattmann et al., 2023b; Ho et al., 2022; Blattmann et al., 2023a; Wang et al., 063 2023b; Lin et al., 2024) have a limited flexibility, as they can only handle fixed interpolation ratio 064 (*i.e.*, predicting three frames between two consecutive frames) or fixed upscaling factors (*i.e.*, $4 \times$). 065 Second, directly training diffusion models on synthesized video pairs may limit its generalization 066 ability, as these models tailored on classic super-resolution task hallucinate high-frequency details, 067 but cannot semantically improve the visual quality of input videos, such as eliminating distortions, 068 artifacts, or recreating new contents (see LaVie-SR's results in Figure 1).

069 Another common approach is training another diffusion model to remove artifacts and to refine 070 distorted content in the generated videos (Henschel et al., 2024; Zhang et al., 2023b). One example is 071 I2VGen-XL (Zhang et al., 2023b), which follows the idea of image generation model SDXL (Podell 072 et al., 2023)- it first upscales videos to higher resolutions using bilinear interpolation, and trains 073 another diffusion refinement model, which will be used for video regeneration through a noising-074 denoising process (Meng et al., 2021). However, this method usually produces over-smoothed videos 075 without realistic texture details (see Figure 1, I2VGen-XL), since the adopted bilinear upsampling cannot generate more spatial details. More importantly, the noising-denoising process (i.e., starting 076 from t = 600) will substantially change the original video content, which cannot always be acceptable 077 in practical applications. 078

In short, current generative video enhancement methods face several challenges. First, sequentially applying temporal and spatial super-resolution is redundant, as they are independently trained, but using similar training datasets. Thus, such design is both sub-optimal and inefficient during the inference. Second, existing refinement methods struggle in balancing between video quality and fidelity to the original content. More importantly, they cannot perform effective super-resolution for increasing spatial and temporal details, which limits their practicality. Third, previous generative video enhancement methods lack the flexibility in dealing with different upscaling factors and refinement strengths for spatial or temporal super-resolution and video refinement.

087 To this end, we propose VEnhancer, a generative space-time enhancement method that supports both 880 spatial and temporal super-resolution with flexible space and time up-sampling scales, as well as has the ability to remove visual artifacts and flickering with good maintenance of video content. It 089 is built upon a pretrained and fixed generative video prior (Zhang et al., 2023b), which supplies the 090 generative ability for video enhancement. To condition the video generation on low-frame-rate and 091 low-resolution videos, we design Space-Time Controller (ST-Controller) for effective conditioning in 092 both spatial and time dimensions. Furthermore, to handle different up-sampling scales and reduce 093 artifacts or flickering with different degrees, we propose a space-time data augmentation algorithm 094 to construct the training data. In particular, at the training stage, we sample different step sizes for skipping frames, downscaling factors, and noise levels to synthesize diversified condition videos. To 096 ensure the proposed ST-Controller be aware of the associated data augmentation applied to each input 097 video, we propose the video-aware conditioning. In particular, for key frame, the condition latent, 098 the embeddings of the associated downscaling factor s and noise level σ by noise augmentation are incorporated into ST-Controller through video-aware conditioning. 099

100 With these designs, VEnhancer is a single end-to-end trainable network that can handle both spatial 101 and temporal super-resolution, as well as video refinement. Moreover, it also supports arbitrary 102 space and time up-sampling scales, and also supports flexible control on refinement strength and 103 generative strength as user may prefer. Extensive experiments have demonstrated VEnhancer's 104 ability in enhancing generated videos (see Figure 1). In these experiments, it outperforms stateof-the-art real-world and generative video super-resolution methods for spatial super-resolution 105 only. In the space-time super-resolution, VEnhancer also surpasses state-of-the-art methods as well 106 as cascaded spatial and temporal diffusion super-resolution models. At last, on the public video 107

generation benchmark, VBench (Huang et al., 2023), VEnhancer can significantly improve the overall performance of existing text-to-video algorithms.

Our contributions can be summarized as below:

- 1. We propose VEnhancer, a generative space-time enhancement method that can achieve generative spatial and temporal super-resolution for different upsampling factors, as well as controllable video refinement in one video diffusion model for the first time.
- 2. To achieve the unified generative space-time enhancement, we devise ST-Controller for effective multi-frame condition injection based on a pretrained and fixed generative video prior. Besides, space-time data augmentation and the associated video-aware conditioning are proposed for training ST-Controller in an end-to-end manner.
- 3. VEnhancer surpasses existing state-of-the-art video super-resolution methods and spacetime super-resolution methods in enhancing generated videos. Also, it could improve the performance of open-source text-to-video methods on public video generation benchmark.
- 121 122 123

124

112

113

114

115

116

117

118

119

120

2 RELATED WORK

125 2.1 VIDEO GENERATION

126 Recently, there have been substantial efforts in training large-scale T2V (Wang et al., 2024; Ho 127 et al., 2022; Guo et al., 2023; Chen et al., 2023c; Gupta et al., 2023; Blattmann et al., 2023b; Wang 128 et al., 2023b;a) models on large scale datasets. Some works (Blattmann et al., 2023b; Wang et al., 129 2023b;a) inflate a pre-trained text-to-image (T2I) model by inserting temporal layers and fine-tuning 130 them or all parameters on video data, or adopts a joint image-video training strategy. In order to 131 achieve high-quality video generation, (Ho et al., 2022; Wang et al., 2023b; Blattmann et al., 2023b) 132 adopts multi-stage pipelines. In particular, cascaded video diffusion models are designed: One T2V base model that is followed by one or more frame interpolation and video super-resolution models. 133 VideoLDM (Blattmann et al., 2023b), LaVie (Wang et al., 2023b), and Upscale-A-Video (Zhou et al., 134 2023) all develop the video super-resolution model based on $4 \times$ sd (StableDiffusion (Rombach et al., 135 2022)) upscaler, which has an additional downsampled image for conditioning the generation. One 136 drawback of this base model is losing quite a lot generative ability compared with T2I base models. 137 On the contrary, I2VGEN-XL follows SDXL (Podell et al., 2023) and uses noising-denoising process 138 (Meng et al., 2021) to refine the generated artifacts. However, this strategy could improve stability 139 but cannot increase the space-time resolution. VEnhancer is based on a generative video prior, and 140 could address temporal/spatial super-resolution and refinement in a unified model.

