

ConditionVideo: Training-Free Condition-Guided Text-to-Video Generation

Bo Peng^{1,2} Xinyuan Chen¹ Yaohui Wang¹ Chaochao Lu¹ Yu Qiao¹

¹Shanghai Artificial Intelligence Laboratory ²Shanghai Jiao Tong University

{pengbo, chenxinyuan, wangyaohui, luchaochao, qiaoyu}@pjlab.org.cn

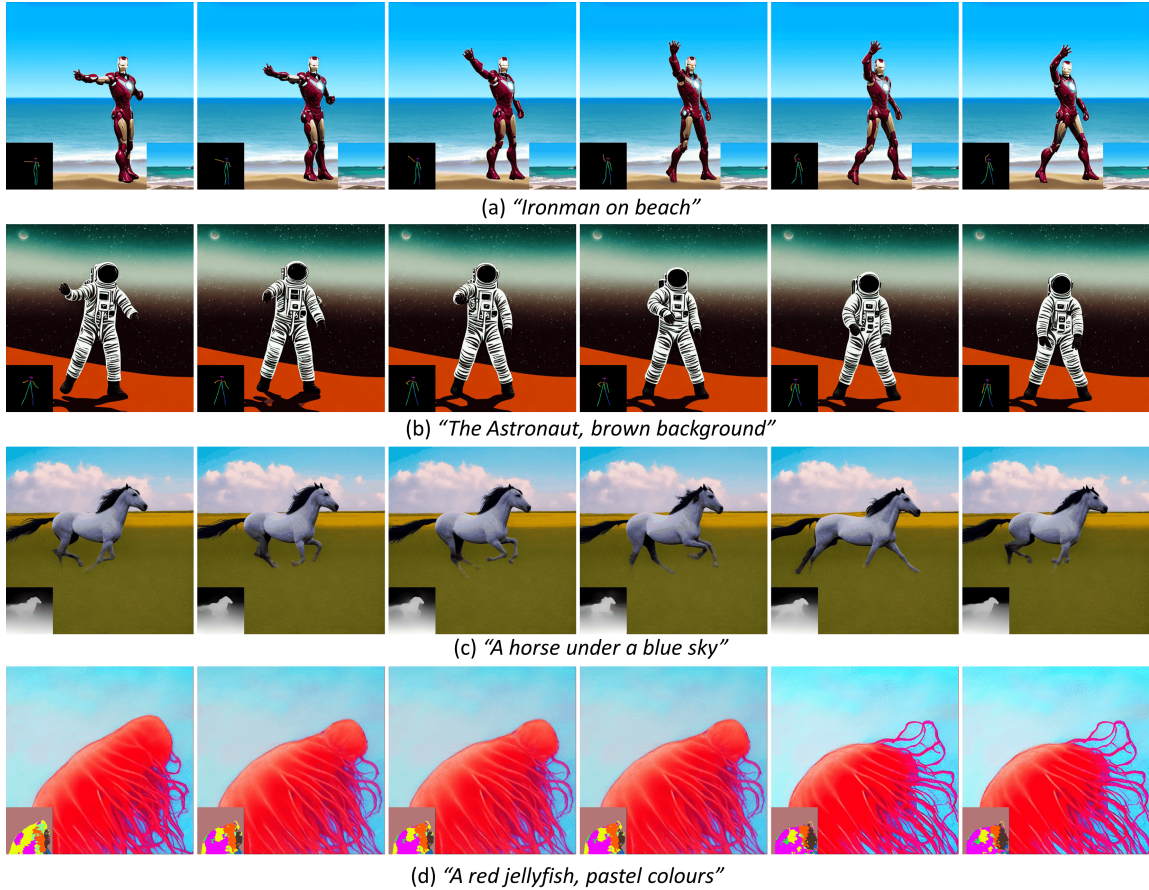


Figure 1: **Our training-free method generates videos conditioned on different inputs.** In (a), the illustration showcases the process of generation using provided scene videos and pose information, with the background wave exhibiting a convincingly lifelike motion. (b), (c), and (d) are generated based on condition only, which are pose, depth, and segmentation, respectively.

