# FreeInit: Bridging Initialization Gap in Video Diffusion Models

Tianxing Wu    Chenyang Si    Yuming Jiang    Ziqi Huang    Ziwei Liu✉

S-Lab, Nanyang Technological University

{tianxing001, chenyang.si, yuming002, ziqi002, ziwei.liu}@ntu.edu.sg

https://github.com/TianxingWu/FreeInit

Figure 1. **FreeInit for Video Generation**. we propose *FreeInit*, a concise yet effective method to significantly improve temporal consistency of videos generated by diffusion models. FreeInit requires no additional training and introduces no learnable parameters, and can be easily incorporated into arbitrary video diffusion models at inference time.

## Abstract

*Though diffusion-based video generation has witnessed rapid progress, the inference results of existing models still exhibit unsatisfactory temporal consistency and unnatural dynamics. In this paper, we delve deep into the noise initialization of video diffusion models, and discover an implicit training-inference gap that attributes to the unsatisfactory inference quality. Our key findings are: 1) the spatial-temporal frequency distribution of the initial latent at inference is intrinsically different from that for training, and 2) the denoising process is significantly influenced by the low-frequency components of the initial noise. Moti-vated by these observations, we propose a concise yet effective inference sampling strategy, FreeInit, which significantly improves temporal consistency of videos generated by diffusion models. Through iteratively refining the spatial-temporal low-frequency components of the initial latent during inference, FreeInit is able to compensate the initialization gap between training and inference, thus effectively improving the subject appearance and temporal consistency of generation results. Extensive experiments demonstrate that FreeInit consistently enhances the generation results of various text-to-video generation models without additional training.*

## 1. Introduction

Recently, diffusion models have demonstrated impressive generative capabilities in text-to-image generation [31, 32, 35]. These advancements have attracted substantial attention, highlighting the potential of generative models to create diverse and realistic images based on textual descriptions. In light of these achievements, researchers are now exploring the application of diffusion models in text-to-video (T2V) generation [1, 4, 11, 13, 15, 36, 45, 46, 51], with the goal of synthesizing visually appealing and contextually coherent videos from textual descriptions. Most of these video diffusion models are built upon powerful pretrained image diffusion models, *e.g.*, Stable Diffusion (SD) [32]. Through the incorporation of temporal layers and fine-tuning on extensive video datasets, these models are capable of generating video clips that align with the given text prompts.

Similar to image-based diffusion models, when training video diffusion models, Gaussian Noise is gradually added to the input video at the diffusion process, aiming at corrupting it into noise. Then at the denoising process, the diffusion model learns to predict noise and reconstruct the original clean video from its noisy states. Ultimately at inference stage, this model is tasked to synthesize videos by iteratively denoising starting from pure Gaussian noise.

However, there exists a gap between the corrupted latents at training and the Gaussian initial noise at inference. The diffusion process does not fully corrupt the clean latents into pure Gaussian Noise. In Figure 2, we visualize the frames decoded from the noisy latents at multiple diffusion steps, and apply spatio-temporal frequency decomposition to analyze the corruption on different frequency bands. Notably, the low-frequency components are corrupted at a much lower speed compared to its high-frequency counterparts. As a result, the noisy latents at the last timestep ($t$=1000) still contain considerable low-frequency information from the clean input video. This eventually leads to an implicit gap at inference: on one hand, since the clean input frames are temporally correlated in nature, their noisy latents at training would also be temporally correlated, especially in their low-frequency band, which clearly differs from the i.i.d Gaussian initial noise at inference. On the other hand, the initial noise at inference, especially its low-frequency component, can substantially affect the generation quality, as revealed in our observations in Figure 5, 6. Hence, when applying the diffusion models trained with the correlated initial noises to non-correlated Gaussian initial noise at inference, the performance deteriorates.

Motivated by these observations, we propose a novel inference-time sampling method, denoted as *FreeInit*, to bridge the initialization gap between training and inference without any additional training or fine-tuning. Specifically, during the inference process, we first initialize independent

Gaussian noise, which then undergoes the DDIM denoising process to yield a clean video latent. Subsequently, we obtain the noisy version of the generated video latent through the forward diffusion process. Since the noisy latents are obtained from denoised clean latents, the low-frequency components of these noisy latents have improved temporal consistency. With these noisy latents, we proceed to reinitialize the noise by combining the low-frequency components of these noisy latents with the high-frequency components of random Gaussian noise. Finally, this reinitialized noise serves as the starting point for new DDIM sampling, facilitating the generation of frames with enhanced temporal consistency and visual appearance. We iterate the aforementioned reinitialization several times to generate the final videos.

Extensive experiments across a range of diverse evaluation prompts demonstrate the steady enhancement brought about by FreeInit for various text-to-video generation models. As illustrated in Figure 1, FreeInit plays a significant role in improving temporal consistency and the visual appearance of generated frames. This method can be readily applied during inference without the need for parameter tuning. Furthermore, to achieve superior generation quality, the frequency filter can be conveniently adjusted for each customized base model. We summarize our contributions as follows:

- We systematically investigate the noise initialization of video diffusion models, and identify an implicit training-inference gap in frequency domain that contributes to the inference quality drop.
- We propose a concise yet effective sampling strategy, referred to as *FreeInit*, which iteratively refines the initial noise without the need for additional training or fine-tuning.
- Extensive quantitative and qualitative experiments demonstrate that FreeInit can be effectively applied to various text-to-video models. It consistently improves the inference quality of generated videos.

## 2. Related Work

**Video Generative Models.** There are mainly three types of video generation models, namely GAN-based [9], transformer-based [43], and diffusion-based [14]. StyleGAN-V [37], MoCoGAN-HD [42], and [2] utilize the powerful StyleGAN [20–22] to generate videos. Transformer-based models [17, 19, 44, 47, 48] such as Phenaki [44], CogVideo [17], and NÜWA [48] encode videos as visual tokens and train transformer models to auto-regressively generate the visual tokens. Recently, diffusion models [5, 14, 38, 40] have made remarkable progress in text-to-image generation [26, 28, 32, 35], and have enabled a line of works that extends these

(a) Frames Decoded from $z_t$

*DDPM Forward*

(b) Frames Decoded from Low Frequency of $z_t$

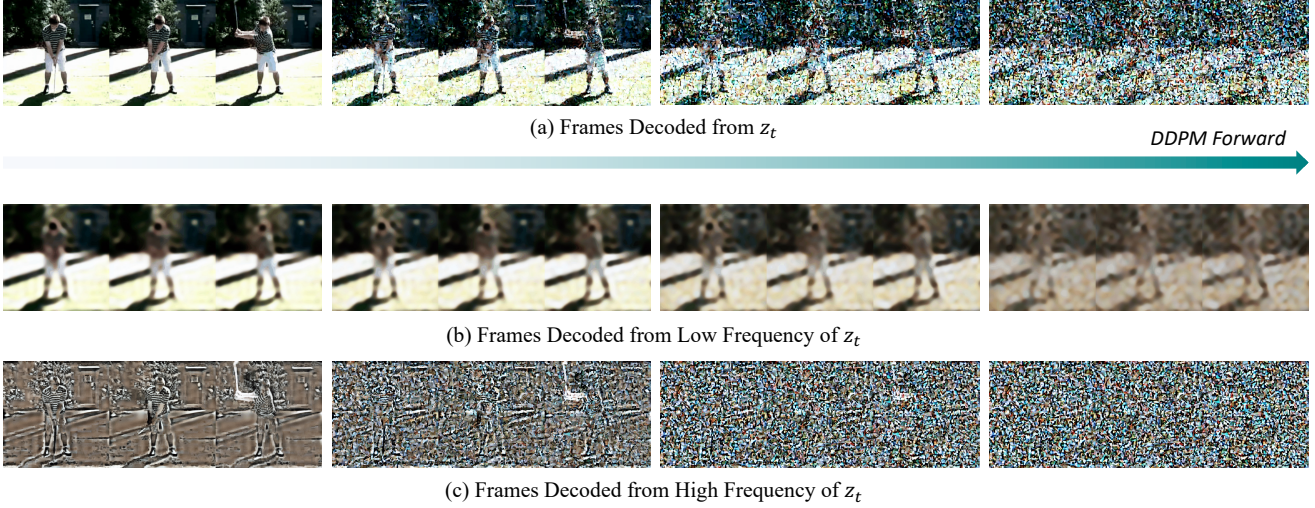(c) Frames Decoded from High Frequency of $z_t$

