

I4VGEN: IMAGE AS FREE STEPPING STONE FOR TEXT-TO-VIDEO GENERATION

Anonymous authors

Paper under double-blind review

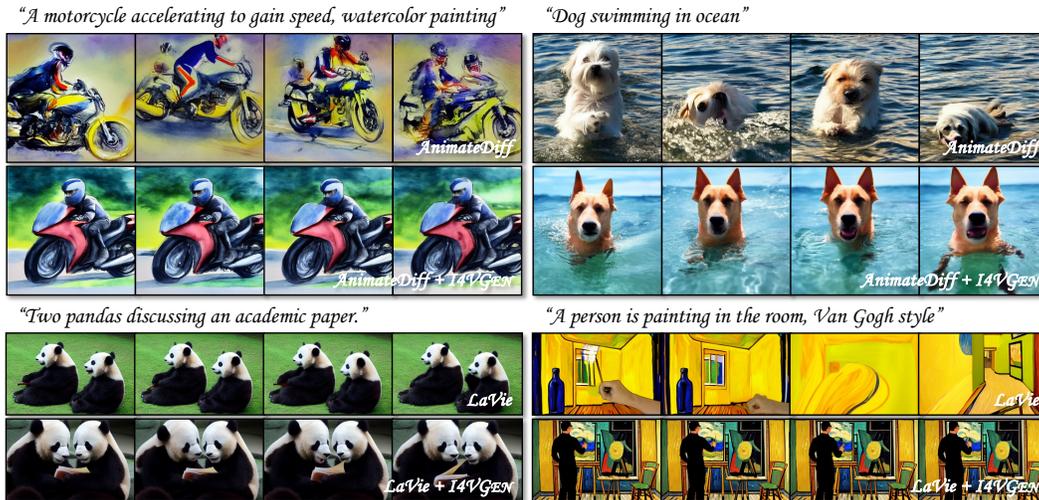


Figure 1: **Example results** synthesized by the proposed I4VGEN. I4VGEN is seamlessly integrated into existing pre-trained text-to-video diffusion models without additional training, significantly improving the temporal consistency (e.g., top-left and bottom-right), visual realism (e.g., top-right), and semantic fidelity (e.g., bottom-left) of the synthesized videos.

ABSTRACT

Text-to-video generation has trailed behind text-to-image generation in terms of quality and diversity, primarily due to the inherent complexities of spatio-temporal modeling and the limited availability of video-text datasets. Recent text-to-video diffusion models employ the image as an intermediate step, significantly enhancing overall performance but incurring high training costs. In this paper, we present I4VGEN, a novel video diffusion inference pipeline to leverage advanced image techniques to enhance pre-trained text-to-video diffusion models, which requires no additional training. Instead of the vanilla text-to-video inference pipeline, I4VGEN consists of two stages: anchor image synthesis and anchor image-augmented text-to-video synthesis. Correspondingly, a simple yet effective generation-selection strategy is employed to achieve visually-realistic and semantically-faithful anchor image, and an innovative noise-invariant video score distillation sampling (NI-VSDS) is developed to animate the image to a dynamic video by distilling motion knowledge from video diffusion models, followed by a video regeneration process to refine the video. Extensive experiments show that the proposed method produces videos with higher visual realism and textual fidelity. Furthermore, I4VGEN also supports being seamlessly integrated into existing image-to-video diffusion models, thereby improving overall video quality.

1 INTRODUCTION

Recent advances in large-scale text-to-image diffusion models (Esser et al., 2021; Balaji et al., 2022; Ramesh et al., 2022; Nichol et al., 2022; Saharia et al., 2022; Feng et al., 2023; Gu et al., 2023; Xue

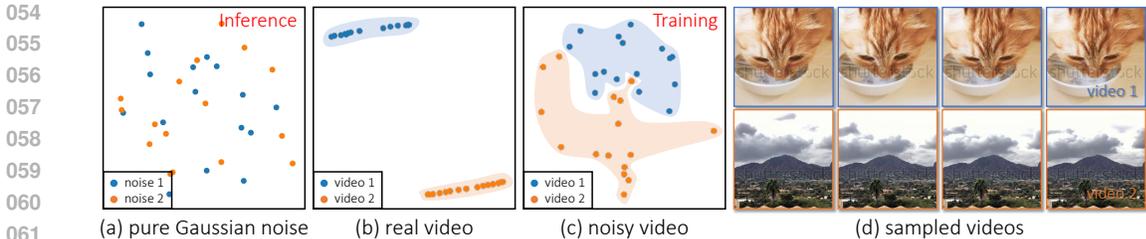


Figure 2: **Illustration of non-zero terminal signal-to-noise ratio.** We employ t-SNE to visualize the distributions of pure Gaussian noise, real video, and noisy video at the timestep T , where each data point represents an independently sampled noise point or video frame. The noise schedule of AnimateDiff (Guo et al., 2024b) is used, and all operations are performed in the latent space of the video autoencoder. (a) The distribution of pure Gaussian noise exhibits a disordered and diffuse nature; (b) real videos are temporally-correlated and different videos can be clearly distinguished from each other; (c) noisy videos preserve a certain degree of temporal correlation and maintain separability between different videos; (d) sampled videos for visualization.

et al., 2023) have demonstrated the capability to generate diverse and high-quality images from extensive web-scale image-text pair datasets. Efforts to extend these diffusion models to text-to-video synthesis (Ho et al., 2022a; Zhou et al., 2022; Chen et al., 2023a; Singer et al., 2023; Wang et al., 2023b; Blattmann et al., 2023a; Girdhar et al., 2023; Guo et al., 2024b; Bao et al., 2024) have involved leveraging video-text pairs and temporal modeling. However, text-to-video generation remains inferior to image counterpart in terms of both quality and diversity, primarily due to the complex nature of spatio-temporal modeling and the limited size of video-text datasets, which are often an order of magnitude smaller than image-text datasets.

This paper explores a novel video diffusion inference pipeline that leverages advanced image techniques to enhance pre-trained text-to-video diffusion models, focusing on the following two insights:

Image conditioning for text-to-video generation. Recent methods (Blattmann et al., 2023a; Zhang et al., 2023b; Girdhar et al., 2023; Chen et al., 2024a; Li et al., 2023; Hu et al., 2023) have adopted image-guided text-to-video generation, where an initial image generation step significantly enhances video output quality. This paradigm benefits from the strong capabilities of text-to-image models by using the generated images as detailed references for video synthesis. While effective, these approaches incur additional high training costs. This paper builds on this insight but innovates by designing a novel video diffusion inference pipeline to leverage image information, thereby enhancing text-to-video generation performance without additional training expense.

Zero terminal-SNR noise schedule. A prevalent issue in diffusion models is the non-zero terminal signal-to-noise ratio (SNR) (Guttenberg; Lin et al., 2024). The mismatch between the training phase, where residual signals persist in noisy videos at the terminal diffusion timestep T , and the inference phase, which uses pure Gaussian noise at the timestep T , creates a gap that degrades the model performance. As illustrated in Fig. 2, noisy videos exhibit temporal correlation that is distinctly different from the independent and identically distributed pure Gaussian noise. This paper is dedicated to reconfiguring the inference pipeline to circumvent this issue.

Motivated by these insights, we propose a novel video diffusion inference pipeline, called I4VGEN, which enhances pre-trained text-to-video diffusion models by incorporating image information into the inference process. This method requires no additional learnable parameters and training costs, and can be seamlessly integrated into existing text-to-video diffusion models, circumventing the non-zero terminal SNR issue and improving output quality.

Specifically, instead of the vanilla text-to-video inference pipeline, which fails to leverage image reference information, I4VGEN decomposes the inference process into two stages: anchor image synthesis and anchor image-augmented text-to-video synthesis. For the former, a simple yet effective generation-selection strategy is introduced, which involves synthesizing candidate images and selecting the most suitable one using a reward-based mechanism, thereby obtaining a visually-realistic anchor image that is closely aligned with the text prompt. For the latter, we develop an innovative noise-invariant video score distillation sampling (NI-VSDS) to animate the anchor image

to a dynamic video by extracting motion knowledge from text-to-video diffusion models, followed by a video regeneration process, *i.e.*, diffusion-denoising, to refine the video. This inference pipeline avoids the issue of non-zero terminal SNR.

Extensive quantitative and qualitative analyses demonstrate that I4VGEN can be effectively applied to various text-to-video diffusion models, significantly improving the temporal consistency, visual realism, and semantic fidelity of the synthesized videos (see Fig. 1). Moreover, our method can also be seamlessly integrated into existing image-to-video diffusion models, thereby enhancing the temporal consistency and visual quality of the generated videos (see Fig. 6).

The main novelties and contributions are as follows:

- We propose a novel video diffusion inference pipeline, called I4VGEN, which enhances pre-trained text-to-video diffusion models by incorporating image reference information into the inference process, without requiring additional training or learnable parameters.
- We employ a simple yet effective generation-selection strategy to achieve high-quality image, and design a novel noise-invariant video score distillation sampling for image animation.
- We comprehensively evaluate our approach with representative text-to-video diffusion models, and demonstrate I4VGEN significantly improves the quality of generated videos. Furthermore, I4VGEN can also be adapted to image-to-video diffusion models, leading to improved results.

2 PRELIMINARIES

Video diffusion models. Aligned with the framework of image diffusion models, Video diffusion models (VDMs) predominantly utilize the paradigm of latent diffusion models (LDMs). Unlike traditional methods that operate directly in the pixel space, VDMs function within the latent space defined by a video autoencoder. Specifically, a video encoder $\mathcal{E}(\cdot)$ learns the mapping from an input video $\mathbf{v} \in \mathcal{V}$, $\mathbf{v} = \{\mathbf{f}^1, \mathbf{f}^2, \dots, \mathbf{f}^F\}$ to a latent code $\mathbf{z} = \mathcal{E}(\mathbf{v}) = \{\mathbf{z}^1, \mathbf{z}^2, \dots, \mathbf{z}^f\}$. Subsequently, a video decoder $\mathcal{D}(\cdot)$ reconstructs the input video, aiming for $\mathcal{D}(\mathcal{E}(\mathbf{v})) \approx \mathbf{v}$. Typically, image autoencoder is used in a frame-by-frame processing manner instead of the video one, where $F = f$.