141 142 143

2.2 VIDEO ENHANCEMENT

144 Video Super-Resolution. Video Super-Resolution (VSR) is proposed to enhance video quality 145 by upsampling low-resolution (LR) frames into high-resolution (HR) ones. Traditional VSR ap-146 proaches(Cao et al., 2021; Chan et al., 2021; 2022a; Isobe et al., 2020a;b;c; Liang et al., 2024; 2022; 147 Wang et al., 2019; Xue et al., 2019) often rely on fixed degradation models to synthesize training data pairs, which leads to a noticeable performance drop in real-world scenarios. To bridge this 148 gap, recent advances(Chan et al., 2022b; Xie et al., 2023) in VSR have embraced more diversified 149 degradation models to better simulate real-world low-resolution videos. To achieve photo-realistic 150 reconstruction, Upscale-A-Video(Zhou et al., 2023) integrates diffusion prior to produce detailed 151 textures, upgrading VSR performance into next level. Space-Time Super-Resolution. Space Time 152 Video Super-Resolution (STVSR) aims to simultaneously increase the resolutions of video frames 153 in both spatial and temporal dimensions. Deep-learning based approaches(Haris et al., 2020; Kim 154 et al., 2020; Xiang et al., 2020; Chen et al., 2022) have achieved remarkable results on STVSR. 155 STARNet(Haris et al., 2020) increases the spatial resolution and frame rate by leveraging the mutual 156 information between space and time. FISR(Kim et al., 2020) propose a joint framework with a multi-157 scale temporal loss to upscale the spatial-temporal resolution of videos. (Xiang et al., 2020) proposes 158 a one-stage STVSR framework, which incorporates different sub-modules for LR frame features 159 interpolation, temporal information aggregation and HR reconstruction. VideoINR(Chen et al., 2022) utilize the continuous video representation to achieve STVSR at arbitrary spatial resolution and 160 frame rate. Although these methods obtain smooth and high-resolution output videos, but they fail in 161 generating realistic texture details.



Figure 2: The overall framework. It consists of space-time data augmentation for constructing training data, and the associated *video-aware conditioning* for realizing diversified conditions across frames, as well as ST-Controller for multi-frame condition injection based on generative video prior.

PRELIMINARIES: VIDEO DIFFUSION MODELS

Our method is built on a pretrained video diffusion model (Zhang et al., 2023b), which is developed based on Stable Diffusion 2.1. Given an video $\mathbf{x} \in \mathbb{R}^{F \times H \times W \times 3}$, the encoder \mathcal{E} first encodes it into latent representation $\mathbf{z} = \mathcal{E}(\mathbf{x})$ frame-by-frame, where $\mathbf{z} \in \mathbb{R}^{F \times H' \times W' \times C}$. Then, the forward diffusion and reverse denoising are conducted in the latent space. In the forward process, the noise is gradually added to the latent vector \mathbf{z} in total T steps. And for each time-step t, the diffusion process is formulated as follows:

$$\mathbf{z}_t = \alpha_t \mathbf{z} + \sigma_t \boldsymbol{\epsilon},\tag{1}$$

where $\epsilon \in \mathcal{N}(0, \mathbf{I})$, and α_t, σ_t specify the noise schedule in which the corresponding log signal-to-noise-ratio $(\log[\alpha_t^2/\sigma_t^2])$ decreases monotonically with t. And at time-step $T, q(\mathbf{z}_T) = \mathcal{N}(\mathbf{0}, \mathbf{I})$. As for backward pass, a diffusion model is used for iteratively denoising under the guidance of the text prompt c_{text} . By adopting v-prediction parameterization (Salimans & Ho, 2022), the U-Net denoiser f_{θ} learns to make predictions of $\mathbf{v}_t \equiv \alpha_t \boldsymbol{\epsilon} - \sigma_t \mathbf{z}$. The optimization objective is simply formulated as:

$$\mathcal{L}_{LDM} = \mathbb{E}_{\mathbf{z}, c_{text}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), t} \Big[\|\mathbf{v} - f_{\theta}(\mathbf{z}_t, t, c_{text})\|_2^2 \Big].$$
(2)

At the end, the generated videos are obtained through the VAE decoder: $\hat{\mathbf{x}} = \mathcal{D}(\mathbf{z})$.

METHODOLOGY

In this section, we introduce the main components of our method. The overall framework is illustrated in Figure 2. First, we present our architecture design in Section 4.1. Then we elaborate on the proposed space-time data augmentation in Section 4.2. In Section 4.3, we give a detailed description on our designed video-aware conditioning.

4.1 ARCHITECTURE DESIGN

The architecture is designed based on a pretrained video diffusion model. This video diffusion model is able to generate temporal-coherent content and high-quality texture details through iterative denoising. To upsample and refine a low-frame-rate and low-resolution videos in both spatial and temporal dimensions, the visual information should be incorporated into the video diffusion model carefully in order to obtain high-quality results with good fidelity to the input videos. High-quality generated videos stem from powerful generative models, while fidelity requires the algorithm to preserve the visual information of the input. Inspired by (Zhang et al., 2023a), we keep the pretrained video diffusion model untouched for preserving generative capability, and create a Space-Time Controller (ST-Controller) to obtain effective multi-frame condition injection for generative video enhancement. The architecture is illustrated in Figure 3.