Figure 2. **Visualization of Decoded Noisy Latent from Different Frequency Bands at Training Stage.** (a) Video frames decoded from the entire frequency band of the noisy latent $z_t$ in DDPM Forward Process. (b) We perform spatio-temporal frequency decomposition on $z_t$, and visualize the frames decoded from the low-frequency components of $z_t$. It is evident that the diffusion process cannot fully corrupt the semantics, leaving substantial spatio-temporal correlations in the low-frequency components. (c) We visualize frames decoded from the high-frequency components of $z_t$, and observe that each frame degenerates rapidly with the noising process.

pre-trained diffusion models towards text-to-video generation [1, 8, 11–13, 15, 16, 23, 25, 36, 45, 46, 50–52]. In this work, our method is built on top of diffusion-based text-to-video methods. We propose to iteratively refine the initial noise to improve temporal consistency of pre-trained video diffusion models. We demonstrate the effectiveness of our method on various diffusion models, including VideoCrafter, AnimateDiff, and ModelScope. VideoCrafter [13] employs the pre-trained text-to-image model Stable Diffusion [32] and incorporates newly initialized temporal layers to enable video generation. AnimateDiff [11] trains motion modeling modules and inserts them into personalized text-to-image diffusion models to achieve animated videos of customized concepts (*e.g.*, characters). ModelScope [25, 45] learns video generation in a noise decomposition fashion, where a base noise is shared among all frames and a residual noise is predicted in a per-frame manner.

**Noise in Diffusion Models.** Only a few previous works have mentioned the limitations of the noise schedule of current diffusion models. In the image domain, [24] points out common diffusion noise schedules cannot fully corrupt information in natural images, limiting the model to only generate images with medium brightness. A rescaled training schedule is then proposed to alleviate this problem through fine-tuning. Recently, [6] makes further discussions on the signal leakage issue, and propose to explicitly model the signal leakage for better inference noise distribution, which produces images with more diverse brightness and colours. In the video domain, PYoCo [8] carefully designs the progressive video noise prior to achieve a better video gener-

ation performance. Similar to [24], PYoCo also focuses on the noise schedule at training stage and requires massive fine-tuning on video datasets. In contrast, we focus on the initial noise at inference stage and proposes a concise inference-time sampling strategy that bridges the training-inference discrepancy with no fine-tuning required. Some recent works [10, 29] also pay attention to the inference initial noise, but aiming at generating long videos. We instead focus on improving inference quality, and further design specific frequency-domain-based operations to modulate different frequency components of the initial noise.

## 3. Preliminaries and Observations

### 3.1. Preliminaries

Similar to image diffusion models, video diffusion models also involve a diffusion process and a denoising process, and operate in the latent space of an autoencoder. The diffusion process includes a sequence of $T$ steps. At each step $t$, Gaussian noise is incrementally added to the video latent $z_0$, following a predefined variance schedule $\beta_1, \ldots, \beta_T$:

$$q(z_{1:T}|z_0) = \prod_{t=1}^{T} q(z_t|z_{t-1}), \tag{1}$$

$$q(z_t|z_{t-1}) = \mathcal{N}(z_t; \sqrt{1-\beta_t}z_{t-1}, \beta_t\mathbf{I}). \tag{2}$$

Let $\alpha_t = 1 - \beta_t, \overline{\alpha_t} = \prod_{s=1}^{t} \alpha_s$:

$$q(z_t|z_0) = \mathcal{N}(z_t; \sqrt{\overline{\alpha}_t}z_0, (1-\overline{\alpha}_t)\mathbf{I}). \tag{3}$$
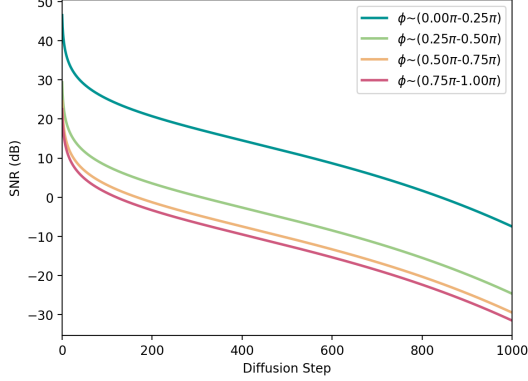
Figure 3. **Variation of SNR in Diffusion.** The figure shows the SNR pattern of the latent code during the diffusion process. Each curve corresponds to a spatio-temporal frequency band of $z_t$, whose frequency range is normalized to 0 to 1, where 0 is the lowest frequency and 1 is the highest frequency.
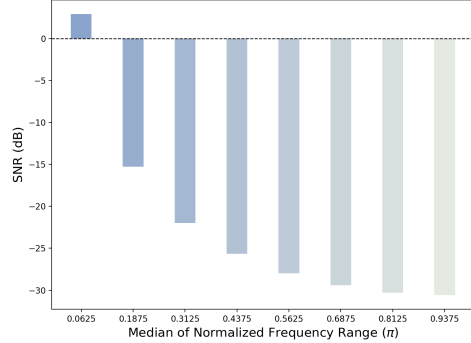


Figure 4. **Frequency Distribution of the SNR in Initial Noise.** We randomly sample 500 videos from UCF-101 and apply diffusion process with Stable Diffusion's $\beta$ scheme, obtaining 500 training initial noises. Then the average spatio-temporal frequency distribution of SNR is computed. Surprisingly, the SNR in low-frequency components is larger than 0 dB, indicating a severe information leak.

As a result, the noisy latent $z_t$ at each timestep $t$ can be directly sampled as:

$$z_t = \sqrt{\overline{\alpha}_t} z_0 + \sqrt{1 - \overline{\alpha}_t} \epsilon, \qquad (4)$$

where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is a Gaussian white noise with the same shape as $z_t$.

In the reverse process, the network learns to recover the clean latent $z_0$ by iterative denoising with U-Net [33], starting from the initial noise $z_T$:

$$p_\theta(z_{0:T}) = p(z_T) \prod_{t=1}^{T} p_\theta(z_{t-1}|z_t), \qquad (5)$$

$$p_\theta(z_{t-1}|z_t) = \mathcal{N}(z_{t-1}; \mu_\theta(z_t, t), \Sigma_\theta(z_t, t)), \qquad (6)$$

where $\mu_\theta$ and $\Sigma_\theta$ are predicted by the denoising U-Net $\epsilon_\theta$.

During inference, an initial latent $\hat{z}_T$ is first initialized, typically as a Gaussian noise sampled from a normal distribution:

$$\hat{z}_T = \epsilon' \sim \mathcal{N}(0, \mathbf{I}). \qquad (7)$$

Then the trained network $\epsilon_\theta$ is used to iteratively denoise the noisy latent to a clean latent $\hat{z}_0$ through DDIM sampling [39], which is then decoded with decoder $\mathcal{D}$ to obtain video frames $\hat{x}_0$.

### 3.2. The Initialization Gap

**Signal-to-Noise Ratio During Forward Diffusion Process.** To better understand the forward diffusion process, we conduct an investigation utilizing the Signal-to-Noise Ratio (SNR) as a metric to analyze progress of information corruption. Figure 3 shows the SNR measurements of the noisy latent $z_t$ (as defined in Eqn. 4) on the UCF101 dataset. Our

observations reveal a consistent pattern wherein the low-frequency components exhibit a significantly higher SNR when compared to the high-frequency components. This finding underscores the inability of the diffusion process to fully corrupt the information within the spatio-temporal low-frequency components of the video latent. Consequently, this inadequacy results in an implicit leakage of signal information at training stage.