Upon training the autoencoder, a Denoising Diffusion Probabilistic Model (DDPM) (Ho et al., 2020) is employed within the latent space to generate a denoised version of an input latent \mathbf{z}_t at each timestep t . During denoising, the diffusion model can be conditioned on additional inputs, such as a text embedding $\mathbf{c} = f_{\text{CLIP}}(\mathbf{y})$ generated by a pre-trained CLIP text encoder (Radford et al., 2021), corresponding to the input text prompt \mathbf{y} . The DDPM model $\epsilon_\theta(\cdot)$, a 3D U-Net parametrized by θ , optimizes the following loss:

$$\mathcal{L} = \mathbb{E}_{\mathbf{z} \sim \mathcal{E}(\mathbf{v}), \mathbf{c} = f_{\text{CLIP}}(\mathbf{y}), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{1}), t} [\|\epsilon - \epsilon_\theta(\mathbf{z}_t, \mathbf{c}, t)\|_2^2], \quad (1)$$

During inference, a latent variable \mathbf{z}_T is sampled from the standard Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{1})$ and subjected to sequential denoising procedures of the DDPM to derive a refined latent \mathbf{z}_0 . This denoised latent \mathbf{z}_0 is then fed into the decoder to synthesize the corresponding video $\mathcal{D}(\mathbf{z}_0)$.

Score distillation sampling. Score distillation sampling (SDS) (Poole et al., 2023; Wang et al., 2023a) employs the priors of pre-trained text-to-image models to facilitate text-conditioned 3D generation. Specifically, given a pre-trained diffusion model $\epsilon_\theta(\cdot)$ and the conditioning embedding $\mathbf{c} = f_{\text{CLIP}}(\mathbf{y})$ corresponding to the text prompt \mathbf{y} , SDS optimizes a set of parameters ϕ of a differentiable parametric image generator $\mathcal{G}(\cdot)$ (*e.g.*, NeRF (Mildenhall et al., 2020)) using the gradient of the SDS loss \mathcal{L}_{SDS} :

$$\nabla_\phi \mathcal{L}_{\text{SDS}} = w(t) (\epsilon_\theta(\mathbf{z}_t, \mathbf{c}, t) - \epsilon) \frac{\partial \mathbf{x}}{\partial \phi}, \quad (2)$$

where ϵ is sampled from $\mathcal{N}(\mathbf{0}, \mathbf{1})$, \mathbf{x} is an image rendered by \mathcal{G} , \mathbf{z}_t is obtained by adding Gaussian noise ϵ to \mathbf{x} corresponding to the timestep t of the diffusion process, $w(t)$ is a constant that depends on the noising schedule. Inspired by this method, we propose a noise-invariant video score distillation sampling (NI-VSDS) strategy to efficiently harness the motion prior learned by the text-to-video diffusion model.

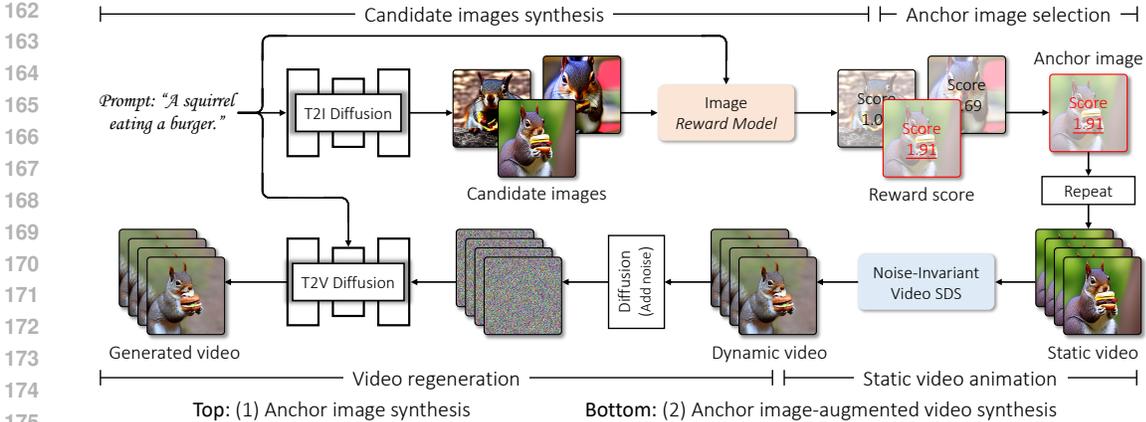


Figure 3: **Illustration of I4VGEN.** I4VGEN is a novel video diffusion inference pipeline, which enhances pre-trained text-to-video diffusion models by incorporating image reference information into the inference process. Instead of the vanilla text-to-video inference pipeline, I4VGEN consists of two stages: (1) anchor image synthesis and (2) anchor image-augmented text-to-video synthesis. Firstly, a simple yet effective generation-selection strategy is applied to synthesize candidate images and select the most suitable image using a reward-based mechanism, thereby obtaining high-quality anchor image. Subsequently, an innovative noise-invariant video scoring distillation sampling (NI-VSDS) is developed, which extracts motion prior from the text-to-video diffusion model to animate the anchor image into dynamic video, followed by a video regeneration process to refine the video.

3 I4VGEN

This section introduces **I4VGEN**, a novel video diffusion inference pipeline designed for enhancing the capabilities of pre-trained text-to-video diffusion models. As illustrated in Fig. 3, we factorize the inference process into two stages: (1) anchor image synthesis to generate the anchor image \mathbf{x} given the text prompt \mathbf{y} , and (2) anchor image-augmented video synthesis to generate the video \mathbf{v} by leveraging the text prompt \mathbf{y} and the anchor image \mathbf{x} . This section provides the detailed explanations of both stages in Sec. 3.1 and 3.2, respectively.

3.1 ANCHOR IMAGE SYNTHESIS

The goal of this stage is to synthesize visually-realistic anchor images \mathbf{x} that accurately correspond to the given text prompts \mathbf{y} . This image serves as a foundation to provide appearance information for enhancing the performance of the subsequent video generation. As illustrated in Fig. 3 (Top), a simple yet effective generation-selection pipeline is employed to produce the anchor image, which consist of candidate images synthesis and reward-based anchor image selection.

Candidate images synthesis. Instead of generating a single image, our approach produces a set of candidate images to ensure the selection of the best example. Utilizing a pre-trained image diffusion model $\mathcal{D}_{\text{img}}(\cdot)$, we construct the candidate image set as follows:

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N = \mathcal{D}_{\text{img}}(\mathbf{y}, \mathbf{z}_1), \mathcal{D}_{\text{img}}(\mathbf{y}, \mathbf{z}_2), \dots, \mathcal{D}_{\text{img}}(\mathbf{y}, \mathbf{z}_N), \quad (3)$$

where N denotes the number of candidate images, and \mathbf{z}_i represents Gaussian noise.

Reward-based anchor image selection. With the help of the image reward model $\mathcal{R}(\cdot)$ (Xu et al., 2023), a promising automatic text-to-image evaluation metric aligned with human preferences, the candidate image with the highest reward score s is selected as the anchor image \mathbf{x} , as defined by:

$$\mathbf{x} = \mathbf{x}_i, \quad \text{where } i = \arg \max_i s_i = \arg \max_i \mathcal{R}(\mathbf{x}_i). \quad (4)$$

The generation-selection design facilitates the acquisition of a high-quality anchor image, particularly beneficial for complex text prompts (see Fig. 5). Notably, our method accommodates both user-provided and retrieved images, extending its applicability to a variety of custom scenarios, as discussed in Sec. 4.5.

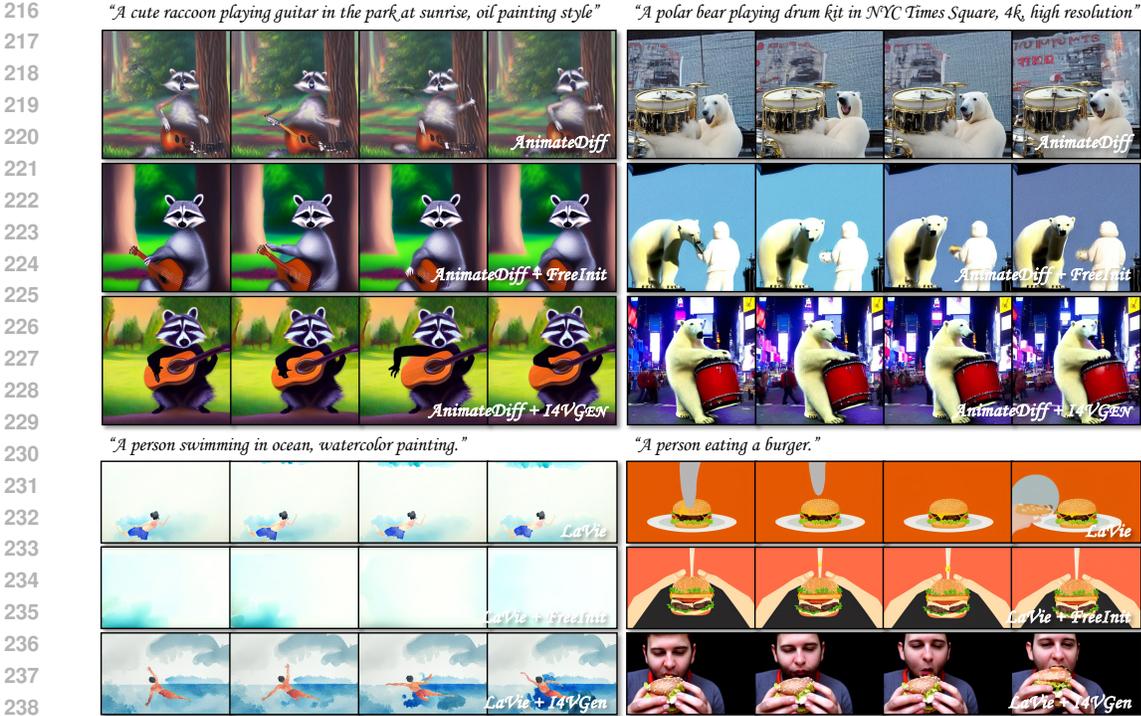


Figure 4: **Qualitative comparison.** Each video is generated with the same text prompt and random seed for all methods. Our approach significantly improves the quality of the generated videos while showing excellent alignment with text prompts.