The pretrained video diffusion model follows the design of stacking a sequence of interleaved spatial and temporal layers within the 3D-UNet (Blattmann et al., 2023b) architecture (gray blocks in Figure 3). Specifically, each spatial convolution layer (or attention layer) is followed by a temporal convolution layer (or attention layer). The spatial layers are the same as those in Stable Diffusion 2.1, including ResBlocks (He et al., 2016), self-attention (Vaswani et al., 2017) layers, and cross-attention layers. The temporal convolution and attention layers are incorporated with their output layers initialized to zero and finetuned with video datasets. Specifically, the temporal convolution is one-dimensional convolution layer with a kernel size of 3, and the temporal attention is one-dimensional attention layer (Wang et al., 2023a). In this 3D-UNet, the video features that aligned by temporal layers in encoder will be skipped to the decoder, in which concatenation operation will be performed to combine skipped features with decoder features.

224 To build our proposed ST-Controller, we make 225 a copy (both the architectures and weights) of 226 the multi-frame encoder and middle block in 3D-UNet (orange blocks in Fig. 3) as the train-227 able condition network. This condition network 228 takes low-frame-rate and low-resolution condi-229 tion latents as well as full frames of noisy latents 230 as inputs. Specifically, the condition latents and 231 the associated augmentation parameters are in-232 corporated into the condition network through 233 our proposed video-aware conditioning. The 234 output multi-scale temporal-coherent video fea-235 tures will be injected into the original 3D-UNet 236 through newly added zero convolutions (yellow 237 blocks in Fig. 3). The output features of the mid-238 dle block in condition network will be added





back to the features of the middle block in 3D-UNet. While for output features of encoder blocks in
condition network, their features will be added to the skipped video features in 3D-UNet, which are
also produced by encoder blocks. The copied condition network, video-aware conditioning, and the
newly added zero convolutions are trained simultaneously.

243244 4.2 Space-Time Data Augmentation

In this section, we discuss about how to achieve unified space-time super-resolution with arbitrary up-sampling space and time scales, as well as refinement with varying degrees. To this end, we propose a novel data augmentation strategy for both space and time axes. Details are discussed below.

Time axis. Given a sequence of high-frame-rate and high-resolution video frames $I^{1:f} = [I^1, I^2, ..., I^f]$ with frame length f, we use a sliding window across time axis to select frames. The frame sliding window size m is randomly sampled from a predefined set, ranging from 1 to 8. This corresponds to time scales from $1 \times$ to $8 \times$. Note that $1 \times$ time scale requires no frame interpolation, thus the multi-task problem downgrades to video super-resolution. After the frame skipping, we obtain a sequence of key frames $I^{1:m:f} = [I^1, I^{1+m}, I^{1+2 \times m}, ..., I^f]$.

Space axis. Then, we perform spatial downsampling for these obtained key frames. Specifically, 255 the downscaling factor s is randomly sampled from [1, 8], which represents $1 \times \sim 8 \times$ space super-256 resolution. When s = 1, there is no need to perform spatial super-resolution. All frames in one 257 sequence are downsampled with the same downscaling factor s. Thus, we arrive at low-frame-rate 258 and low-resolution video frames: $I_{\downarrow s}^{1:m:f}$. In practice, we should upsample them back to the original 259 spatial sizes by bilinear interpolation before being passed to the networks, so we obtain $I_{\downarrow s,\uparrow s}^{1:m:f}$. Note 260 that each space or time scale corresponds to different difficulty level, and thus the sampling is not 261 uniform. Particularly, we set sampling probabilities of scales $4 \times$ and $8 \times$ based on a ratio of 1:2, 262 which is determined by their associated scale values. 263

Then, we use the encoder part of a pretrained variational autoencoder (VAE) \mathcal{E} to project the input sequence to the latent space frame-wisely:

$$z_s^{1:m:f} = [\mathcal{E}(I^1_{\downarrow s,\uparrow s}), \mathcal{E}(I^{1+m}_{\downarrow s,\uparrow s}), \mathcal{E}(I^{1+2\times m}_{\downarrow s,\uparrow s}), ..., \mathcal{E}(I^f_{\downarrow s,\uparrow s})].$$
(3)

267 268

266

269 *Noise augmentation in latent space*. At this stage, we conduct noise augmentation to noise the latent condition information in varying degrees in order to achieve controllable refinement. This noise

270 augmentation process is the same as the diffusion process equation 1 used in the video diffusion 271 model. Specifically, the condition latent sequence is corrupted by: 272

$$z_{s\,t'}^{1:m:f} = \alpha_{t'} z_s^{1:m:f} + \sigma_{t'} \epsilon^{1:m:f}, \tag{4}$$

274 where $\alpha_{t'}, \sigma_{t'}$ determine the signal-to-noise-ratio at time-step t', and $t' \in \{1, ..., T'\}$. Note that the 275 pretrained video diffusion model adopts 1,000 steps (T = 1000 in equation 1). While the noise 276 augmentation only needs to corrupt the low-level information, T' is set to 300 empirically. For more intuitive denotation, we use σ instead of t'. Finally, we arrive at $z_{s,\sigma}^{1:m:f} = \mathcal{E}(I_{\downarrow s,\uparrow s})_{\sigma}^{1:m:f}$. 277 278

The whole process of space-time data augmentation is summarized as follows:

$$I^{1:f} \to I^{1:m:f} \to I^{1:m:f}_{\downarrow s} \to I^{1:m:f}_{\downarrow s,\uparrow s} \to \mathcal{E}(I_{\downarrow s,\uparrow s})^{1:m:f} \to \mathcal{E}(I_{\downarrow s,\uparrow s})^{0}_{\sigma}.$$
(5)