Furthermore, as seen in Figure 4, it is noteworthy that the SNR of low-frequency components consistently remains above 0 dB, even at $t$=1000. This observation demonstrates the existence of a noticeable gap between the training and inference processes. Specifically, the noise introduced during training is insufficient to completely corrupt latent information, and the low-frequency components of the initial noise (*i.e.*, latent at $t$=1000) persistently contain spatio-temporal correlations. However, during the inference process, the video generation model is tasked with generating coherent frames from independent Gaussian noise. This presents a considerable challenge for the denoising network, as its initial noise lacks spatio-temporal correlations at inference. For instance, as illustrated in Figure 6, the "biking" video generated from Gaussian noise exhibits unsatisfactory temporal consistency. In contrast, when utilizing the noise latent obtained through the forward diffusion process from real videos as initial noise, the generated "biking" video showcases improved temporal consistency. This observation indicates the presence of an initialization gap between training and inference, which significantly affects the quality of the generated content.

**Influence of Initial Low-frequency Components.** Considering the SNR gaps of initial noise between training and inference, we further investigate the influence of the low-

Figure 5. **Role of Initial Low-Frequency Components.** We progressively remove the spatio-temporal high frequency components in the initial noise maps $z_T$, and visualize the roles of the remaining low-frequency components. We observe that even if the majority (*e.g.*, "Remove 80%") of high frequencies are removed, the generated results still remain largely similar to the "Full $z_T$" videos, indicating that overall distribution of the generated results is determined by the low-frequency components of the initial noise $z_T$.



(a) Inference with Gaussian Noise



(b) Inference with Training Initial Noise

Figure 6. **Initialization Gap.** (a) With randomly initialized Gaussian noise for different frames, the sampled video exhibits inconsistency among frames. (b) When we start from noisy latent obtained from the diffusion process from real videos, the generated video is temporally consistent. This is because the initial noise is aligned with training stage and it contains correlated information among different frames in nature.

frequency components of initial noise. Through the forward diffusion process from real videos, we can obtain its noise latent $z_T$. We gradually remove its high-frequency components and mix it with a random Gaussian Noise $\epsilon'$, only keeping its low-frequency components in the initial noise for inference. As shown in Figure 5, it is evident that variations in high-frequency have a negligible impact on the overall generation results. Remarkably, the overall distribution of the generated outcomes remains stable, even when

employing only 20% of the original initial latent information from the low-frequency band. When all information is removed, the denoising process is initialized with pure Gaussian noise, which leads to relatively poor generation results. This observation highlights two key conclusions: 1) the low-frequency components of the initial noisy latent plays an essential role in inference, and 2) the quality of the low-frequency components is crucial for the generation quality. These conclusions motivate us to propose a concise yet effective strategy for enhancing the inference quality of video diffusion models.

## 4. FreeInit

Motivated by the above analysis, we propose a method for relieving this gap by progressively refining the low-frequency components of the initial noise using the inherent power of the diffusion model. We refer to this method as **FreeInit**, which substantially improves the generation quality without additional training or fine-tuning.

As illustrated in Figure 7, during the inference process, independent Gaussian noise $\epsilon$ is first initialized, which then undergoes the DDIM sampling process to yield a primary denoised latent $z_0$. Subsequently, we obtain the noise latent $z_T$ of the generated latent $z_0$ through the forward DDPM diffusion process, *i.e.*, adding noise to diffuse $z_0$ to $z_T$. Hence, the low-frequency components of noise latent $z_T$ have a better spatio-temporal correlation compared to $\epsilon$. It is worth noting that, during this forward diffusion process, we
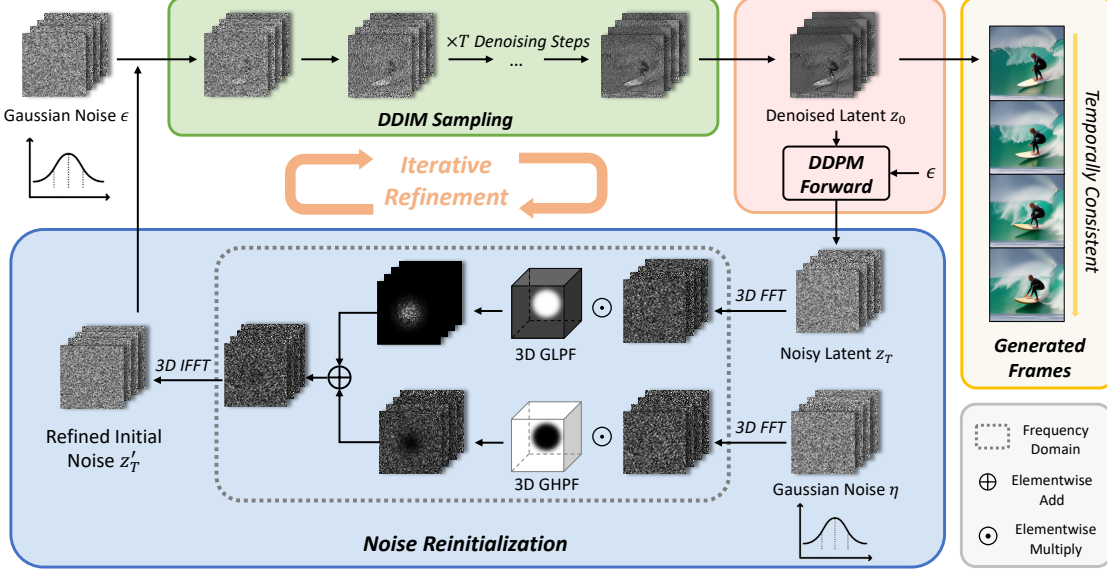
Figure 7. **Framework of FreeInit**. FreeInit refines the low-frequency components of the inference initial noise in an iterative manner. During inference, a Gaussian Noise is first initialized and goes through the standard DDIM sampling process. The resulting denoised latent $z_0$ is then diffused using the original Gaussian Noise $\epsilon$, through a DDPM forward process. With the obtained noisy latent $z_T$ which contains richer low-frequency information, a noise reinitialization process is further performed: $z_T$ is firstly transformed into frequency domain through 3D FFT, then its spatio-temporal low-frequency components are fused with the high-frequency from a randomly sampled Gaussian noise $\eta$, bringing flexibility for refinement in the higher frequency band. After transforming back to the time domain, the refined initial noise $z_T'$ is used as the initial noise for the next iteration.

have observed that adding randomly sample Gaussian noise could introduce significant uncertainty in the mid-frequency band, compromising the spatio-temporal correlation. Consequently, we opt to utilize the same original noise $\epsilon$ used in DDIM sampling when diffusing $z_0$ to $z_T$. The mathematical representation of this process is as follows:

$$z_T = \sqrt{\overline{\alpha}_T} z_0 + \sqrt{1 - \overline{\alpha}_T} \epsilon \qquad (8)$$
$$= \sqrt{\overline{\alpha}_T}(DDIM_{sample}(\epsilon)) + \sqrt{1 - \overline{\alpha}_T} \epsilon,$$

where $\overline{\alpha}_T$ is aligned with the $\beta$ schedule used at training, *e.g.*, Stable Diffusion schedule.

**Noise Reinitialization.** To maintain alignment with the SNR distribution at training stage, we propose a noise reinitialization strategy, which is essential for the improvement of temporal consistency. Specifically, we employ a spatio-temporal frequency filter to combine the low-frequency components of the noise latent $z_T$ with the high-frequency components of a random Gaussian noise $\eta$, resulting in the reinitialized noisy latent $z_T'$. This approach allows us to preserve essential information contained in $z_T$ while introducing sufficient randomness in the high-frequency components to enhance visual details, complementing the improved low-frequency components. The mathematical op-

erations are performed as follows:

$$\mathcal{F}_{z_T}^L = \mathcal{FFT}_{3D}(z_T) \odot \mathcal{H}, \qquad (9)$$
$$\mathcal{F}_{\eta}^H = \mathcal{FFT}_{3D}(\eta) \odot (1 - \mathcal{H}), \qquad (10)$$
$$z_T' = \mathcal{IFFT}_{3D}(\mathcal{F}_{z_T}^L + \mathcal{F}_{\eta}^H), \qquad (11)$$

where $\mathcal{FFT}_{3D}$ is the Fast Fourier Transformation operated on both spatial and temporal dimensions, $\mathcal{IFFT}_{3D}$ is the Inverse Fast Fourier Transformation that maps back the blended representation $F_{z_T'}$ from the frequency domain back. $\mathcal{H} \in^{4,N,H',W'}$ is the spatial-temporal Low Pass Filter (LPF), which is a tensor of the same shape as the latent.