3.2 ANCHOR IMAGE-AUGMENTED VIDEO SYNTHESIS

Upon obtaining the anchor image \mathbf{x} , we replicate it F times to create an initial static video $\hat{\mathbf{v}} \in \mathcal{V}$, $\hat{\mathbf{v}} = \{\mathbf{x}, \mathbf{x}, \dots, \mathbf{x}\}$. The goal of this stage is to convert this static video into a high-quality video reflecting the text prompt \mathbf{y} . As illustrated in Fig. 3 (Bottom), we introduce static video animation and video regeneration.

Static video animation. A straightforward approach to animate the static video involves applying a diffusion-denoising process to transition from the static to dynamic state. However, this approach still encounters a training-inference gap, as the text-to-video diffusion model is trained on dynamic real-world videos but tested on static videos, leading to sub-optimal motion quality due to the introduction of static priors, as discussed in Sec. 4.4.

To address this limitation, we propose a novel approach leveraging the motion prior from the pre-trained text-to-video diffusion model to animate static videos. Drawing inspiration from score distillation sampling (SDS) as introduced in (Poole et al., 2023; Wang et al., 2023a), we develop the noise-invariant video score distillation sampling (NI-VSDS). Unlike vanilla SDS, which optimizes a parametric image generator, our approach directly parameterizes the static video $\hat{\mathbf{v}}$ and applies targeted optimization to it. The NI-VSDS loss function is defined as follows:

$$\nabla_{\hat{\mathbf{v}}} \mathcal{L}_{\text{NI-VSDS}} = w(t) (\epsilon_{\theta}(\hat{\mathbf{v}}_t, \mathbf{c}, t) - \epsilon), \tag{5}$$

where $\hat{\mathbf{v}}_t$ represents the noisy video at timestep t perturbed by Gaussian noise ϵ . Furthermore, we incorporate three strategic modifications:

- Instead of resampling the Gaussian noise at each iteration as in traditional SDS, we maintain a constant noise across the optimization, enhancing convergence speed.
- Optimization is confined to the initial stages of the denoising process, where noise levels are higher, focusing on dynamic information distillation.
- We implement a coarse-to-fine optimization strategy, evolving from high to low noise levels, specifically from timestep T to $\tau_{\text{NI-VSDS}}$, where $T > \tau_{\text{NI-VSDS}} > 0$. This approach stabilizes the optimization trajectory and yields superior motion quality.

Table 1: **VBench evaluation results per dimension.** This table compares the performance of I4VGEN with other counterparts across each of the 16 VBench dimensions.

Methods	Subj. Cons.	Back. Cons.	Tem. Flick.	Moti. Smo.	Dyna. Degr.	Aest. Qual.	Imag. Qual.	Obj. Class
AnimateDiff	87.11%	95.22%	95.99%	93.12%	74.89%	56.07%	64.29%	83.69%
+ FreeInit	90.45%	96.57%	96.89%	95.66%	70.17%	59.25%	63.51%	87.55%
+ I4VGEN	95.17%	97.73%	98.51%	96.45%	57.72%	64.68%	66.18%	92.59%
LaVie	91.65%	96.30%	98.03%	95.73%	71.94%	59.64%	65.13%	91.25%
+ FreeInit	92.32%	96.35%	98.06%	95.83%	71.11%	59.41%	63.89%	89.13%
+ I4VGEN	94.12%	96.90%	98.55%	96.37%	70.55%	60.88%	66.55%	92.26%
Methods	Mult. Obj.	Hum. Acti.	Color	Spat. Rela.	Scene	Appe. Style	Tem. Style	Over. Cons.
AnimateDiff	22.61%	90.40%	81.73%	31.55%	45.61%	24.40%	24.49%	25.71%
+ FreeInit	26.92%	93.00%	86.39%	30.71%	44.61%	23.98%	25.03%	25.61%
+ I4VGEN	57.22%	95.80%	91.98%	45.20%	54.67%	25.07%	26.11%	28.01%
LaVie	24.02%	94.80%	83.64%	26.27%	52.89%	23.67%	24.94%	27.25%
+ FreeInit	22.59%	94.20%	84.34%	27.46%	52.70%	23.61%	24.85%	26.89%
+ I4VGEN	32.77%	96.20%	88.59%	33.81%	55.64%	24.35%	25.62%	27.68%

The implementation of noise-invariant video score distillation sampling (NI-VSDS) algorithm is detailed in Algorithm 1, which outlines the process of converting a static video into a dynamic video using the defined NI-VSDS loss. Notably, we only perform a single update from timestep T to $\tau_{\text{NI-VSDS}}$, requiring fewer than 50 iterations, this is a significant reduction compared to the thousands of iterations typically required for text-to-3D synthesis in SDS. α is a scalar that defines the step size of the gradient update. We empirically set $\tau_{\text{NI-VSDS}} = \text{Int}(T \times p_{\text{NI-VSDS}})$.

Video regeneration. After animating the static video, we further enhance the appearance detail quality of the video through a diffusion-denoising process. This stage is not affected by the aforementioned training-inference gap, thereby achieving more refined generation results.

Notably, we can flexibly add noise up to any timestep τ_{re} , calculated as $\tau_{\text{re}} = \text{Int}(T \times p_{\text{re}})$, followed by the corresponding denoising process. This strategy not only preserves the fine appearance textures but also reduces the required denoising steps, thus streamlining the video synthesis process and elevating the overall quality of the resulting video.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETTINGS

Implementation details. I4VGEN a novel video diffusion inference pipeline that leverages advanced image techniques to enhance pre-trained text-to-video diffusion models without requiring additional training, and can be seamlessly integrated into existing text-to-video diffusion models. To ascertain the efficacy and adaptability of I4VGEN, we apply it to two well-regarded text-to-video diffusion models: AnimateDiff (Guo et al., 2024b) and LaVie (Wang et al., 2023e).

For AnimateDiff, the mm-sd-v15-v2 motion module¹, alongside Stable Diffusion v1.5, is utilized to synthesize 16 consecutive frames at a resolution of 512×512 pixels for evaluation. For LaVie, the

Algorithm 1: NI-VSDS

Input: T2V diffusion model $\epsilon_{\theta}(\cdot)$, text prompt \mathbf{y} , static video $\hat{\mathbf{v}}$, timestep $\tau_{\text{NI-VSDS}}$.

Output: Dynamic video.

- 1 Sampling $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$; $\mathbf{c} = f_{\text{CLIP}}(\mathbf{y})$
- 2 **for** $t = T, \dots, \tau_{\text{NI-VSDS}}$ **do**
- 3 $\hat{\mathbf{v}}_t \leftarrow \text{AddNoise}(\hat{\mathbf{v}}, \epsilon, t)$
- 4 $\nabla_{\hat{\mathbf{v}}} \mathcal{L}_{\text{NI-VSDS}} \leftarrow w(t) (\epsilon_{\theta}(\hat{\mathbf{v}}_t, \mathbf{c}, t) - \epsilon)$
- 5 $\hat{\mathbf{v}} \leftarrow \hat{\mathbf{v}} - \alpha \cdot \nabla_{\hat{\mathbf{v}}} \mathcal{L}_{\text{NI-VSDS}}$
- 6 **return** $\hat{\mathbf{v}}$

¹<https://github.com/guoyww/AnimateDiff>

Table 2: **Ablation study.** Orange highlights generation-selection, while yellow highlights NI-VSDS.

Methods	Subj. Cons.	Back. Cons.	Tem. Flick.	Moti. Smo.	Dyna. Degr.	Aest. Qual.	Imag. Qual.	Obje. Class
AnimateDiff	87.11%	95.22%	95.99%	93.12%	74.89%	56.07%	64.29%	83.69%
+ I4VGEN (w/o gen.-sel.)	94.89%	97.80%	98.28%	96.99%	55.91%	62.23%	64.18%	90.95%
+ I4VGEN (w/o NI-VSDS)	96.47%	98.82%	98.99%	97.56%	28.24%	65.17%	65.52%	92.66%
+ I4VGEN	95.17%	97.73%	98.51%	96.45%	57.72%	64.68%	66.18%	92.59%

Methods	Mult. Obj.	Hum. Acti.	Color	Spat. Rela.	Scene	Appe. Style	Tem. Style	Over. Cons.
AnimateDiff	22.61%	90.40%	81.73%	31.55%	45.61%	24.40%	24.49%	25.71%
+ I4VGEN (w/o gen.-sel.)	40.68%	94.40%	90.55%	37.79%	53.72%	24.76%	26.03%	26.62%
+ I4VGEN (w/o NI-VSDS)	62.84%	94.80%	91.95%	47.57%	55.80%	24.88%	25.72%	27.91%
+ I4VGEN	57.22%	95.80%	91.98%	45.20%	54.67%	25.07%	26.11%	28.01%

base-version² is employed to generate 16 consecutive frames at 320×512 pixels for evaluation. All other inference details adhere to the original settings described in Guo et al. (2024b) and Wang et al. (2023e), respectively. Notably, both AnimateDiff and LaVie possess inherent text-to-image generation capabilities when excluding the motion module. To avoid introducing additional GPU storage requirements, we leverage their corresponding image versions for text-to-image generation in I4VGEN. For AnimateDiff, we empirically set $N = 16$, $p_{\text{NI-VSDS}} = 0.4$, $\alpha = 1$, and $p_{\text{re}} = 1$. For LaVie, we empirically set $N = 16$, $p_{\text{NI-VSDS}} = 0.4$, $\alpha = 1$, and $p_{\text{re}} = 0.8$. All experiments are conducted on a single NVIDIA V100 GPU (32 GB).