4.3 VIDEO-AWARE CONDITIONING

273

279 280 281

283

315

316

284 Besides data augmentation, the corresponding conditioning mechanism should also be designed in order to influence the model training and avoid averaging 287 performance for different space or time scales and noise augmentation. In practice, the condition latent 288 sequence $z_{s,\sigma}^{1:m:f}$, the corresponding downscaling fac-289 tor s, and augmented noises σ are all considered as 290 for conditioning. Please refer to Figure. 4 for more 291 intuitive demonstration. 292

Given the synthesized condition latent sequence 293 $z_{s,\sigma}^{1:m:f}$, we use one convolution with zeroinitialization -Convzero for connecting it to the con-295 dition network. Specifically, we have: 296

297
$$f_{out}^{1:f} = \text{Conv}(z_t^{1:f}),$$
 (6)
298 $f_{out}^{1:m:f} = \text{Conv}(z_t^{1:m:f}) + \text{Conv}_{zero}(z_{s,\sigma}^{1:m:f}),$ (7)

300 where Conv is the first convolution in the condition network, $z_t^{1:f}$ and $z_t^{1:m:f}$ denote the full frames and 301 key frames of noisy latents at timestep t, respectively. 302 Note that Conv and Convzero share the same hyper-303 parameter configuration (i.e., kernel size, padding, 304 et.al.), As it is shown, only key-frame features in 305 condition network will be added with the condition 306 features, while others remain unchanged. This strat-307 egy enables progressive condition injection as the 308 weights of Convzero grows from zero starting point. 309



Figure 4: Video-aware conditioning. For frame that has condition image as input (key frame), we add the condition latent to the noisy latent after one convolution layer. Besides, the embeddings of noise level σ and downscaling factor s are added to the existing t embedding, which will be broadcast to all spatial layers.

For conditioning regarding downscaling factor s and noise augmentation σ , we incorporate them 310 to the existing time embedding in the condition network. Specifically, for timestep t, sinusoidal 311 encoding (Ho et al., 2020; Rombach et al., 2022; Vaswani et al., 2017) is used to provide the model 312 with a positional encoding for time. Then, one MLP (two linear layers with a SiLU (Elfwing et al., 313 2018) activation layer in between) is applied. Specifically, we have: 314

$$t_{emb} = \text{MLP}_t(\text{Sinusoidal}(t)), \quad t_{emb}^{1:f} = \text{Repeat}(t_{emb}, f), \tag{8}$$

where $t_{emb}^{1:f}$ is obtained by Repeat t_{emb} by f times in the frame axis. This time embedding sequence 317 will be broadcast to all ResBlocks in the condition network for timestep injection. 318

319 Also, we elucidate the conditioning for noise augmentation. As mentioned in equation 4, noise 320 augmentation shares the same way as diffusion process, but with much smaller maximum timestep 321 (*i.e.*, T' = 300). Thus, we reuse the encoding and mapping for timestep t in diffusion process. After this, we add a linear layer with zero initialization (denoted as Linearzero). To conclude, we have: 322 323 σ

$$m_{emb} = \text{Linear}_{\text{zero},\sigma}(\text{MLP}_t(\text{Sinusoidal}(\sigma)).$$
(9)

324 To achieve video-aware conditioning, we add σ_{emb} only to the key frames. We repeat σ_{emb} by k 325 times to obtain $\sigma_{emb}^{1:k}$, where k is the number of key frames. The video-aware controlling is presented 326 as follows:

$$t_{emb}^{1:m:f} = t_{emb}^{1:m:f} + \sigma_{emb}^{1:k},$$
(10)

where the addition operation is performed frame-wisely.

330 Regarding downscaling factor s, the corresponding encoding, mapping and controlling are similar as 331 above. In particular, we newly introduce one MLP_{zero}, in which the output layer is zero-initialized. 332 The video-aware conditioning is performed as:

$$s_{emb} = MLP_{zero,s}(Sinusoidal(s)), \quad s_{emb}^{1:k} = Repeat(s_{emb}, k),$$
 (11)

334 335 336

337

338

339

340

353

354

355 356 357

358

333

327 328

$$t_{emb}^{1:m:f} = t_{emb}^{1:m:f} + s_{emb}^{1:k}.$$
(11)

$$t_{emb}^{...m.j} = t_{emb}^{1..m.j} + s_{emb}^{1.k}.$$
 (12)

With our proposed space-time data augmentation and video-aware conditioning, VEnhancer can be well-trained in an end-to-end manner, and yields great performance in handling diversified conditions for generative enhancement. Here we provide a demonstration in Figure 5. For $4 \times$ video superresolution, we modify the input downscaling factor s to produce different results. It is shown that more texture details are generated as s grows (from smooth to sharp). This indicates that s can determine how many details are generated through our proposed video-aware conditioning.



Figure 5: The effectiveness of video-aware conditioning. For video super-resolution $(4\times)$, we modify the input downscaling factor from s = 1 to s = 8, and more texture details are generated. Zoom in for best view.

5 **EXPERIMENTS**

Datasets. We collect around 350k high-quality and high-resolution video clips from Panda-70M 360 (Chen et al., 2024b) dataset and the Internet to constitute our training set. We train VEnhancer on 361 resolution 1280×720 with center cropping, and the target FPS is fixed to 24 by frame skipping. 362 Regarding test dataset, we collect generated videos from state-of-the-art text-to-video methods. Practically, we select videos with large motions and diverse contents. This test dataset is denoted as AIGC2023, which is used to evaluate VEnhancer and baselines for video super-resolution and 364 space-time super-resolution tasks. For evaluation on VBench, all generated videos based on the provided prompt suite are considered, resulting in around 5k videos. 366

367 Implementation Details. The batch size is set to 256. AdamW (Loshchilov & Hutter, 2017) is 368 used as the optimizer, and the learning rate is set to 10^{-5} . During training, we dropout the text 369 prompt with a probability of 10%. The training process lasts about four days with 16 NVIDIA A100 GPUs. During inference, we use DPM-Solver (Lu et al., 2022) and perform 15 sampling steps with 370 classifier-free guidance (cfg) (Ho & Salimans, 2022). 371