Finally, this reinitialized noise $z_T'$ serves as the starting point for new DDIM sampling, facilitating the generation of frames with enhanced temporal consistency and visual appearance.

**Iterative Refinement of Initial Noise.** It is important to note that the aforementioned operations can be iteratively applied. At each iteration, the latent code undergoes improvements in spatio-temporal consistency by preserving low-frequency information for noise reinitialization. Simultaneously, it gains flexibility in the high-frequency domain through reinitialization, resulting in an improved initial noise for the subsequent iteration. In this iterative manner, the quality of the initial noise is progressively refined, effectively bridging the gap between training and inference.

Table 1. **Quantitative Comparisons on UCF-101 and MSR-VTT.** FreeInit significantly improves the temporal consistency.

| Method | DINO ↑ | |
| --- | --- | --- |
| | UCF-101 | MSR-VTT |
| AnimateDiff [11] | 85.24 | 83.24 |
| AnimateDiff+FreeInit | **92.01** | **91.86** |
| ModelScope [45] | 88.16 | 88.95 |
| ModelScope+FreeInit | **91.11** | **93.28** |
| VideoCrafter [4] | 85.62 | 84.68 |
| VideoCrafter+FreeInit | **89.27** | **88.72** |

Ultimately, this iterative process contributes to the overall enhancement of generation quality.

## 5. Experiments

### 5.1. Implementation Details

To evaluate the effectiveness and generalization of our proposed method, we apply the FreeInit strategy to three publicly available diffusion based text-to-video models: AnimateDiff [11], ModelScope [45] and VideoCrafter [4]. Following [8, 36], we evaluate the inference performance with prompts from UCF-101 [41] and MSR-VTT [49] dataset. For UCF-101, we use the same prompt list as proposed in [8]. For MSR-VTT, we randomly sample 100 prompts from the test set for evaluation. We also incorporate diverse prompts from [18] for qualitative evaluations.

During inference, the parameters of frequency filter for each model are kept the same for fair comparison. Specifically, we use a Gaussian Low Pass Filter (GLPF) $\mathcal{H}_{\mathcal{G}}$ with a normalized spatio-temporal stop frequency of $D_0 = 0.25$. For each prompt, we first adopt the default inference settings of each model for a single inference pass, then apply 4 extra FreeInit iterations and evaluate the progress of generation quality. All FreeInit metrics in Quantitative Comparisons are computed at the $4^{th}$ iteration.

**Evaluation Metrics.** To measure the temporal consistency of the generated video, we compute frame-wise similarity between the first frame and all succeeding $N - 1$ frames. Notably, one typical failure case in current video diffusion models is semantically close but visually inconsistent generation result. For example in Figure 6 (a), all frames are semantically aligned ("biking"), but the appearance of the subject and background exhibits unsatisfactory consistency. Consequently, semantic-based features like CLIP [30] is not appropriate for evaluating the visual temporal consistency in video generation. Following previous studies [34], we utilize ViT-S/16 DINO [3, 27] to measure the visual similarities, denoted as the **DINO** metric. The metrics is averaged on all frames.

### 5.2. Quantitative Comparisons

The quantitative results on UCF-101 and MSR-VTT are reported in Table 1. According to the metrics, FreeInit signifi-

Table 2. **Ablation Study on Noise Reinitialization (NR).** Removing NR or changing Gaussian Low Pass Filter (GLPF) to Ideal Low Pass Filter (ILPF) leads to non-optimal results. *ModelName\** refers to *Model+FreeInit*.

| Method | DINO ↑ | |
| --- | --- | --- |
| | UCF-101 | MSR-VTT |
| AnimateDiff* w/o NR | 86.77 | 85.18 |
| AnimateDiff* w/ NR-ILPF | 87.53 | 86.17 |
| AnimateDiff* w/ NR-GLPF | **92.01** | **91.86** |
| ModelScope* w/o NR | 88.20 | 90.90 |
| ModelScope* w/ NR-ILPF | 89.04 | 90.93 |
| ModelScope* w/ NR-GLPF | **91.11** | **93.28** |
| VideoCrafter* w/o NR | 86.09 | 87.11 |
| VideoCrafter* w/ NR-ILPF | 87.53 | 88.01 |
| VideoCrafter* w/ NR-GLPF | **89.27** | **89.33** |

cantly improves the temporal consistency of all base models on both prompt sets, by a large margin from 2.92 to 8.62. We also conduct a User Study to evaluate the results through Temporal Consistency, Text Alignment and Overall Quality, which can be refer to the Supplementary File.

### 5.3. Qualitative Comparisons

Qualitative comparisons are shown in Figure 8. Our proposed FreeInit significantly improves the temporal consistency as well as visual quality. For example, with text prompt "a musician playing the flute", performing FreeInit effectively fix the temporally unstable artifacts exhibited in vanilla AnimateDiff. More qualitative results are listed in the Supplementary File.

### 5.4. Ablation Study

We evaluate the influence of different design choices and parameters of FreeInit with ablation studies quantitatively and qualitatively.

**Influence of Noise Reinitialization and Filter Selection.**
To evaluate the importance of Noise Reinitialization in the frequency domain and the choice of filter, we run two FreeInit variants on both datasets with all three base models. Firstly, Noise Reinitialization is totally skipped, *i.e.*, the noisy latent $z_T$ after DDPM Forward Pass is directly used as initial noise for sampling. Secondly, the frequency filter used for Noise Reinitialization is changed from GLPF to ILPF, with the same stop frequency 0.25. The metrics in Table 2 clearly demonstrate that Noise Reinitialization is crucial for improving temporal consistency. Also, replacing the soft Gaussian filter GLPF with the hard Ideal filter ILPF will also lead to a performance drop in the metrics, which reveals the importance of also introducing moderate randomness into mid-frequency and low-frequency components. More detailed discussions can be referred to the Supplementary File.

**Influence of Iteration Steps.** We show the influence of FreeInit iteration step number in Figure 9. It can be observed that the temporal consistency consistently increases

Figure 8. **Qualitative Comparisons.** We apply FreeInit to different base models and inference with diverse text prompts. By refining the initial noise at inference time, FreeInit significantly improves the temporal consistency and the subject appearance of the generated videos.
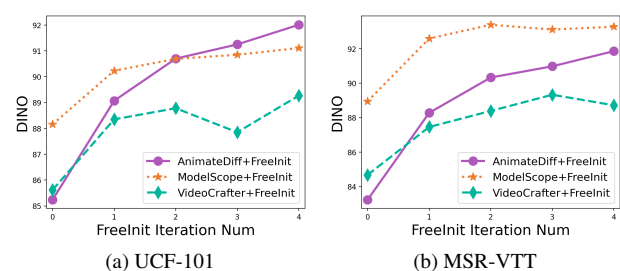


(a) UCF-101      (b) MSR-VTT

Figure 9. **Ablation Study on Iteration Number.** We report the DINO scores under different FreeInit iteration numbers on (a) UCF-101 and (b) MSR-VTT. More iteration steps mostly leads to better temporal consistency, and the most significant improvement is observed at the $1^{st}$ iteration.