Benchmark. I4VGEN is assessed using VBench (Huang et al., 2024), a comprehensive benchmark that evaluates video generation models across 16 disentangled dimensions, which is more authoritative than FVD. These dimensions provide a detailed analysis of generation quality from two overarching perspectives: video quality³, focusing on the perceptual quality of the generated videos, and video-condition consistency⁴, assessing how well the generated videos align with the provided conditions.

4.2 QUALITATIVE COMPARISON

Fig. 4 presents a comparative analysis of our results against state-of-the-art counterparts using identical text prompts and random seeds. I4VGEN excels in enhancing both the temporal consistency and the frame-wise quality, alongside superior alignment with the text prompts. For instance, in the case of “playing guitar”, AnimateDiff suffers from poor video quality, and FreeInit encounters an incomplete guitar in the middle of the video. In contrast, our method effectively addresses these issues, maintaining stable temporal consistency. Furthermore, while baseline methods struggle with accurate synthesis of all text-described components, e.g., “NYC Times Square”, I4VGEN generates videos that are visually realistic and closely aligned with the text prompts by utilizing anchor images obtained by the generation-selection strategy.

4.3 QUANTITATIVE COMPARISON

Objective evaluation. Following the protocols established by VBench, we evaluate I4VGEN in terms of both video quality and video-text consistency. As detailed in Table 1, I4VGEN outperforms all other approaches in temporal quality (higher background and subject consistency, less flickering, and better smoothness), frame-wise quality (higher aesthetic and imaging quality), and video-text

²<https://github.com/Vchitect/LaVie>

³Video quality includes 7 evaluation dimensions: Subject Consistency, Background Consistency, Temporal Flickering, Motion Smoothness, Dynamic Degree, Aesthetic Quality, and Imaging Quality. The first 5 evaluate temporal quality, and the last 2 evaluate frame-wise quality.

⁴Video-condition consistency includes 9 evaluation dimensions: Object Class, Multiple Objects, Human Action, Color, Spatial Relationship, Scene, Appearance Style, Temporal Style, Overall Consistency. The first 6 evaluate semantics, the 7 and 8-th evaluate style, and the 9-th evaluates overall consistency.



388 Figure 5: **Intermediate results visualization.** We provide visualizations of the candidate images with reward scores, the dynamic video, and the corresponding generated video.

389 consistency (greater semantics, style, and overall consistency). Although counterparts occasionally produce videos with more dynamic motion, they are often linked to inappropriate or excessive movements. I4VGEN strikes a more effective balance between motion intensity and overall video quality, which is further verified in the user study.

390

391

392 **User study.** We conduct a subjective user study involving 20 volunteers with expertise in image and video processing, with each participant answering 15 questions. Specifically, participants are asked to select the video with the highest quality across three dimensions: video quality, video-condition consistency, and overall score. As shown in Table 3, our approach outperforms the other methods favorably.

393 Table 3: **User study.**

Method	Video Quality	Vid.-Cond. Consistency	Overall score
AnimateDiff	6.00%	10.67%	6.33%
+ FreeInit	27.67%	15.67%	25.00%
+ I4VGEN	66.33%	73.67%	68.67%
LaVie	27.67%	21.33%	22.33%
+ FreeInit	22.67%	18.33%	19.67%
+ I4VGEN	49.67%	60.33%	58.00%

394

395

396 **Inference time.** We define the time cost of a single denoising iteration for a video in a video diffusion model as c . For AnimateDiff (Guo et al., 2024b), following the original inference setting, the time cost to generate a single 16-frame video is $25c$. FreeInit requires 5 rounds of diffusion-denoising to generate a single video, taking a time of $5 \times 25c = 125c$. The time cost for I4VGEN to generate a single video is: $< 25c$ (for synthesizing 16 candidate images) + $0.6 \times 25c$ (for NI-VSDS) + $\leq 25c$ (for video regeneration) = $< 65c$ (**total cost**), making it more efficient compared to FreeInit. LaVie (Wang et al., 2023e) shares the same conclusion.

397 Table 4: **Inference time.**

Method	Time
AnimateDiff	21.73s
AnimateDiff + FreeInit	113.67s
AnimateDiff + I4VGEN	53.78s

398 We also provide the inference time for a single video in Table 4, evaluated on a single NVIDIA V100 GPU (32 GB), where 50 videos are randomly generated to obtain an average inference time. Our method performs better than FreeInit.

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

4.4 ABLATION STUDY

On generation-selection strategy. We adopt a generation-selection strategy to create visually-realistic and semantically-faithful anchor images, which serve as a foundation for providing appearance information to enhance subsequent video generation performance. As shown in Table 2, highlighted in orange, compared to randomly synthesizing a single anchor image, the generation-selection strategy significantly improves the quality of the generated videos in terms of frame-wise quality and consistency with the text. Fig. 5 provides a visualization of the candidate images, where the reward-based selection strategy eliminates unsatisfactory images, leading to better results.

On NI-VSDS. Directly applying the video regeneration process to static videos introduces static priors, resulting in suboptimal motion quality. As shown in Table 2, highlighted in yellow, while direct diffusion-denoising improves the temporal consistency of the generated videos, it severely sacrifices the motion dynamics, adversely affecting the motion style. In contrast, our method achieves an effective balance between motion intensity and overall video quality.



442 Figure 6: **Adaptation on SparseCtrl.** I4VGEN can be seamlessly integrated into SparseCtrl by
 443 replacing the anchor image with the provided image, leading to improved results.
 444

445

446 **On video regeneration.** Fig. 5 visualizes the intermediate results, demonstrating that the video
 447 regeneration process is essential for refining appearance details.
 448

449 4.5 MORE APPLICATIONS

450

451 **Adaptation on real image.** Our method
 452 adapts to user-provided images, as shown
 453 in Fig. 7, where we use real images as
 454 anchor images, resulting in high-fidelity
 455 videos that are semantically consistent
 456 with the real images. Notably, our ap-
 457 proach differs from vanilla image-to-video
 458 generation, as the synthesized videos are
 459 not completely aligned with the provided
 460 images. NI-VSDS is designed to ani-
 461 mate static images and is implemented as
 a spatio-temporal co-optimization.



479 Figure 7: **Adaptation on real image.**

480

481 **Adaptation on image-to-video diffusion models.** I4VGEN can be seamlessly integrated into ex-
 482 isting image-to-video diffusion models by replacing the anchor images with the provided images,
 483 thereby enhancing the overall video quality. As shown in Fig. 6, integrating I4VGEN into SparseC-
 484 ctrl (Guo et al., 2023) significantly improves the quality of the generated videos in terms of temporal
 485 consistency and appearance fidelity.

486 5 CONCLUSION

487

488 The paper introduces I4VGEN, a novel video diffusion inference pipeline to leverage advanced im-
 489 age techniques to enhance pre-trained text-to-video diffusion models, which requires no additional
 490 learnable parameters and training costs. I4VGEN decomposes the text-to-video inference process
 491 into anchor image synthesis and anchor image-augmented video synthesis. Correspondingly, a sim-
 492 ple yet effective generation-selection strategy is applied to produce a high-quality anchor image, and
 493 an innovative noise-invariant video score distillation sampling (NI-VSDS) is designed to animate the
 494 image, followed by a video regeneration process to enhance the final output. I4VGEN effectively
 495 alleviates non-zero terminal signal-to-noise ratio issues and demonstrates improved visual realism
 496 and textual fidelity when integrated with existing video diffusion models.

497

498 **Limitation and discussion.** I4VGEN improves the video diffusion model but requires more in-
 499 ference cost. As discussed in Sec. 4.3, the inference time of I4VGEN is over double the baseline.
 500 Enhancing inference efficiency remains a future goal, with distillation techniques as a potential ap-
 501 proach. Furthermore, removing the generation-selection strategy can reduce inference costs to some
 502 extent. As shown in Table 2, our method still significantly outperforms the baseline under this set-
 503 ting. Additionally, although our method and FreeInit are orthogonal, integrating both by replacing
 504 video regeneration with FreeInit fails to produce notable benefits.
 505

REFERENCES

- 486
487
488 Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika
489 Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, et al. ediffi: Text-to-image diffusion models
490 with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022.
- 491 Fan Bao, Chendong Xiang, Gang Yue, Guande He, Hongzhou Zhu, Kaiwen Zheng, Min Zhao,
492 Shilong Liu, Yaole Wang, and Jun Zhu. Vidu: a highly consistent, dynamic and skilled text-to-
493 video generator with diffusion models. *arXiv preprint arXiv:2405.04233*, 2024.
- 494
495 Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik
496 Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling
497 latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023a.
- 498 Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler,
499 and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion
500 models. In *CVPR*, 2023b.
- 501
502 Tim Brooks, Janne Hellsten, Miika Aittala, Ting-Chun Wang, Timo Aila, Jaakko Lehtinen, Ming-Yu
503 Liu, Alexei A Efros, and Tero Karras. Generating long videos of dynamic scenes. In *NeurIPS*,
504 2022.
- 505 Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing,
506 Yaofang Liu, Qifeng Chen, Xintao Wang, et al. Videocrafter1: Open diffusion models for high-
507 quality video generation. *arXiv preprint arXiv:2310.19512*, 2023a.
- 508
509 Nanxin Chen, Yu Zhang, Heiga Zen, Ron J Weiss, Mohammad Norouzi, and William Chan. Wave-
510 grad: Estimating gradients for waveform generation. In *ICLR*, 2021.
- 511 Tsai-Shien Chen, Chieh Hubert Lin, Hung-Yu Tseng, Tsung-Yi Lin, and Ming-Hsuan Yang. Motion-
512 conditioned diffusion model for controllable video synthesis. *arXiv preprint arXiv:2304.14404*,
513 2023b.
- 514
515 Xinyuan Chen, Yaohui Wang, Lingjun Zhang, Shaobin Zhuang, Xin Ma, Jiashuo Yu, Yali Wang,
516 Dahua Lin, Yu Qiao, and Ziwei Liu. Seine: Short-to-long video diffusion model for generative
517 transition and prediction. In *ICLR*, 2024a.
- 518 Xuweiyi Chen, Tian Xia, and Sihan Xu. Unictrl: Improving the spatiotemporal consistency
519 of text-to-video diffusion models via training-free unified attention control. *arXiv preprint*
520 *arXiv:2403.02332*, 2024b.
- 521
522 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In
523 *NeurIPS*, 2021.
- 524 Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image
525 synthesis. In *CVPR*, 2021.
- 526
527 Patrick Esser, Johnathan Chiu, Parmida Atighehchian, Jonathan Granskog, and Anastasis Germani-
528 dis. Structure and content-guided video synthesis with diffusion models. In *ICCV*, 2023.
- 529 Martin Nicolas Everaert, Athanasios Fitsios, Marco Bocchio, Sami Arpa, Sabine Süssstrunk, and
530 Radhakrishna Achanta. Exploiting the signal-leak bias in diffusion models. In *WACV*, 2024.
- 531
532 Mengyang Feng, Jinlin Liu, Kai Yu, Yuan Yao, Zheng Hui, Xiefan Guo, Xianhui Lin, Haolan Xue,
533 Chen Shi, Xiaowen Li, et al. Dreamoving: A human video generation framework based on
534 diffusion models. *arXiv preprint arXiv:2312.05107*, 2023.
- 535 Gereon Fox, Ayush Tewari, Mohamed Elgharib, and Christian Theobalt. Stylevideogan: A temporal
536 generative model using a pretrained stylegan. In *BMVC*, 2021.
- 537
538 Tsu-Jui Fu, Licheng Yu, Ning Zhang, Cheng-Yang Fu, Jong-Chyi Su, William Yang Wang, and Sean
539 Bell. Tell me what happened: Unifying text-guided video completion via multimodal masked
video generation. In *CVPR*, 2023.