372 Metrics. Regarding evaluation for video super-resolution and space-time super-resolution on 373 AIGC2023 test dataset, we use both image quality assessment (IQA) and video quality assess-374 ment (VQA) metrics. Specifically, MUSIQ (Ke et al., 2021) and DOVER (Wu et al., 2023a) are 375 adopted. Moreover, we refer to video generation benchmark, VBench (Huang et al., 2023), for more comprehensive evaluation. Specifically, we choose **Dynamic Degree** (*i.e.*, whether it contains large 376 motions), Motion smoothness (*i.e.*, how smooth the video is), and Aesthetic Quality for evaluation. 377 Regarding evaluation for video generation, we consider all 16 evaluation dimensions from VBench.

3785.1 COMPARISON WITH VIDEO SUPER-RESOLUTION METHODS379

For video super-resolution, VEnhancer is compared with the state-of-the-art real-world video super-resolution method, RealBasicVSR (Chan et al., 2022b), and the state-of-the-art generative video
 super-resolution method, LaVie-SR (Wang et al., 2023b) (super-resolution). For more comprehensive comparison, we also include I2VGen-XL' (Zhang et al., 2023b) refinement model as our baseline.

Table 1: Quantitative comparison for video super-resolution $(4\times)$ on AIGC2023 test dataset. Red and blue indicate the best and second best performance. The top 3 results are marked as gray.

	DOVER↑	MUSIQ↑	Aesthetic Ouality	Dynamic Degree	Motion Smoothness
LaVie-SR Wang et al. (2023b)	0.8427	55.8428	0.6692	0.525	0.9710
I2VGen-XL(refiner) Zhang et al. (2023b)	0.5603	25.5988	0.6439	0.475	0.9835
RealBasicVSR Chan et al. (2021)	0.8252	50.5978	0.6622	0.550	0.9729
Ours	0.8586	59.4474	0.6671	0.550	0.9781
RealBasicVSR					
LaVie-SR					
I2VGen-XL (refiner)		20er			
Ours					

Figure 6: Visual comparison for video super-resolution $(4\times)$ on AIGC2023 test dataset. Input resolution: 512×312 ; output resolution: 2048×1280 . Prompt: *Iron Man flying in the sky*.

As shown in Table 1, VEnhancer outperforms both generative video super-resolution method (LaVie-SR), real-world video super-resolution method (RealBasicVSR), and generative video refinment method (I2VGen-XL, refiner) in most metrics, suggesting its outstanding enhancement ability for videos. Note that LaVie-SR surpasses RealBasicVSR in image/video quality (MUSIQ, DOVER, and Aesthetic Quality), as diffusion-based methods are better at generating sharp details. But LaVie-SR achieves the worst Motion Smoothness, indicating its insufficient capability in balancing video smoothness and video quality. I2VGen-XL's refiner could obtain the highest score in Motion **Smoothness**, but sacrifices the magnitude of motion significantly (worst **Dynamic Degree**). Moreover, it achieves very unsatisfactory results in metrics evaluating image and video quality. Because it uses bilinear interpolation for $\times 4$ upsampling, which produces very blurry results. Nevertheless, VEnhancer achieves overall best results with good balance among image/video quality, motion smoothness, and motion magnitude.

432 The visual comparison is presented in Figure. 6. The prompt is "Iron man flying in the sky". The input 433 video is already consistent with the prompt, but lacks details on the iron man suit. RealBasicVSR 434 could remove some noises or artifacts of the generated videos as it incorporates complex degradation 435 for model training. However, it fails in generating realistic details but produces over-smoothed results, 436 since its generative ability is limited. On the other hand, the results of LaVie-SR contains more artifacts than input. Without successfully removing artifacts, the generative super-resolution model 437 will enlarge the existing defects. The refiner of I2VGen-XL could achieve successful refinement but 438 produces over-smoothed results. In contrast, VEnhancer could first remove unpleasing artifacts and 439 refine the distorted content (e.g., head region), and then generate faithfuls details (e.g., helmet and 440 armor) that are consistent with the text prompt. 441

442 443

444

445

446

447

448

449

5.2 COMPARISON WITH SPACE-TIME SUPER-RESOLUTION METHODS

For space-time super-resolution task, we compare two state-of-the-art space-time super-resolution methods: VideoINR (Chen et al., 2022) and Zooming-Slow-Mo (Xiang et al., 2020) (Zoom for short). We also consider LaVie's cascaded T-SR and S-SR DM-based pipeline: LaVie-FI (frame interpolation) + LaVie-SR (super-resolution) for more thorough comparison.

Table 2: Quantitative comparison for space-time super-resolution $(4\times)$ on AIGC2023 test dataset. **Red** and blue indicate the best and second best performance. The top 3 results are marked as gray.

	DOVER↑	MUSIQ↑	Aesthetic Quality	Dynamic Degree	Motion Smoothness
LaVie-FI + LaVie-SR Wang et al. (2023b)	0.8159	53.2128	0.6566	0.60	0.9857
VideoINR Chen et al. (2022)	0.7608	34.1060	0.6624	0.60	0.9933
Zooming Slow-Mo Xiang et al. (2020)	0.7328	33.8470	0.6624	0.55	0.9908
Ours	0.8609	51.2940	0.6710	0.60	0.9937