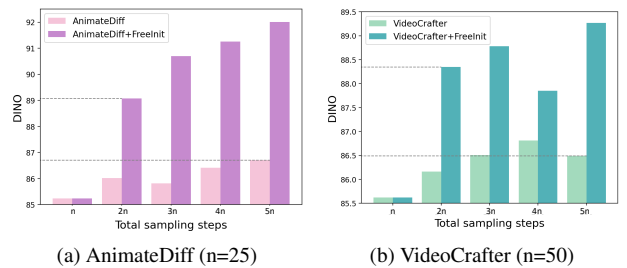(a) AnimateDiff (n=25)      (b) VideoCrafter (n=50)

Figure 10. **Comparison with Same Sampling Steps without FreeInit.** We analyze if increasing the DDIM sampling steps for baseline methods would help to improve the temporal consistency on UCF-101. For all base models, the Vanilla inference with 5n steps is inferior to incorporating FreeInit with 2n steps. This indicates that FreeInit is not equivalent to trivially increasing the DDIM sampling steps.

with the iteration step, thanks to the gradually refined initial noise. Notably, the largest temporal consistency improvement for each model comes from the 1st iteration, where FreeInit is applied for the first time. This is because at the 0-th iteration, the initial noise is non-correlated Gaussian noise, while at the 1st iteration, low-frequency information is injected into the noise for the first time, largely eliminating the gap between inference noise and training noise.

### 5.5. Further Discussion

**Comparison with Same Inference Step without FreeInit.**
Since FreeInit uses more than one DDIM sampling pass, it

is natural to ask if the quality improvement is due to the increased sampling steps. To answer this question, we compare FreeInit with the typical DDIM sampling strategy using the same total inference steps performed in a single iteration. As shown in Figure 10, trivially increasing the DDIM sampling steps only brings little improvement in temporal consistency. This further proves the importance of refining initial noise at inference time: a good beginning matters more than spending time struggling with a bad initial state.

**Limitations.** Since FreeInit is an iterative method, a natural drawback is the increased sampling time. This issue can be mitigated through a coarse-to-fine sampling strategy. We

explain the details of this strategy and more discussions on its limitations and potential negative societal impacts in the Supplementary File.

## 6. Conclusion

In this paper, we identify an implicit training-inference gap in the noise initialization of video diffusion models that causes degenerated inference quality: 1) the frequency distribution of the initial latent's SNR is different between training and inference; 2) the denoising process is significantly affected by the low-frequency components of initial noise. Based on these observations, we propose FreeInit, which improves temporal consistency through the iterative refinement of the spatial-temporal low-frequency component of the initial latent during inference. This narrows the initialization gap between training and inference. Extensive quantitative and qualitative experiments on various text-to-video models and text prompts demonstrate the effectiveness of our proposed FreeInit.

## References

[1] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *CVPR*, 2023. 2, 3

[2] Tim Brooks, Janne Hellsten, Miika Aittala, Ting-Chun Wang, Timo Aila, Jaakko Lehtinen, Ming-Yu Liu, Alexei Efros, and Tero Karras. Generating long videos of dynamic scenes. In *NeurIPS*, 2022. 2

[3] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021. 7

[4] 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*, 2023. 2, 7, 11, 12

[5] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat GANs on image synthesis. In *NeurIPS*, 2021. 2

[6] Martin Nicolas Everaert, Athanasios Fitsios, Marco Bocchio, Sami Arpa, Sabine Süsstrunk, and Radhakrishna Achanta. Exploiting the signal-leak bias in diffusion models. *arXiv preprint arXiv:2309.15842*, 2023. 3

[7] Hugging Face. Diffusers. https://huggingface.co/docs/diffusers/index. 12

[8] Songwei Ge, Seungjun Nah, Guilin Liu, Tyler Poon, Andrew Tao, Bryan Catanzaro, David Jacobs, Jia-Bin Huang, Ming-Yu Liu, and Yogesh Balaji. Preserve your own correlation: A noise prior for video diffusion models. In *ICCV*, 2023. 3, 7

[9] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, 2014. 2

[10] Jiaxi Gu, Shicong Wang, Haoyu Zhao, Tianyi Lu, Xing Zhang, Zuxuan Wu, Songcen Xu, Wei Zhang, Yu-Gang Jiang, and Hang Xu. Reuse and diffuse: Iterative denoising for text-to-video generation. *arXiv preprint arXiv:2309.03549*, 2023. 3

[11] 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. 2, 3, 7, 11, 12

[12] William Harvey, Saeid Naderiparizi, Vaden Masrani, Christian Weilbach, and Frank Wood. Flexible diffusion modeling of long videos. *arXiv preprint arXiv:2205.11495*, 2022.

[13] Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent video diffusion models for high-fidelity video generation with arbitrary lengths. *arXiv preprint arXiv:2211.13221*, 2022. 2, 3, 12

[14] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 2

[15] 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. 2, 3

[16] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *arXiv preprint arXiv:2204.03458*, 2022. 3

[17] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. CogVideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022. 2

[18] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, Yaohui Wang, Xinyuan Chen, Limin Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. VBench: Comprehensive benchmark suite for video generative models. *arXiv preprint arXiv:2311.17982*, 2023. 7

[19] Yuming Jiang, Shuai Yang, Tong Liang Koh, Wayne Wu, Chen Change Loy, and Ziwei Liu. Text2Performer: Text-driven human video generation. In *ICCV*, 2023. 2

[20] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019. 2

[21] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *CVPR*, 2020.

[22] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. In *NeurIPS*, 2021. 2

[23] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators. *arXiv preprint arXiv:2303.13439*, 2023. 3

[24] Shanchuan Lin, Bingchen Liu, Jiashi Li, and Xiao Yang. Common diffusion noise schedules and sample steps are flawed. *arXiv preprint arXiv:2305.08891*, 2023. 3

[25] Zhengxiong Luo, Dayou Chen, Yingya Zhang, Yan Huang, Liang Wang, Yujun Shen, Deli Zhao, Jingren Zhou, and Tieniu Tan. VideoFusion: Decomposed diffusion models for high-quality video generation. In *CVPR*, 2023. 3

[26] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. GLIDE: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 2

[27] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 7

[28] 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. 2

[29] Haonan Qiu, Menghan Xia, Yong Zhang, Yingqing He, Xintao Wang, Ying Shan, and Ziwei Liu. Freenoise: Tuning-free longer video diffusion via noise rescheduling. *arXiv preprint arXiv:2310.15169*, 2023. 3

[30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 7

[31] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with CLIP latents. *arXiv preprint arXiv:2204.06125*, 2022. 2

[32] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 2, 3

[33] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 4

[34] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *CVPR*, 2023. 7

[35] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. 2

[36] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022. 2, 3, 7

[37] Ivan Skorokhodov, Sergey Tulyakov, and Mohamed Elhoseiny. Stylegan-v: A continuous video generator with the price, image quality and perks of stylegan2. In *CVPR*, 2022. 2

[38] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, 2015. 2

[39] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021. 4

[40] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *ICLR*, 2021. 2

[41] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012. 7

[42] Yu Tian, Jian Ren, Menglei Chai, Kyle Olszewski, Xi Peng, Dimitris N Metaxas, and Sergey Tulyakov. A good image generator is what you need for high-resolution video synthesis. *arXiv preprint arXiv:2104.15069*, 2021. 2

[43] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 2

[44] Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kindermans, Hernan Moraldo, Han Zhang, Mohammad Taghi Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. Phenaki: Variable length video generation from open domain textual description. *arXiv preprint arXiv:2210.02399*, 2022. 2

[45] Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modelscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023. 2, 3, 7, 11, 12

[46] 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*, 2023. 2, 3

[47] Chenfei Wu, Lun Huang, Qianxi Zhang, Binyang Li, Lei Ji, Fan Yang, Guillermo Sapiro, and Nan Duan. Godiva: Generating open-domain videos from natural descriptions, 2021. 2

[48] Chenfei Wu, Jian Liang, Lei Ji, Fan Yang, Yuejian Fang, Daxin Jiang, and Nan Duan. Nüwa: Visual synthesis pretraining for neural visual world creation, 2021. 2

[49] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *CVPR*, 2016. 7

[50] David Junhao Zhang, Jay Zhangjie Wu, Jia-Wei Liu, Rui Zhao, Lingmin Ran, Yuchao Gu, Difei Gao, and Mike Zheng Shou. Show-1: Marrying pixel and latent diffusion models for text-to-video generation, 2023. 3

[51] Daquan Zhou, Weimin Wang, Hanshu Yan, Weiwei Lv, Yizhe Zhu, and Jiashi Feng. Magicvideo: Efficient video generation with latent diffusion models. *arXiv preprint arXiv:2211.11018*, 2022. 2

[52] Daquan Zhou, Weimin Wang, Hanshu Yan, Weiwei Lv, Yizhe Zhu, and Jiashi Feng. Magicvideo: Efficient video generation with latent diffusion models, 2023. 3

# FreeInit: Bridging Initialization Gap in Video Diffusion Models

## Supplementary Material

In this ***Supplementary File***, we first explain the SNR distribution computation in Section A, then tabulate the user study results in Section B. More detailed discussions on Noise Reinitialization, filter selection and iteration steps are provided in Section C. We list more implementation details in Section D. More qualitative comparisons are illustrated in Section E to visualize the performance of FreeInit. We further discuss some limitations of FreeInit and its possible social impact in Section F and Section G.