- 540 Songwei Ge, Thomas Hayes, Harry Yang, Xi Yin, Guan Pang, David Jacobs, Jia-Bin Huang, and
541 Devi Parikh. Long video generation with time-agnostic vqgan and time-sensitive transformer. In
542 *ECCV*, 2022.
- 543 Songwei Ge, Seungjun Nah, Guilin Liu, Tyler Poon, Andrew Tao, Bryan Catanzaro, David Jacobs,
544 Jia-Bin Huang, Ming-Yu Liu, and Yogesh Balaji. Preserve your own correlation: A noise prior
545 for video diffusion models. In *ICCV*, 2023.
- 546 Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Tokenflow: Consistent diffusion features
547 for consistent video editing. In *ICLR*, 2024.
- 548 Rohit Girdhar, Mannat Singh, Andrew Brown, Quentin Duval, Samaneh Azadi, Sai Saketh Ramb-
549 hatla, Akbar Shah, Xi Yin, Devi Parikh, and Ishan Misra. Emu video: Factorizing text-to-video
550 generation by explicit image conditioning. *arXiv preprint arXiv:2311.10709*, 2023.
- 551 Jiatao Gu, Shuangfei Zhai, Yizhe Zhang, Josh Susskind, and Navdeep Jaitly. Matryoshka diffusion
552 models. *arXiv preprint arXiv:2310.15111*, 2023.
- 553 Xiefan Guo, Jinlin Liu, Miaomiao Cui, Jiankai Li, Hongyu Yang, and Di Huang. Initno: Boosting
554 text-to-image diffusion models via initial noise optimization. In *CVPR*, 2024a.
- 555 Yuwei Guo, Ceyuan Yang, Anyi Rao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Sparsectrl:
556 Adding sparse controls to text-to-video diffusion models. *arXiv preprint arXiv:2311.16933*, 2023.
- 557 Yuwei Guo, Ceyuan Yang, Anyi Rao, Yaohui Wang, Yu Qiao, Dahua Lin, and Bo Dai. Animatediff:
558 Animate your personalized text-to-image diffusion models without specific tuning. In *ICLR*,
559 2024b.
- 560 Agrim Gupta, Lijun Yu, Kihyuk Sohn, Xiuye Gu, Meera Hahn, Li Fei-Fei, Irfan Essa, Lu Jiang,
561 and José Lezama. Photorealistic video generation with diffusion models. *arXiv preprint*
562 *arXiv:2312.06662*, 2023.
- 563 Nicholas Guttenberg. Diffusion with offset noise. [https://www.crosslabs.org/blog/
564 diffusion-with-offset-noise](https://www.crosslabs.org/blog/diffusion-with-offset-noise).
- 565 Jiawei He, Andreas Lehrmann, Joseph Marino, Greg Mori, and Leonid Sigal. Probabilistic video
566 generation using holistic attribute control. In *ECCV*, 2018.
- 567 Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent video diffusion mod-
568 els for high-fidelity video generation with arbitrary lengths. *arXiv preprint arXiv:2211.13221*,
569 2022.
- 570 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*,
571 2020.
- 572 Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P
573 Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition
574 video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022a.
- 575 Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J
576 Fleet. Video diffusion models. In *NeurIPS*, 2022b.
- 577 Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pre-
578 training for text-to-video generation via transformers. In *ICLR*, 2023.
- 579 Li Hu, Xin Gao, Peng Zhang, Ke Sun, Bang Zhang, and Liefeng Bo. Animate anyone:
580 Consistent and controllable image-to-video synthesis for character animation. *arXiv preprint*
581 *arXiv:2311.17117*, 2023.
- 582 Ziqi Huang, Yanan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianx-
583 ing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for
584 video generative models. In *CVPR*, 2024.
- 585 Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile
586 diffusion model for audio synthesis. In *ICLR*, 2021.

- 594 Yichong Leng, Zehua Chen, Junliang Guo, Haohe Liu, Jiawei Chen, Xu Tan, Danilo Mandic, Lei
595 He, Xiangyang Li, Tao Qin, et al. Binauralgrad: A two-stage conditional diffusion probabilistic
596 model for binaural audio synthesis. In *NeurIPS*, 2022.
- 597 Mingxiao Li, Tingyu Qu, Wei Sun, and Marie-Francine Moens. Alleviating exposure bias in diffu-
598 sion models through sampling with shifted time steps. In *ICLR*, 2024.
- 600 Xin Li, Wenqing Chu, Ye Wu, Weihang Yuan, Fanglong Liu, Qi Zhang, Fu Li, Haocheng Feng,
601 Errui Ding, and Jingdong Wang. Videogen: A reference-guided latent diffusion approach for
602 high definition text-to-video generation. *arXiv preprint arXiv:2309.00398*, 2023.
- 603 Yitong Li, Martin Min, Dinghan Shen, David Carlson, and Lawrence Carin. Video generation from
604 text. In *AAAI*, 2018.
- 606 Shanchuan Lin, Bingchen Liu, Jiashi Li, and Xiao Yang. Common diffusion noise schedules and
607 sample steps are flawed. In *WACV*, 2024.
- 608 Jinglin Liu, Chengxi Li, Yi Ren, Feiyang Chen, and Zhou Zhao. Diffsinger: Singing voice synthesis
609 via shallow diffusion mechanism. In *AAAI*, 2022.
- 611 Zhengxiong Luo, Dayou Chen, Yingya Zhang, Yan Huang, Liang Wang, Yujun Shen, Deli Zhao,
612 Jingren Zhou, and Tieniu Tan. Videofusion: Decomposed diffusion models for high-quality video
613 generation. In *CVPR*, 2023.
- 614 Shijie Ma, Huayi Xu, Mengjian Li, Weidong Geng, Meng Wang, and Yaxiong Wang. Optimal noise
615 pursuit for augmenting text-to-video generation. *arXiv preprint arXiv:2311.00949*, 2023.
- 617 Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon.
618 Sdedit: Guided image synthesis and editing with stochastic differential equations. In *ICLR*, 2022.
- 619 Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and
620 Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.
- 622 Gaurav Mittal, Tanya Marwah, and Vineeth N Balasubramanian. Sync-draw: Automatic video
623 generation using deep recurrent attentive architectures. In *ACM MM*, 2017.
- 624 Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob
625 McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and
626 editing with text-guided diffusion models. In *ICML*, 2022.
- 628 Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d
629 diffusion. In *ICLR*, 2023.
- 630 Vadim Popov, Ivan Vovk, Vladimir Gogoryan, Tasnima Sadekova, and Mikhail Kudinov. Grad-tts:
631 A diffusion probabilistic model for text-to-speech. In *ICML*, 2021.
- 633 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
634 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
635 models from natural language supervision. In *ICML*, 2021.
- 636 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-
637 conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- 638 Weiming Ren, Harry Yang, Ge Zhang, Cong Wei, Xinrun Du, Stephen Huang, and Wenhui
639 Chen. Consisti2v: Enhancing visual consistency for image-to-video generation. *arXiv preprint*
640 *arXiv:2402.04324*, 2024.
- 642 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
643 resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- 644 Ludan Ruan, Yiyang Ma, Huan Yang, Huiguo He, Bei Liu, Jianlong Fu, Nicholas Jing Yuan, Qin
645 Jin, and Baining Guo. Mm-diffusion: Learning multi-modal diffusion models for joint audio and
646 video generation. In *CVPR*, 2023.
- 647