456 As shown in Table 2, we observe that our method surpasses all baselines in DOVER and Aesthetic 457 Quality, showing its superior capability in generating sharp and realistic video content. Besides, it obtains highest scores in Motion Smoothness and Dynamic Degree, indicating VEnhancer's 458 excellent ability in synthesizing stable temporal details. The cascaded T-SR and S-SR approach 459 (LaVie-FI + LaVie-SR) obtains good scores in DOVER and MUSIQ, demonstrating DM-based 460 methods' advantage in generation. However, its performance in temporal aspect is unsatisfactory 461 due to its inferior capability in temporal refinement. We notice that state-of-the-art space-time 462 super-resolution methods (VideoINR and Zooming Slow-Mo) behave well in Motion Smoothness. 463 As both of them are optimized with reconstruction loss, the produced results are very smooth across 464 frames. At a cost, they perform poorly in metrics regarding quality, such as DOVER and MUSIQ. 465

The visual comparison is illustrated in Figure. 7. The first and third columns present the low-resolution 466 key frames. Note that the input frames are not consistent especially in the region of guitar strings. 467 The cascaded T-SR and S-SR approach, LaVie-FI + LaVie-SR, can produce very sharp results for all 468 frames (key and predicted ones). However, it generates messy contents which are not semantically 469 aligned with prompt. Moreover, the generated details are changing across time, indicating severe 470 flickering. For reconstruction-based methods (VideoINR and Zoom), the produced results are similar: 471 lacking details and failing in improving the consistency of the original input frames. On the contrary, 472 VEnhancer is not only able to achieve unified space-time super-resolution, but can also improve the 473 temporal consistency of the generated videos by refinement (*i.e.*, guitar strings and raccoon hands).

474 475

476

5.3 EVALUATION ON IMPROVING VIDEO GENERATION

Here we evaluate VEnhancer's ability in improving state-of-the-art T2V methods. The baselines
includes open-source T2V methods–VideoCrafter-2 (Chen et al., 2024a) (VC-2 for short), Lavie
(Wang et al., 2023b), Open-Sora, CogVideoX (Yang et al., 2024), and professional video generation
products–Pika and Gen-2. In particular, we enhance T2V results of CogVideoX-5B and VC-2.

The quantitative results are organized in Table 3. Before enhancement, CogVideoX-5B and VC-2 achieve the best and second best in *Semantic* compared with other baselines, demonstrating their superiority in generating video contents that are highly consistent to the VBench's prompt suite. Regarding *Quality*, they lag behind Pika. But with VEnhancer, VC-2 and CogVideoX-5B are able to achieve the highest and second highest scores in *Quality*. Besides, their scores regarding *Semantic* improves a lot, especially VC-2' (3.31% increase). This indicates that VEnhancer can improve the Table 3: VBench Evaluation Results. This table compares the performance of open-source T2V methods and professional video generation products regarding two aspects (*Quality* and *Semantic*). A higher score indicates better performance. Red and blue indicate the best and second best performance. The top 3 results are marked as gray.

	LaVie Wang et al. (2023b)	Open-Sora	Pika	Gen-2	CogVideoX-5B Yang et al. (2024)	VC-2 Chen et al. (2024a)	CogVideoX-5B +Venhancer	VC-2 +VEnhancer
Quality	78.78%	80.71%	82.92%	82.47%	82.75%	82.20%	82.99%	83.28 %
Semantic	70.31%	73.30%	71.77%	73.03%	77.04%	73.42%	77.52%	76.73%
Overall	77.08%	79.23%	80.69%	80.58%	81.61%	80.44%	81.90%	81.97 %

semantic content and video quality at the same time, showing a powerful generative enhancement ability. More importantly, the advantage of adopting a two-stage pipeline is observed: the first T2V model focuses on generating semantic content and motions with good fidelity to the prompts, while the following enhancement model can improve the semantic in low-level and image quality, as well as temporal consistency. For visual results, **please see the video demonstration in supp.**

495

496

497

5.4 ABLATION STUDIES

502 The Effectiveness of Noise Augmentation. During training, the noise level regarding noise augmentation is randomly sampled within a predefined range. While during inference, one can change the 504 noise level to achieve refinement with different strengths. In general, higher noise corresponds to 505 stronger refinement and regeneration. We present the visual comparison among different noise levels 506 in Figure 8. The first frame of one AI-generated video is presented in the left. It is of low-resolution 507 and lacks details. Also, the original video has very obvious flickering. If we set $\sigma = 0$, VEnhancer 508 will generate unpleasing noises in the background. As there is domain mismatch between the training 509 data and testing data, the enhancement fails in handling unseen and challenging scenarios. Fortunately, we can mitigate this by adding noise in the condition latents for corrupting the noisy and unknown 510 low-level details. As we increase the noise level, the artifacts are gradually vanishing. When $\sigma = 250$, 511 the result is noise-clean, and has abundant semantic details. 512

Arbitrary Up-sampling Scales for Spatial Super-Resolution. Here we show that VEnhancer is able to up-sample videos with arbitrary scales. From Figure 9, we observe that VEnhancer could produce satisfactory results on different scales $(2.5 \times, 3 \times, 3.5 \times, 4 \times, \text{ and } 4.5 \times)$, suggesting its flexibility and generalization in adapting to different tasks. In particular, given one frame of the generated video (312×512) , VEnhancer could improve the generated details when the up-sampling scale grows up. When $s = 2.5 \sim 3.5$, the panda's hand is less realistic. But it becomes better when s = 4 or s = 4.5. It is also noticed that the panda' fur is becoming more realistic as s grows.