We also record a demo video and more visual comparisons in video format. Please see the video and project page for more visualizations.

## A. Signal-to-Noise Ratio Distribution

In this section, we explain how we derive the frequency SNR distribution of $z_t$ in manuscript Section 3.2.

Mathematically, SNR is defined as the ratio of the power of the signal $P_{signal}$ to the power of the noise $P_{noise}$:

$$SNR = \frac{P_{signal}}{P_{noise}} = \frac{A_{signal}^2}{A_{noise}^2} \qquad (12)$$

Where $A$ denotes the amplitude of the signal and the noise.

To measure the frequency distribution of the hidden information in the noisy latent $z_t$ during training, we apply 3D Fourier Transformation $\mathcal{FFT}_{3D}$ to both the clean latent $z_0$ and the Gaussian noise $\epsilon$:

$$\mathcal{F}_{z_0} = \mathcal{FFT}_{3D}(z_0), \mathcal{F}_\epsilon = \mathcal{FFT}_{3D}(\epsilon) \qquad (13)$$

The amplitude spectrum can then be derived with the absolute value of the frequency-domain representation:

$$\mathcal{A}_{z_0} = |\mathcal{F}_{z_0}|, \mathcal{A}_\epsilon = |\mathcal{F}_\epsilon| \qquad (14)$$

According to Eqn. 4 in manuscript and the linear property of the Fourier transform, the full-band SNR of $z_t$ is derived as:

$$SNR(z_t) = \frac{(\sqrt{\hat{\alpha}_t}A_{z_0})^2}{(\sqrt{1-\hat{\alpha}_t}A_\epsilon)^2} = \frac{\hat{\alpha}_t}{1-\hat{\alpha}_t}\frac{A_{z_0}^2}{A_\epsilon^2} \qquad (15)$$

Consider a frequency band $\Phi$ with spatio-temporal frequency range in $\{(f_s^L, f_s^H), (f_t^L, f_t^H)\}$, the SNR of $z_t$ in this frequency band can be calculated using a band-pass filter (BPF), or approximate by summing the amplitudes in the corresponding spatio-temporal range. Converting to logarithm scale, the SNR for frequency band $\Phi$ is finally derived as:

$$SNR_{dB}(z_t, \Phi) = 10\log_{10}\frac{\hat{\alpha}_t}{1-\hat{\alpha}_t}\frac{\sum_{f_s^L}^{f_s^H}\sum_{f_t^L}^{f_t^H}A_{z_0}^2}{\sum_{f_s^L}^{f_s^H}\sum_{f_t^L}^{f_t^H}A_\epsilon^2} \qquad (16)$$

Table A3. **User Study.** Each participant vote for the image that they consider superior for Temporal Consistency, Text-Video Alignment and Overall Quality, respectively.

| Method | Temporal Consistency | Text-Video Alignment | Overall Quality |
|---|---|---|---|
| VideoCrafter [4] | 14.29% | 22.62% | 23.81% |
| VideoCrafter+FreeInit | **85.71%** | **77.38%** | **76.19%** |
| ModelScope [45] | 23.81% | 28.57% | 21.43% |
| ModelScope+FreeInit | **76.19%** | **71.43%** | **78.57%** |
| AnimateDiff [11] | 13.89% | 31.94% | 23.61% |
| AnimateDiff+FreeInit | **86.11%** | **68.06%** | **76.39%** |

## B. User Study

We conduct a User Study to further evaluate the influence of FreeInit. We randomly select 12 diverse text prompts for each model (VideoCrafter [4], ModelScope [45] and AnimateDiff [11]), and ask 20 participants to vote for the generation results. Specifically, each of the 20 participants are provided with the text prompt and a pair of synthesized videos, one generated from the vanilla model and the other one with FreeInit. Then the participants vote for the image that they consider superior for Temporal Consistency, Text-Video Alignment and Overall Quality, respectively. The average vote rates are shown in Table A3. The majority of the votes goes to the category using FreeInit under all evaluation metrics, which indicates that FreeInit consistently improves the quality of video generation.

## C. More Discussions on Ablation Study

**Influence of Noise Reinitialization and Filter Selection.** The qualitative results generated by the vanilla AnimateDiff model and three ablation variants of FreeInit are illustrated in Figure A11. Aligned with the conclusion in the quantitative results in the manuscript, the visual results indicate that the Noise Reinitialization process is crucial for the refinement of temporal consistency. As depicted in Figure A11, the color of the boxers' cloths and the background suffer from large inconsistencies in the original generation results. Adding the FreeInit refinement loop without Noise Reinitialization still leads to unsatisfactory appearance and undesirable temporal consistency, as no randomness is provided for developing better visual details that complements with the refined low-frequency components. With Noise Reinitialization performed by Ideal Low Pass Filter, the consistency of the subject appearance gains clear improvement, but the large noises (*e.g.*, the yellow area) are still not removed. In comparison, Reinitialization with Gaussian Low Pass Filter is able to remove the large noises, as it intro-

(a) w/o FreeInit



(b) FreeInit w/o Noise Reinitialization



(c) FreeInit w/ Noise Reinitialization (ILPF)



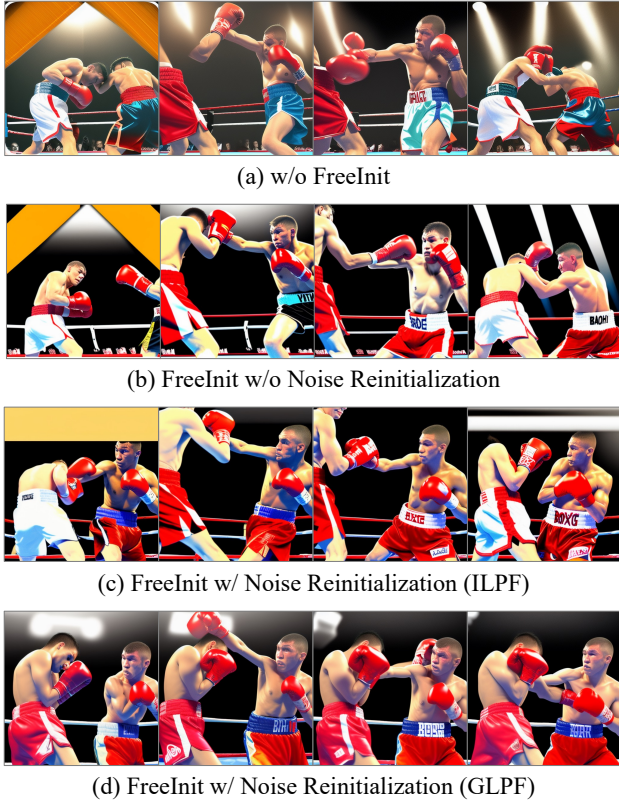(d) FreeInit w/ Noise Reinitialization (GLPF)

Figure A11. **Ablation Study on Noise Reinitialization.** Using Noise Reinitialization with a proper frequency filter is crucial for generating temporally consistent results.



(a) Iteration 0



(b) Iteration 1



(c) Iteration 2

Figure A12. **Ablation Study on Iteration Steps.** The temporal consistency of the results increases with the iteration steps.