- 648 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar
649 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic
650 text-to-image diffusion models with deep language understanding. In *NeurIPS*, 2022.
- 651 Masaki Saito, Eiichi Matsumoto, and Shunta Saito. Temporal generative adversarial nets with sin-
652 gular value clipping. In *ICCV*, 2017.
- 653 Masaki Saito, Shunta Saito, Masanori Koyama, and Sosuke Kobayashi. Train sparsely, generate
654 densely: Memory-efficient unsupervised training of high-resolution temporal gan. *IJCV*, 2020.
- 655 Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Mostgan-v: Video generation with temporal
656 motion styles. In *CVPR*, 2023.
- 657 Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry
658 Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video
659 data. In *ICLR*, 2023.
- 660 Ivan Skorokhodov, Sergey Tulyakov, and Mohamed Elhoseiny. Stylegan-v: A continuous video
661 generator with the price, image quality and perks of stylegan2. In *CVPR*, 2022.
- 662 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*,
663 2021.
- 664 Yu Tian, Jian Ren, Menglei Chai, Kyle Olszewski, Xi Peng, Dimitris N Metaxas, and Sergey
665 Tulyakov. A good image generator is what you need for high-resolution video synthesis. In
666 *ICLR*, 2021.
- 667 Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. Mocogan: Decomposing motion
668 and content for video generation. In *CVPR*, 2018.
- 669 Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kindermans, Hernan Moraldo, Han Zhang,
670 Mohammad Taghi Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. Phenaki: Variable
671 length video generation from open domain textual descriptions. In *ICLR*, 2023.
- 672 Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics.
673 In *NIPS*, 2016.
- 674 Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A Yeh, and Greg Shakhnarovich. Score jacobian
675 chaining: Lifting pretrained 2d diffusion models for 3d generation. In *CVPR*, 2023a.
- 676 Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Mod-
677 elscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023b.
- 678 Xiang Wang, Hangjie Yuan, Shiwei Zhang, Dayou Chen, Jiuniu Wang, Yingya Zhang, Yujun Shen,
679 Deli Zhao, and Jingren Zhou. Videocomposer: Compositional video synthesis with motion con-
680 trollability. *arXiv preprint arXiv:2306.02018*, 2023c.
- 681 Xiang Wang, Shiwei Zhang, Han Zhang, Yu Liu, Yingya Zhang, Changxin Gao, and Nong Sang.
682 Videolcm: Video latent consistency model. *arXiv preprint arXiv:2312.09109*, 2023d.
- 683 Yaohui Wang, Piotr Bilinski, Francois Bremond, and Antitza Dantcheva. G3an: Disentangling
684 appearance and motion for video generation. In *CVPR*, 2020.
- 685 Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yanan
686 He, Jiashuo Yu, Peiqing Yang, et al. Lavie: High-quality video generation with cascaded latent
687 diffusion models. *arXiv preprint arXiv:2309.15103*, 2023e.
- 688 Yuhan Wang, Liming Jiang, and Chen Change Loy. Styleinv: A temporal style modulated inversion
689 network for unconditional video generation. In *ICCV*, 2023f.
- 690 Chenfei Wu, Jian Liang, Lei Ji, Fan Yang, Yuejian Fang, Daxin Jiang, and Nan Duan. Nüwa: Visual
691 synthesis pre-training for neural visual world creation. In *ECCV*, 2022.

- 702 Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu,
703 Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion
704 models for text-to-video generation. In *ICCV*, 2023a.
- 705
706 Tianxing Wu, Chenyang Si, Yuming Jiang, Ziqi Huang, and Ziwei Liu. Freeinit: Bridging initial-
707 ization gap in video diffusion models. *arXiv preprint arXiv:2312.07537*, 2023b.
- 708 Jinbo Xing, Menghan Xia, Yong Zhang, Haoxin Chen, Wangbo Yu, Hanyuan Liu, Gongye Liu,
709 Xintao Wang, Ying Shan, and Tien-Tsin Wong. Dynamicrafter: Animating open-domain images
710 with video diffusion priors. In *ECCV*, 2024.
- 711
712 Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao
713 Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation.
714 In *NeurIPS*, 2023.
- 715
716 Zeyue Xue, Guanglu Song, Qiushan Guo, Boxiao Liu, Zhuofan Zong, Yu Liu, and Ping Luo.
717 Raphael: Text-to-image generation via large mixture of diffusion paths. In *NeurIPS*, 2023.
- 718
719 Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. Videogpt: Video generation using
720 vq-vae and transformers. *arXiv preprint arXiv:2104.10157*, 2021.
- 721
722 Jaehoon Yoo, Semin Kim, Doyup Lee, Chiheon Kim, and Seunghoon Hong. Towards end-to-end
723 generative modeling of long videos with memory-efficient bidirectional transformers. In *CVPR*,
724 2023.
- 725
726 Jiwen Yu, Xiaodong Cun, Chenyang Qi, Yong Zhang, Xintao Wang, Ying Shan, and Jian
727 Zhang. Animatezero: Video diffusion models are zero-shot image animators. *arXiv preprint*
728 *arXiv:2312.03793*, 2023a.
- 729
730 Lijun Yu, Yong Cheng, Kihyuk Sohn, José Lezama, Han Zhang, Huiwen Chang, Alexander G
731 Hauptmann, Ming-Hsuan Yang, Yuan Hao, Irfan Essa, et al. Magvit: Masked generative video
732 transformer. In *CVPR*, 2023b.
- 733
734 Sihyun Yu, Jihoon Tack, Sangwoo Mo, Hyunsu Kim, Junho Kim, Jung-Woo Ha, and Jinwoo Shin.
735 Generating videos with dynamics-aware implicit generative adversarial networks. In *ICLR*, 2022.
- 736
737 Sihyun Yu, Kihyuk Sohn, Subin Kim, and Jinwoo Shin. Video probabilistic diffusion models in
738 projected latent space. In *CVPR*, 2023c.
- 739
740 David Junhao Zhang, Jay Zhangjie Wu, Jia-Wei Liu, Rui Zhao, Lingmin Ran, Yuchao Gu, Difei
741 Gao, and Mike Zheng Shou. Show-1: Marrying pixel and latent diffusion models for text-to-
742 video generation. *arXiv preprint arXiv:2309.15818*, 2023a.
- 743
744 Shiwei Zhang, Jiayu Wang, Yingya Zhang, Kang Zhao, Hangjie Yuan, Zhiwu Qin, Xiang Wang,
745 Deli Zhao, and Jingren Zhou. I2vgen-xl: High-quality image-to-video synthesis via cascaded
746 diffusion models. *arXiv preprint arXiv:2311.04145*, 2023b.
- 747
748 Daquan Zhou, Weimin Wang, Hanshu Yan, Weiwei Lv, Yizhe Zhu, and Jiashi Feng. Magicvideo:
749 Efficient video generation with latent diffusion models. *arXiv preprint arXiv:2211.11018*, 2022.
750
751
752
753
754
755

This appendix is structured as follows:

- In Appendix A, we provide a discussion of related work.
- In Appendix B, we provide additional experiment results and analysis.
- In Appendix C, we provide the code for I4VGEN.

A RELATED WORK

Video Generative Models. The domain of video generation has seen significant advancements through the use of Generative Adversarial Networks (GANs) (Vondrick et al., 2016; Saito et al., 2017; Tulyakov et al., 2018; Wang et al., 2020; Saito et al., 2020; Tian et al., 2021; Fox et al., 2021; Yu et al., 2022; Skorokhodov et al., 2022; Brooks et al., 2022; Shen et al., 2023; Wang et al., 2023f), Variational Autoencoders (VAEs) (Mittal et al., 2017; Li et al., 2018; He et al., 2018), and Autoregressive models (ARs) (Yan et al., 2021; Ge et al., 2022; Wu et al., 2022; Hong et al., 2023; Villegas et al., 2023; Fu et al., 2023; Yoo et al., 2023; Yu et al., 2023b). Despite these developments, synthesizing videos from text prompts remains challenging due to the complexities of modeling spatio-temporal dynamics. Recent innovations driven by the successes of diffusion models (Ho et al., 2020; Dhariwal & Nichol, 2021; Song et al., 2021), which have been applied effectively in image generation (Rombach et al., 2022; Nichol et al., 2022; Ramesh et al., 2022; Saharia et al., 2022; Gu et al., 2023; Balaji et al., 2022; Xue et al., 2023; Meng et al., 2022; Guo et al., 2024a) and audio synthesis (Kong et al., 2021; Chen et al., 2021; Popov et al., 2021; Leng et al., 2022; Liu et al., 2022), and underscore the emergence of substantial headway (Ho et al., 2022b;a; He et al., 2022; Singer et al., 2023; Blattmann et al., 2023b; Yu et al., 2023c; Ruan et al., 2023; Wu et al., 2023a; Chen et al., 2023a;b; Esser et al., 2023; Ge et al., 2023; Chen et al., 2024a; Geyer et al., 2024; Ma et al., 2023; Wang et al., 2023e; Zhang et al., 2023a;b; Hu et al., 2023; Wang et al., 2023d; Feng et al., 2023; Guo et al., 2024b; Girdhar et al., 2023; Blattmann et al., 2023a; Gupta et al., 2023; Wang et al., 2023b;c; Luo et al., 2023) in research endeavors devoted to video synthesis from text input.

The foundational contributions of the Video Diffusion Model (VDM) (Ho et al., 2022b) represents a milestone in leveraging diffusion models for video generation by adapting the 2D U-Net architecture used in image generation to a 3D U-Net capable of temporal modeling. Successive researches, such as Make-A-Video (Singer et al., 2023) and Imagen Video (Ho et al., 2022a), expand video generation capabilities significantly. To enhance efficiency, subsequent models have transitioned the diffusion process from pixel to latent space (He et al., 2022; Zhou et al., 2022; Wang et al., 2023b; Blattmann et al., 2023b;a; Guo et al., 2024b; Wang et al., 2023e), paralleling advancements in latent diffusion for images (Rombach et al., 2022).

However, the direct generation of videos from text prompts remains intrinsically challenging. Recent approaches (Blattmann et al., 2023a; Zhang et al., 2023b; Girdhar et al., 2023; Chen et al., 2023a; 2024a; Li et al., 2023; Hu et al., 2023; Yu et al., 2023a; Ren et al., 2024) have employed text-to-image synthesis as an intermediary step, enhancing overall performance. Despite these advancements, these methods still face the challenge of high computational training costs. In this study, we explore a novel training-free methodology aimed at bridging the existing gap in the field.

In addition, (Chen et al., 2024b) (contemporary researches) introduces additional operations in the attention layer, *i.e.*, cross-frame self-attention control, to enhance the video model. However, this necessitates modifications to the model architecture, whereas our method does not.

Signal-Leak Bias. Diffusion models are designed to generate high-quality visuals from noise through a sequential denoising process, which is consistent in both image and video diffusion models. During training, Gaussian noise corrupts the visual content, challenging the model to restore it to its original form. In the inference phase, the model operates on pure Gaussian noise, transforming it into a realistic visual content step-by-step.