520 Arbitrary Up-sampling Scales for Temporal Super-Resolution. In this part, we show VEnhancer is able to achieve arbitrary up-sampling in time axis. Given two low-resolution key frames, we aim to 521 up-sample them to high-resolution ones, and also interpolate several frames (ranging from 2 to 4) 522 between them. As shown in Figure 10, the results are consistent across frames, showing not flicking 523 or distortions. Besides, the spatial quality has also been significantly improved. As shown in the 524 last row, $5 \times$ frame interpolation yields smooth frames with generated contents: the shadow in the 525 right leg is changing, showing a very natural transition. This indicates that diffusion-based frame 526 interpolation has great capability in both motion and content generation. 527

528 529

530

6 CONCLUSION AND LIMITATION

531 In this work, we propose a generative space-time enhancement method – VEnhancer for video 532 generation. It can achieve spatial super-resolution, temporal super-resolution and video refinement 533 in one video diffusion model. We base on a pretrained generative video prior and build a Space-534 Time Controller (ST-Controller) for effective condition injection. Space-time data augmentation and video-aware conditioning are proposed to train ST-Controller in an end-to-end manner. Extensive 536 experiments have demonstrated our superiority over state-of-the-art video super-resolution and spacetime super-resolution methods in enhancing AI-generated videos. However, our work has several limitations. First, as it is based on diffusion models, the inference takes more time than one-step 538 methods. Second, it may face challenges in handling AI-generated long videos, since the long-term (over 10s) consistency has not been addressed in this work.

540 REFERENCES 541

547

566

567

568

569

- Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and 542 image encoder for end-to-end retrieval. In IEEE International Conference on Computer Vision, 543 2021. 544
- Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik 546 Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. arXiv preprint arXiv:2311.15127, 2023a. 548
- Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and 549 Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. 550 In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 551 22563–22575, 2023b. 552
- 553 Jiezhang Cao, Yawei Li, Kai Zhang, and Luc Van Gool. Video super-resolution transformer. arXiv 554 preprint arXiv:2106.06847, 2021. 555
- Kelvin CK Chan, Xintao Wang, Ke Yu, Chao Dong, and Chen Change Loy. Basicvsr: The search 556 for essential components in video super-resolution and beyond. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 4947-4956, 2021. 558
- 559 Kelvin CK Chan, Shangchen Zhou, Xiangyu Xu, and Chen Change Loy. Basicvsr++: Improving 560 video super-resolution with enhanced propagation and alignment. In Proceedings of the IEEE/CVF 561 conference on computer vision and pattern recognition, pp. 5972–5981, 2022a. 562
- Kelvin CK Chan, Shangchen Zhou, Xiangyu Xu, and Chen Change Loy. Investigating tradeoffs 563 in real-world video super-resolution. In Proceedings of the IEEE/CVF Conference on Computer 564 Vision and Pattern Recognition, pp. 5962–5971, 2022b. 565
 - Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, et al. Videocrafter1: Open diffusion models for high-quality video generation. arXiv preprint arXiv:2310.19512, 2023a.
- Haoxin Chen, Yong Zhang, Xiaodong Cun, Menghan Xia, Xintao Wang, Chao Weng, and Ying 570 Shan. Videocrafter2: Overcoming data limitations for high-quality video diffusion models. arXiv 571 preprint arXiv:2401.09047, 2024a. 572
- 573 Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, 574 James Kwok, Ping Luo, Huchuan Lu, et al. Pixart- α : Fast training of diffusion transformer for 575 photorealistic text-to-image synthesis. arXiv preprint arXiv:2310.00426, 2023b. 576
- Shoufa Chen, Mengmeng Xu, Jiawei Ren, Yuren Cong, Sen He, Yanping Xie, Animesh Sinha, Ping 577 Luo, Tao Xiang, and Juan-Manuel Perez-Rua. Gentron: Delving deep into diffusion transformers 578 for image and video generation. arXiv preprint arXiv:2312.04557, 2023c. 579
- 580 Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, Ekaterina Deyneka, Hsiang-wei Chao, 581 Byung Eun Jeon, Yuwei Fang, Hsin-Ying Lee, Jian Ren, Ming-Hsuan Yang, and Sergey Tulyakov. 582 Panda-70m: Captioning 70m videos with multiple cross-modality teachers. arXiv preprint 583 arXiv:2402.19479, 2024b. 584
- Zeyuan Chen, Yinbo Chen, Jingwen Liu, Xingqian Xu, Vidit Goel, Zhangyang Wang, Humphrey 585 Shi, and Xiaolong Wang. Videoinr: Learning video implicit neural representation for continuous 586 space-time super-resolution. In Proceedings of the IEEE/CVF Conference on Computer Vision and 587 Pattern Recognition, pp. 2047–2057, 2022. 588
- 589 Stefan Elfwing, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network 590 function approximation in reinforcement learning. Neural networks, 107:3-11, 2018. 591

Peng Gao, Le Zhuo, Ziyi Lin, Chris Liu, Junsong Chen, Ruoyi Du, Enze Xie, Xu Luo, Longtian Qiu, 592 Yuhang Zhang, et al. Lumina-t2x: Transforming text into any modality, resolution, and duration via flow-based large diffusion transformers. arXiv preprint arXiv:2405.05945, 2024.