(a) Iteration 0
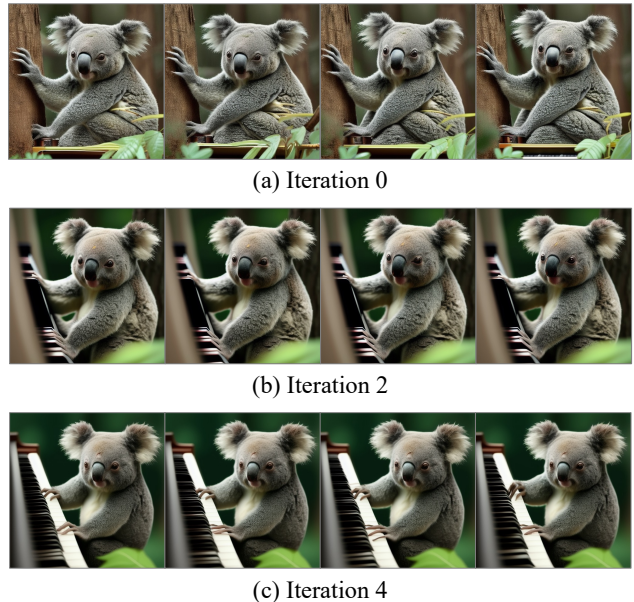


(b) Iteration 2



(c) Iteration 4

Figure A13. **Semantic Growth.** The frames are generated with the text prompt "A koala bear is playing piano in the forest". The missing semantics "playing piano" gradually grows with the FreeInit iteration.

duces an abundant amount of randomness into both mid and high-frequency, creating room for fixing large visual discrepancies. Despite gaining results with better temporal consistency, we also find using Gaussian Low Pass Filter sometimes leads to over-smoothed frames. To alleviate this side effect, the Butterworth Low Pass Filter can be utilized to replace the GLPF for keeping a balance between temporal consistency and visual quality.

**Influence of Iteration Steps.** Since the low-frequency components of the latent are refined during each FreeInit iteration, more iteration steps normally leads to generation results with better temporal consistency, as shown in Figure A12. We also observe that the missing semantics (*e.g.*, "playing piano", as depicted in Figure A13) can gradually grow with the FreeInit iteration, thanks to the refined low-frequency components and the randomness introduced by Noise Reinitialization.

## D. Implementation Details

**Base Models.** Three open-sourced text-to-video models are used as the base models for FreeInit evaluation. For VideoCrafter [4], the VideoCrafter-v0.9 Base T2V model based on the latent video diffusion models (LVDM) [13] is adopted. For ModelScope [45], we utilize the *diffusers* [7]

implementation of the text-to-video pipeline, with the text-to-video-ms-1.7b model. For AnimateDiff [11], we use the mm-sd-v14 motion module with the Realistic Vision V5.1 LoRA model for evaluation.

**Inference Details.** Experiments on VideoCrafter and ModelScope are conducted on $256 \times 256$ spatial scale and 16 frames, while experiments on AnimateDiff is conducted on
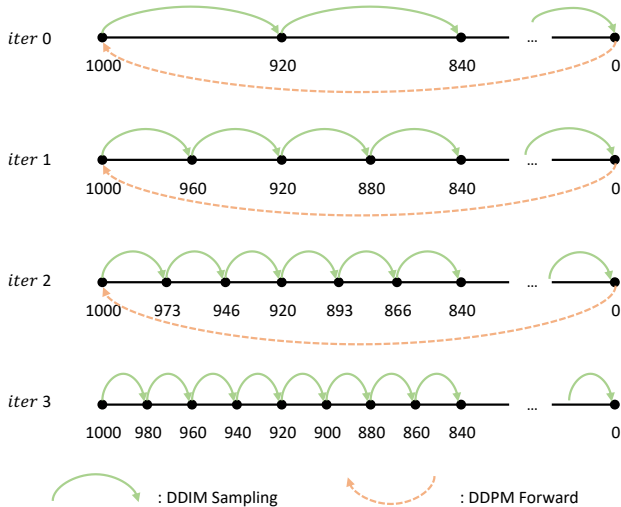
Figure A14. **Coarse-to-Fine Sampling for Fast Inference.** Early FreeInit iterations uses fewer DDIM sampling steps, while latter iterations performs more steps for detailed refinement.



(a) ModelScope



(b) ModelScope + FreeInit

Figure A15. **Failure Case.** With input prompt "Clown fish swimming through the coral reef", performing FreeInit improves temporal consistency but falsely removes the fast, small foreground object (clown fish).

a video size of $512 \times 512$, 16 frames. During the inference process, we use classifier-free guidance for all experiments including the comparisons and ablation studies, with a constant guidance weight 7.5. All experiments are conducted on a single Nvida A100 GPU.

## E. More Qualitative Comparisons

Extra qualitative comparisons with each base model are provided in Figure A16- A21. For more visual results in video format, please refer to the the video and project page.

## F. Limitations

**Inference Time.** Since FreeInit is an iterative method, a natural drawback is the increased sampling time. This issue can be mitigated through a **Coarse-to-Fine Sampling** Strategy. Specifically, the DDIM sampling steps of each FreeInit iteration can be reduced according to the iteration step: early FreeInit iterations use fewer DDIM sampling steps for a coarse refinement of the low-frequency components, while the latter iterations perform more steps for detailed refinement. A simple linear scaling strategy can be used for Coarse-to-Fine Sampling:

$$T_i = [\frac{T}{N}(i+1)] \tag{17}$$

Where $T_i$ is the DDIM steps for iteration $i$, $T$ is the commonly used DDIM steps (*e.g.*, 50), and $N$ is the number of FreeInit iterations. For instance, the fast Sampling process for $T = 50$, $N = 4$ is illustrated in Figure A14.

**Failure Cases.** In some cases where the video includes small and fast-moving foreground objects, performing FreeInit may leads to the distinction of the object. This
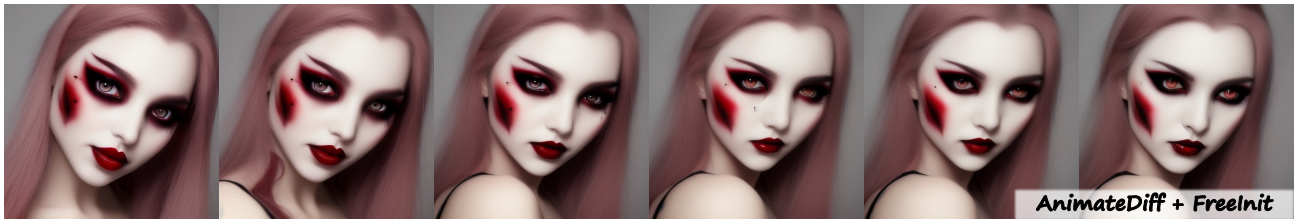
is because the iterative low-frequency refinement strategy tends to guide the generation towards the more stable low-frequency subjects. As the iteration progresses, the enhanced semantics within the initial noise's lower frequencies increasingly takes more control of the generation process, causing a partial loss of control from the text condition. This issue can be alleviated by using less FreeInit iterations, tuning the frequency filter parameters or changing the classifier-free guidance weight.

## G. Potential Negative Societal Impacts

FreeInit is a research focused on improving the inference quality of existing video diffusion models without favoring specific content categories. Nonetheless, its application in aiding other video generation models could potentially be exploited for malicious purposes, resulting in the creation of fake content.
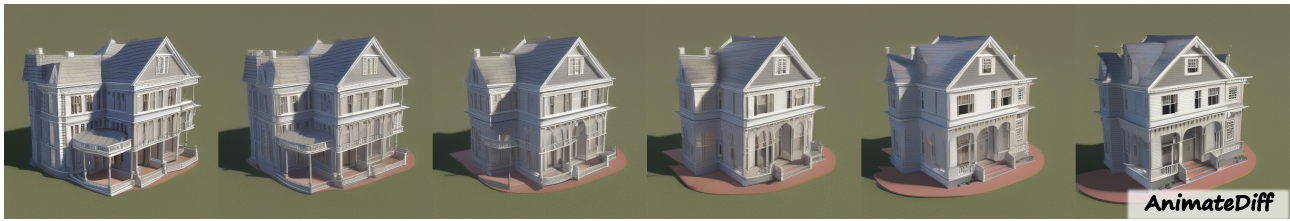
A corgi is playing drum kit.