Unfortunately, most existing diffusion models exhibit a disparity between the corrupted image during training and the pure Gaussian noise during inference. Commencing denoising from pure Gaussian noise in the inference phase deviates from the training process, potentially introducing *signal-leak bias*. For image diffusion models, (Guttenberg; Lin et al., 2024; Li et al., 2024) point out flaws in common diffusion noise schedules and sample steps, and propose to fine-tune the diffusion model

810 to mitigate or eliminate the signal-leak bias during training, leading to improved results. (Everaert
811 et al., 2024) attempts to exploit signal-leak bias to achieve more control over the generated images.
812 For video diffusion models, this issue becomes more pronounced. (Wu et al., 2023b; Ma et al.,
813 2023) invert the retrieved video or generated low-quality to construct initial noise to alleviate the
814 problem of signal-leak, improving inference quality. However, they suffer from limited diversity
815 and cumbersome inference. At the same time, first-round inference of FreeInit (Wu et al., 2023b)
816 still exhibits a training-inference gap.

817 In contrast to existing methods, our approach utilizes images as the stepping stone for text-to-
818 video generation. This novel pathway aims to produce visually-realistic and semantically-reasonable
819 videos while maintaining manageable computational overheads, as detailed in Sec. 4.

821 B EXPERIMENTS

823 B.1 QUALITATIVE COMPARISON

825 We provide more visualization results in Fig. 8, it can be seen that our method generates more
826 semantically plausible and photo-realistic results than its counterparts. We provide the videos shown
827 in the main paper and appendix in mp4 format in the Supplementary material.

829 B.2 QUANTITATIVE COMPARISON

831 **On hyperparameters.** I4VGEN is a training-free method that improves video generation perfor-
832 mance by correcting the inference process. It is obvious that I4VGEN is also a case-wise method,
833 where different cases correspond to different optimal hyperparameters. In this paper, we provide an
834 empirical setting that is mild for most instances, serving as a performance lower bound for I4VGEN,
835 and facilitating large-scale quantitative comparisons. Furthermore, we also provide a visualization
836 of the impact of hyperparameters in Fig. 9, which shows that carefully tuned hyperparameters can
837 achieve higher-quality videos.

838 B.3 FAILURE CASES AND DISCUSSIONS

840 We provide the failure cases in Fig. 10, I4VGEN is designed to fully unleash the potential of ex-
841 isting video diffusion models, but it still fails to synthesize high-quality videos that are out of the
842 distribution.

844 C CODE

846 We also provide the code for I4VGEN in the Supplementary material.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917



Figure 8: **Qualitative comparison.** Each video is generated with the same text prompt and random seed for all methods. Our approach significantly improves the quality of the generated videos while showing excellent alignment with text prompts.

918

919

920

*"A shark swimming in clear Carribean ocean, 2k, high quality"**"A squirrel eating a burger, high quality"**AnimateDiff*

926

*AnimateDiff + I4VGEN - pNI-vSDS = 0.6, p_{re} = 1.0*

931

*AnimateDiff + I4VGEN - pNI-vSDS = 1.0, p_{re} = 1.0*

937

*AnimateDiff + I4VGEN - pNI-vSDS = 0.6, p_{re} = 0.6*

942

Figure 9: **Impact of hyperparameters.** For different texts, the optimal parameter settings are different, and the sensitivity to parameters also varies. However, they all significantly outperform the baseline. In this paper, we provide an empirical setting that is mild for most cases, serving as a performance lower bound for I4VGEN. I4VGEN supports fine-tuning parameters on a per-example, achieving higher-quality videos.

949

950

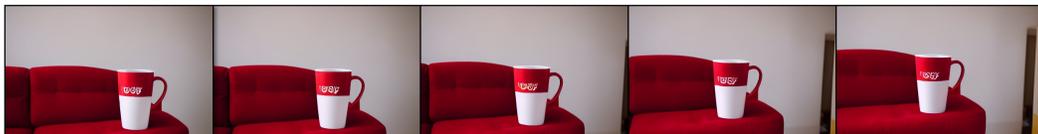
951

952

953

Prompt: "A cat and a dog reading books on the street, 4k, high resolution"

961

Prompt: "A red cup and a white sofa"

966

Figure 10: **Failure cases.** I4VGEN is designed to fully unleash the potential of existing video diffusion models, but it still cannot synthesize high-quality videos that are out of the distribution. For example, the text marked in red.

970

971

D ADDITIONAL EXPERIMENTAL RESULTS

D.1 ADDITIONAL VISUAL RESULTS ON DYNAMICRAFTER

We integrate I4VGEN into DynamiCrafter (Xing et al., 2024), which exhibits state-of-the-art performance on the VBench Image-to-Video Leaderboard. As shown in Fig. 11, the beginnings and endings of videos generated by DynamiCrafter suffer from low quality. For example, the front part of the face video generated by DynamiCrafter exhibits serious artifacts, and the face in the latter part are deformed. Our method alleviates these issues, which demonstrates that I4VGEN can significantly improve the quality of videos synthesized by DynamiCrafter.

D.2 ADDITIONAL VISUAL RESULTS ON SPARSECTRL

We conduct experiments on action instructions. As shown in Fig. 12, we explore two prompt-based motion enhancement strategies:

- By providing static descriptions in negative prompt, the dynamic intensity of the synthesized videos can be further enhanced.
- By providing specific action instruction in the prompt, such as “waving its hands”, the synthesized video accurately renders this action.

These findings indicate that I4VGEN does not compromise the dynamic nature of the synthesized videos but rather depicts more reasonable and accurate motion.

D.3 ADDITIONAL VISUAL RESULTS USING FLUX

We provide the visual results of I4VGEN adapted on FLUX in the Fig. 13. Despite the detailed and realistic images synthesized by FLUX, AnimateDiff + I4VGEN is still constrained by the video baseline, *i.e.*, AnimateDiff, in rendering image details and is unable to synthesize realistic videos. Evidently, the distribution of images synthesized by FLUX exceeds what AnimateDiff can handle, which relies on SD 1.5. However, the layout and composition information of images synthesized by FLUX still provide strong support for video synthesis, resulting in promising outcomes.

D.4 ADDITIONAL VISUAL RESULTS ON ANIMATEDIFF

We provide the visual results on AnimateDiff using the `Realistic Vision V5.1 LoRA` in the Fig. 14. Our method still significantly improves the quality of the generated videos while showing excellent temporal consistency.

D.5 ADDITIONAL INTERMEDIATE RESULTS VISUALIZATION

We provide additional intermediate results in the Fig. 17.

D.6 ADDITIONAL VISUAL RESULTS ON LARGE CAMERA POSE CHANGE

We provide more visual results involving significant changes in camera poses in Fig. 16, which demonstrate that our method can handle this scenario, improving the temporal consistency and smoothness of the synthesized videos.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
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
1076
1077
1078
1079

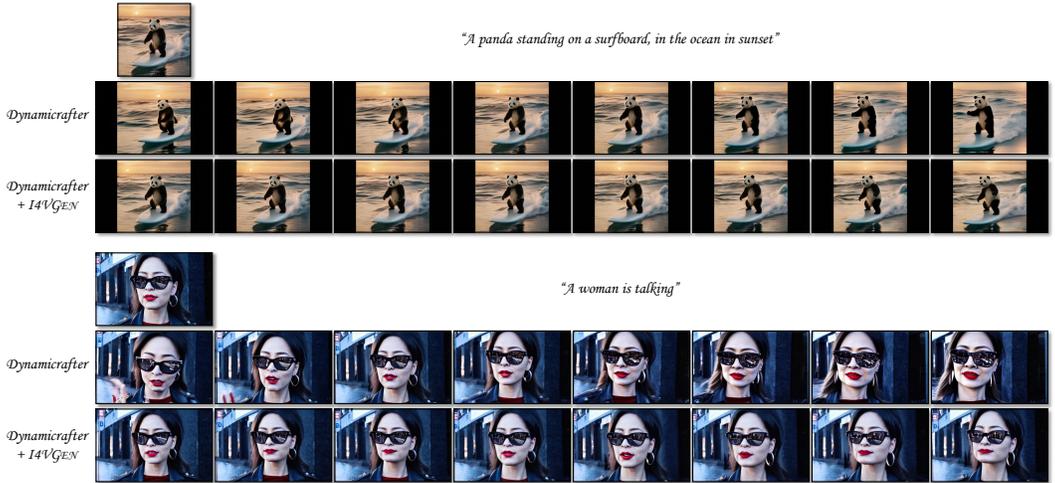


Figure 11: **Additional visual results on DynamiCrafter.** We provide the videos in mp4 format in the supplementary material for better viewing.