630

- Yuwei Guo, Ceyuan Yang, Anyi Rao, Yaohui Wang, Yu Qiao, Dahua Lin, and Bo Dai. Animatediff:
 Animate your personalized text-to-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023.
- Agrim Gupta, Lijun Yu, Kihyuk Sohn, Xiuye Gu, Meera Hahn, Li Fei-Fei, Irfan Essa, Lu Jiang, and José Lezama. Photorealistic video generation with diffusion models. *arXiv preprint arXiv:2312.06662*, 2023.
- Muhammad Haris, Greg Shakhnarovich, and Norimichi Ukita. Space-time-aware multi-resolution
 video enhancement. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2859–2868, 2020.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Roberto Henschel, Levon Khachatryan, Daniil Hayrapetyan, Hayk Poghosyan, Vahram Tadevosyan,
 Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Streamingt2v: Consistent, dynamic,
 and extendable long video generation from text. *arXiv preprint arXiv:2403.14773*, 2024.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. arXiv preprint arXiv:2207.12598, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in neural information processing systems, 33:6840–6851, 2020.
- Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P
 Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition
 video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.
- Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing
 Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. *arXiv preprint arXiv:2311.17982*, 2023.
- Takashi Isobe, Xu Jia, Shuhang Gu, Songjiang Li, Shengjin Wang, and Qi Tian. Video super resolution with recurrent structure-detail network. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XII 16*, pp. 645–660. Springer, 2020a.
- Takashi Isobe, Songjiang Li, Xu Jia, Shanxin Yuan, Gregory Slabaugh, Chunjing Xu, Ya-Li Li, Shengjin Wang, and Qi Tian. Video super-resolution with temporal group attention. In *Proceedings* of the IEEE/CVF conference on computer vision and pattern recognition, pp. 8008–8017, 2020b.
- Takashi Isobe, Fang Zhu, Xu Jia, and Shengjin Wang. Revisiting temporal modeling for video super-resolution. *arXiv preprint arXiv:2008.05765*, 2020c.
- Junjie Ke, Qifei Wang, Yilin Wang, Peyman Milanfar, and Feng Yang. Musiq: Multi-scale image
 quality transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 5148–5157, 2021.
- Soo Ye Kim, Jihyong Oh, and Munchurl Kim. Fisr: Deep joint frame interpolation and superresolution with a multi-scale temporal loss. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 11278–11286, 2020.
- Jingyun Liang, Yuchen Fan, Xiaoyu Xiang, Rakesh Ranjan, Eddy Ilg, Simon Green, Jiezhang Cao,
 Kai Zhang, Radu Timofte, and Luc V Gool. Recurrent video restoration transformer with guided
 deformable attention. *Advances in Neural Information Processing Systems*, 35:378–393, 2022.
- Jingyun Liang, Jiezhang Cao, Yuchen Fan, Kai Zhang, Rakesh Ranjan, Yawei Li, Radu Timofte, and Luc Van Gool. Vrt: A video restoration transformer. *IEEE Transactions on Image Processing*, 2024.
- ⁶⁴⁷ Xinqi Lin, Jingwen He, Ziyan Chen, Zhaoyang Lyu, Bo Dai, Fanghua Yu, Wanli Ouyang, Yu Qiao, and Chao Dong. Diffbir: Towards blind image restoration with generative diffusion prior, 2024.

667

677

684

- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017.
- Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *arXiv preprint arXiv:2206.00927*, 2022.
- Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon.
 Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021.
- Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe
 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image
 synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer- ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modelscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023a.
- Kiang Wang, Hangjie Yuan, Shiwei Zhang, Dayou Chen, Jiuniu Wang, Yingya Zhang, Yujun Shen,
 Deli Zhao, and Jingren Zhou. Videocomposer: Compositional video synthesis with motion
 controllability. *Advances in Neural Information Processing Systems*, 36, 2024.
- Xintao Wang, Kelvin CK Chan, Ke Yu, Chao Dong, and Chen Change Loy. Edvr: Video restoration with enhanced deformable convolutional networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pp. 0–0, 2019.
- Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yinan
 He, Jiashuo Yu, Peiqing Yang, et al. Lavie: High-quality video generation with cascaded latent
 diffusion models. *arXiv preprint arXiv:2309.15103*, 2023b.
- Haoning Wu, Erli Zhang, Liang Liao, Chaofeng Chen, Jingwen Hou Hou, Annan Wang, Wenxiu Sun Sun, Qiong Yan, and Weisi Lin. Exploring video quality assessment on user generated contents from aesthetic and technical perspectives. In *International Conference on Computer Vision (ICCV)*, 2023a.
- Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu,
 Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion
 models for text-to-video generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7623–7633, 2023b.
- Kiaoyu Xiang, Yapeng Tian, Yulun Zhang, Yun Fu, Jan P Allebach, and Chenliang Xu. Zooming slow-mo: Fast and accurate one-stage space-time video super-resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 3370–3379, 2020.
- Liangbin Xie, Xintao Wang, Shuwei Shi, Jinjin Gu, Chao Dong, and Ying Shan. Mitigating artifacts
 in real-world video super-resolution models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 2956–2964, 2023.
- 701 Tianfan Xue, Baian Chen, Jiajun Wu, Donglai Wei, and William T Freeman. Video enhancement with task-oriented flow. *International Journal of Computer Vision*, 127:1106–1125, 2019.

- Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang,
 Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models
 with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024.
- Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3836–3847, 2023a.
 - Shiwei Zhang, Jiayu Wang, Yingya Zhang, Kang Zhao, Hangjie Yuan, Zhiwu Qin, Xiang Wang, Deli Zhao, and Jingren Zhou. I2vgen-xl: High-quality image-to-video synthesis via cascaded diffusion models. *arXiv preprint arXiv:2311.04145*, 2023b.
- Shangchen Zhou, Peiqing Yang, Jianyi Wang, Yihang Luo, and Chen Change Loy. Upscale-a-video: Temporal-consistent diffusion model for real-world video super-resolution. *arXiv preprint* arXiv:2312.06640, 2023.

A APPENDIX

A.1 COMPARISON WITH SPACE-TIME SUPER-RESOLUTION METHODS



Figure 7: Visual comparison for space-time super-resolution on AIGC2023 test dataset. Prompt: *A cute raccoon playing guitar in a boat on the ocean.* **Zoom in for best view.**

A.2 ABLATION STUDIES

Comparison with Space-Time Super-Resolution Methods



- Figure 8: Visual comparison of setting different noise levels in noise augmentation during testing.

Arbitrary Up-sampling Scales for Spatial Super-Resolution.



Figure 9: Visual results of different up-sampling scales $(2.5\times, 3\times, 3.5\times, 4\times, \text{ and } 4.5\times)$ for spatial super-resolution during testing.

Arbitrary Up-sampling Scales for Temporal Super-Resolution.



Figure 10: Visual results of different up-sampling scales $(3\times, 4\times, \text{ and } 5\times)$ for temporal superresolution during testing.