Vampire makeup face of beautiful girl, red contact lenses.



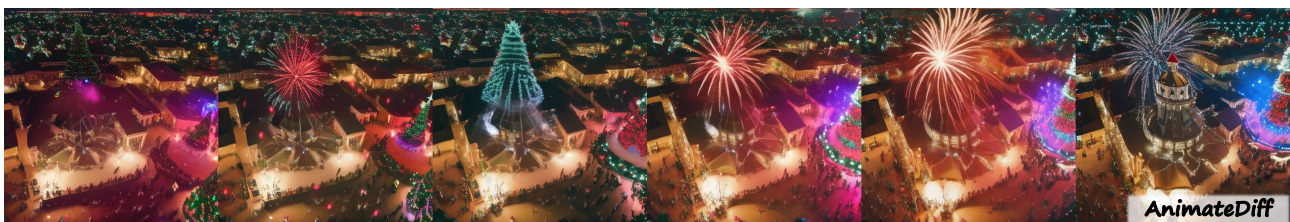A cat wearing sunglasses and working as a lifeguard at a pool.

Figure A16. **More Qualitative Results on AnimateDiff.**
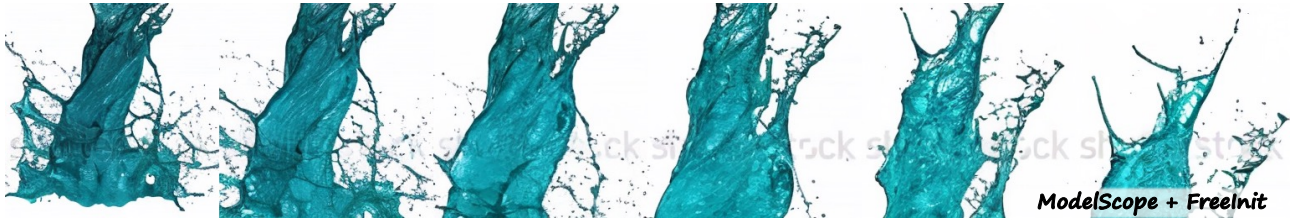
A 3D model of a 1800s victorian house.



A car's special features are being discussed on a car commercial on tv.
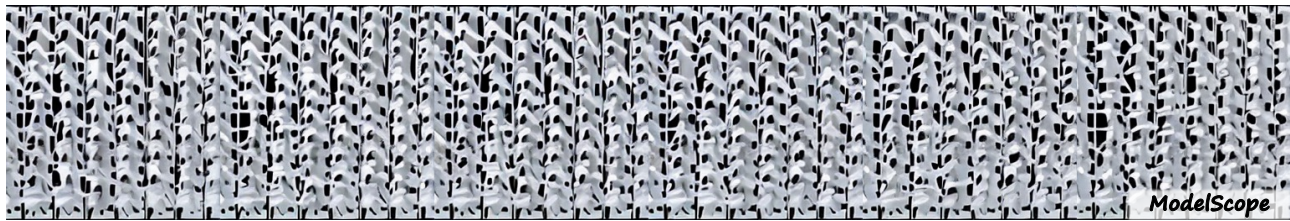


A drone view of celebration with Christmas tree and fireworks, starry sky background.
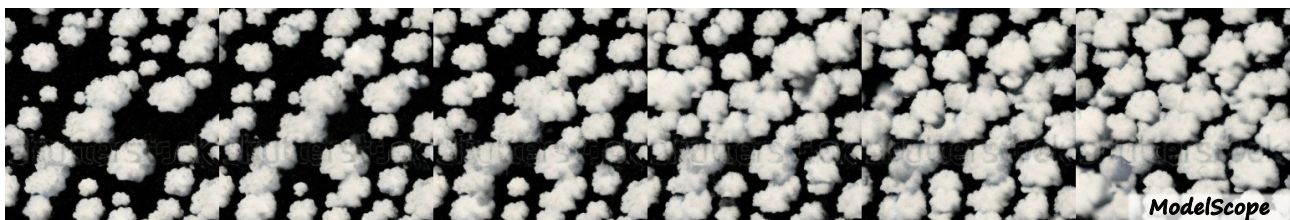
Figure A17. **More Qualitative Results on AnimateDiff.**

Splash of turquoise water in extreme slow motion, alpha channel included.



Origami dancers in white paper, 3D render, on white background, studio shot, dancing modern dance.



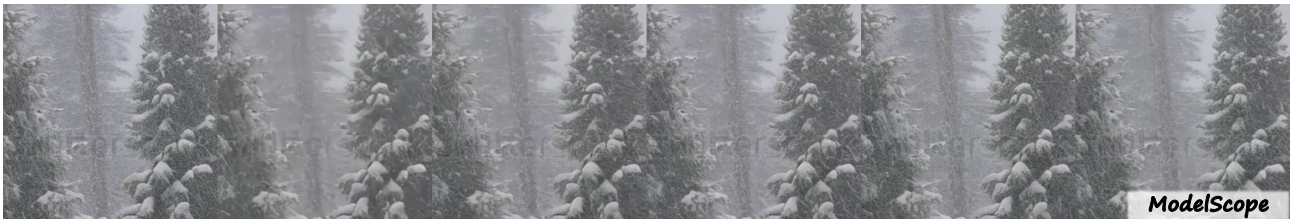An animated painting of fluffy white clouds moving in sky.

Figure A18. **More Qualitative Results on ModelScope.**

A beautiful coastal beach in spring, waves lapping on sand by Vincent van Gogh.



An oil painting of a couple in formal evening wear going home get caught in a heavy downpour with umbrellas.



A bigfoot walking in the snowstorm.

Figure A19. **More Qualitative Results on ModelScope.**
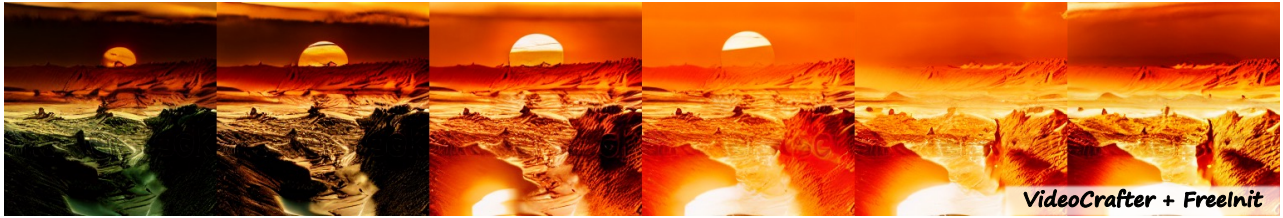
**Turtle swimming in ocean.**



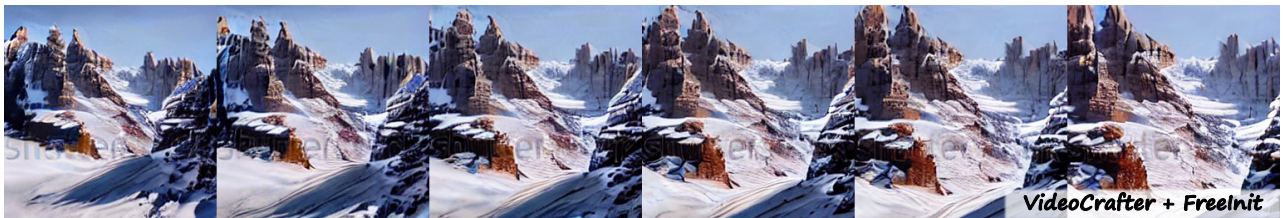**Campfire at night in a snowy forest with starry sky in the background.**



**A cute raccoon playing guitar in a boat on the ocean.**

Figure A20. **More Qualitative Results on VideoCrafter.**

**Time lapse of sunrise on mars.**



**Snow rocky mountains peaks canyon. snow blanketed rocky mountains surround and shadow deep canyons. the canyons twist and bend through the high elevated mountain peaks.**



**A shark swimming in clear Caribbean ocean.**

Figure A21. **More Qualitative Results on VideoCrafter.**