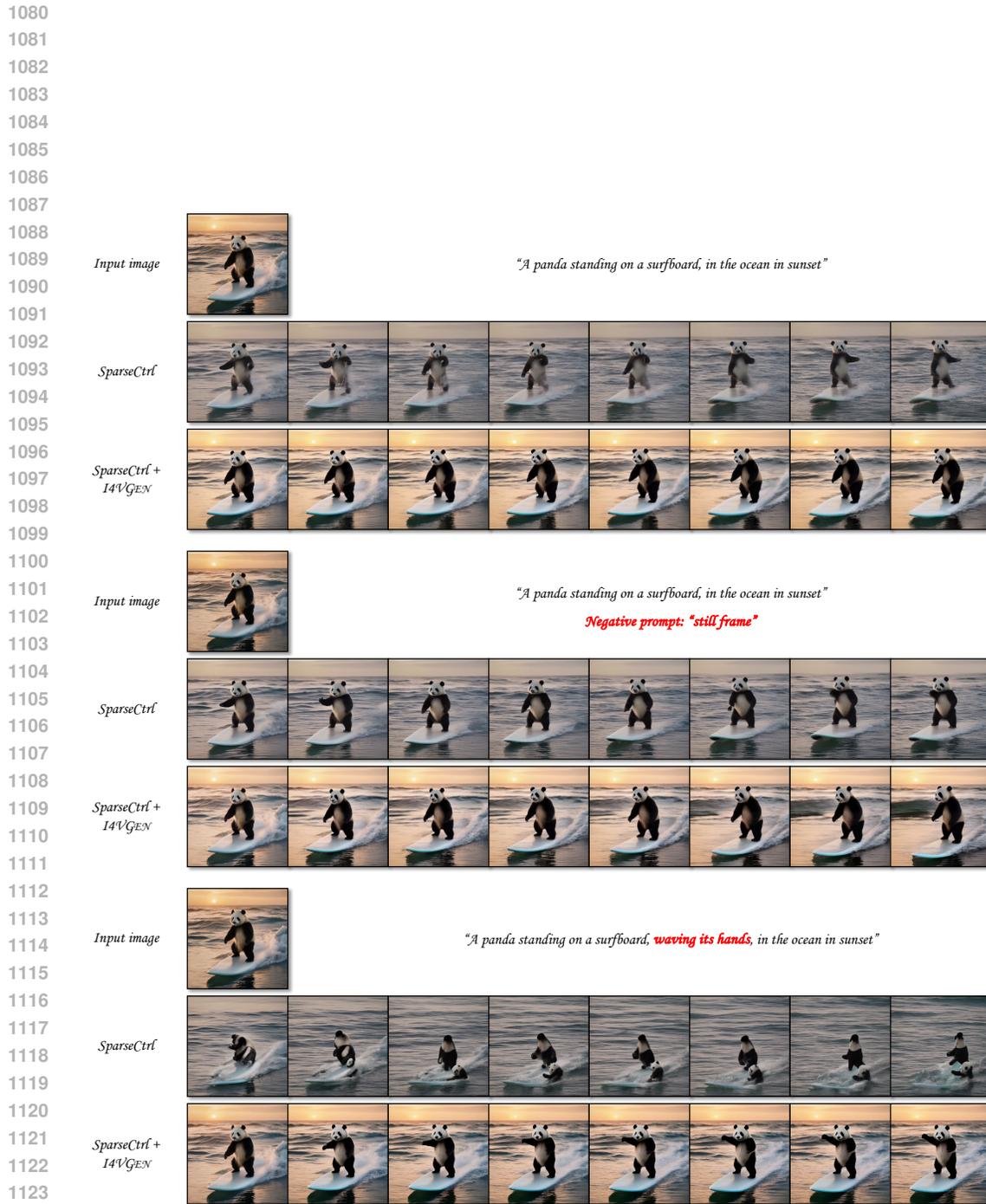
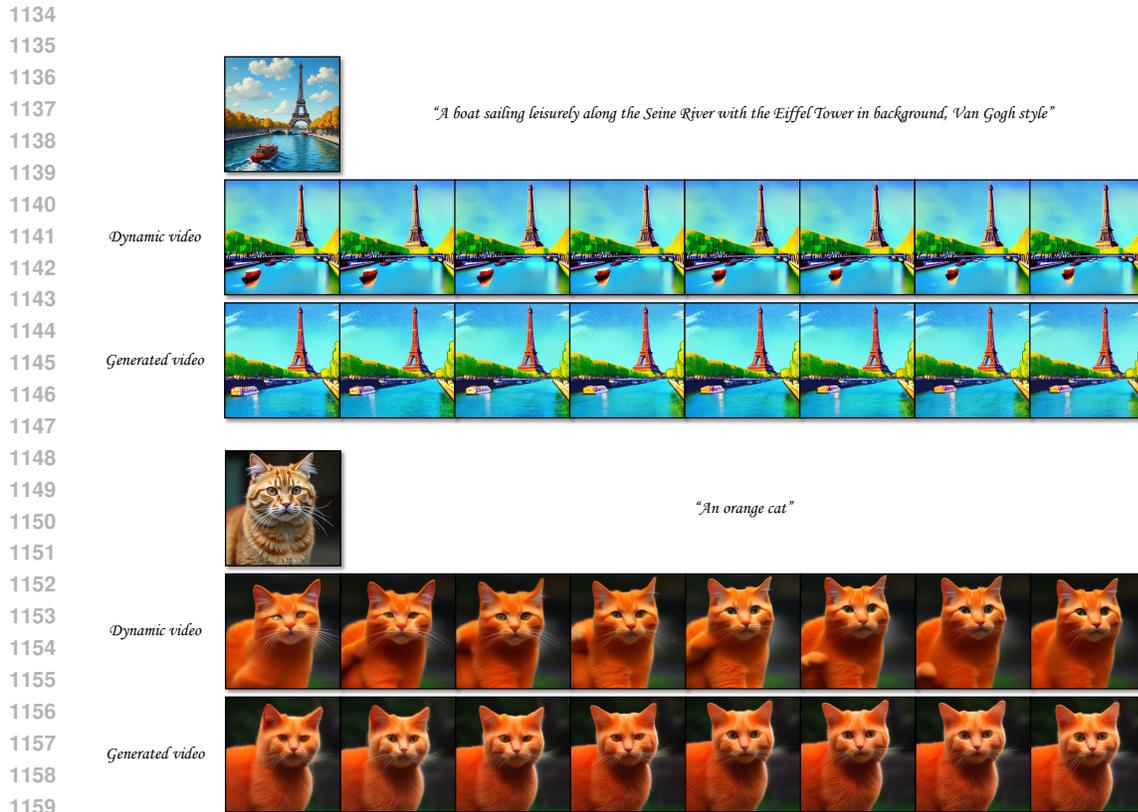
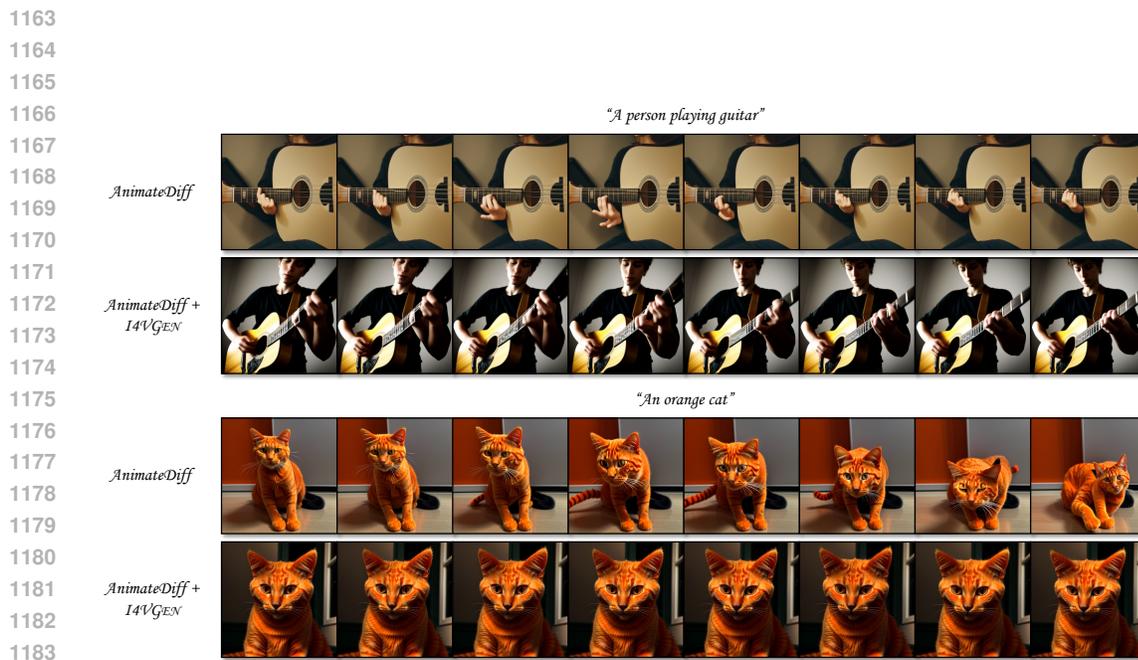


Figure 12: **Additional visual results on SparseCtrl.** We provide the videos in mp4 format in the supplementary material for better viewing.

1125
1126
1127
1128
1129
1130
1131
1132
1133



1161 **Figure 13: Adaptation on image generated by FLUX.** We provide the videos in mp4 format in the
1162 supplementary material for better viewing.



1185 **Figure 14: Additional visual results on AnimateDiff using LoRA.** We provide the videos in mp4
1186 format in the supplementary material for better viewing.

1187



Figure 15: **Additional intermediate results visualization.** We provide the videos in mp4 format in the supplementary material for better viewing.

1242
 1243
 1244
 1245
 1246
 1247
 1248
 1249
 1250
 1251
 1252
 1253
 1254
 1255
 1256
 1257
 1258
 1259
 1260
 1261
 1262
 1263
 1264
 1265
 1266
 1267
 1268
 1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282
 1283
 1284
 1285
 1286
 1287
 1288
 1289
 1290
 1291
 1292
 1293
 1294
 1295

"A drone camera circles around a beautiful historic church built on a rocky outcropping along the Amalfi Coast, the view showcases historic and magnificent architectural details and tiered pathways and patios, waves are seen crashing against the rocks below as the view overlooks the horizon of the coastal waters and hilly landscapes of the Amalfi Coast Italy, several distant people are seen walking and enjoying vistas on patios of the dramatic ocean views, the warm glow of the afternoon sun creates a magical and romantic feeling to the scene, the view is stunning captured with beautiful photography."



"The camera follows behind a white vintage SUV with a black roof rack as it speeds up a steep dirt road surrounded by pine trees on a steep mountain slope, dust kicks up from it's tires, the sunlight shines on the SUV as it speeds along the dirt road, casting a warm glow over the scene. The dirt road curves gently into the distance, with no other cars or vehicles in sight. The trees on either side of the road are redwoods, with patches of greenery scattered throughout. The car is seen from the rear following the curve with ease, making it seem as if it is on a rugged drive through the rugged terrain."



Figure 16: **Additional visual results on large camera pose change.** We provide the videos in mp4 format in the supplementary material for better viewing.



Figure 17: **Impact of p_{re} .** We provide the videos in mp4 format in the supplementary material for better viewing.