Abstract

Recent works have successfully extended large-scale text-to-image models to the video domain, producing promising results but at a high computational cost and requiring a large amount of video data. In this work, we introduce ConditionVideo, a training-free approach to text-to-video generation based on the provided condition, video, and input text, by leveraging the power of off-the-shelf text-to-image generation methods (*e.g.*, Stable Diffusion). ConditionVideo generates realistic dynamic videos from random noise or given scene videos. Our method explicitly disentangles the motion representation into condition-guided and scenery motion components. To this end, the ConditionVideo model is designed with a UNet branch and a control branch. To improve temporal coherence, we introduce sparse bi-directional spatial-temporal attention (sBiST-Attn). The 3D control network extends the conventional 2D controlnet model, aiming to strengthen conditional generation accuracy by additionally leveraging the bi-directional frames in the temporal domain. Our method exhibits superior performance in terms of frame consistency, clip score, and conditional accuracy, outperforming compared methods. See the project website at <https://pengbo807.github.io/conditionvideo-website/>.

1 Introduction

Diffusion-based models (Song, Meng, and Ermon 2021; Song et al. 2021; Ho, Jain, and Abbeel 2020; Sohl-Dickstein et al. 2015) demonstrates impressive results in large-scale text-to-image (T2I) generation (Ramesh et al. 2022; Saharia et al. 2022; Gafni et al. 2022; Rombach et al. 2022). Much of the existing research proposes to utilize image generation models for video generation. Recent works (Singer et al. 2023; Blattmann et al. 2023; Hong et al. 2022) attempt to inflate the success of the image generation model to video generation by introducing temporal modules. While these methods reuse image generation models, they still require a massive amount of video data and training with significant amounts of computing power. Tune-A-Video (Wu et al. 2022b) extends Stable Diffusion (Rombach et al. 2022) with additional attention and a temporal module for video editing by tuning one given video. It significantly decreases the training workload, although an optimization process is still necessary. Text2Video (Khachatryan et al. 2023) proposes training-free generation, however, the generated video fails to simulate natural background dynamics. Consequently, the question arises: *How can we effectively utilize image generation models without any optimization process and embed controlling information as well as modeling dynamic backgrounds for video synthesis?*

In this work, we propose ConditionVideo, a training-free conditional-guided video generation method that utilizes off-the-shelf text-to-image generation models to generate realistic videos without any fine-tuning. Specifically, aiming at generating dynamic videos, our model disentangles the representation of motion in videos into two distinct components: conditional-guided motion and scenery motion, enabling the generation of realistic and temporally consistent frames. By leveraging this disentanglement, we propose a pipeline that consists of a UNet branch and a control branch, with two separate noise vectors utilized in the sampling process. Each noise vector represents conditional-guided motion and scenery motion, respectively. To further enforce temporal consistency, we introduce sparse bi-directional spatial-temporal attention (sBiST-Attn) and a 3D control branch that leverages bi-directional adjacent frames in the temporal dimension to enhance conditional accuracy. These components strengthen our model’s ability to generate high-quality conditional-guided videos. Our ConditionVideo method outperforms the baseline methods in terms of frame consistency, conditional accuracy, and clip score.

Our key contributions are as follows. (1) We propose ConditionVideo, a training-free video generation method that leverages off-the-shelf text-to-image generation models to generate conditional-guided videos with realistic dynamic backgrounds. (2) Our method disentangles motion representation into conditional-guided and scenery motion components via a pipeline that includes a U-Net branch and a conditional-control branch. (3) We introduce sparse bi-directional spatial-temporal attention (sBiST-Attn) and a 3D conditional-control branch to improve conditional accuracy and temporal consistency, generating high-quality videos that outperform compared methods.

2 Related work

2.1 Diffusion Models

Image diffusion models have achieved significant success in the field of generation (Ho, Jain, and Abbeel 2020; Song, Meng, and Ermon 2021; Song et al. 2021), surpassing numerous generative models that were once considered state-of-the-art (Dhariwal and Nichol 2021; Kingma et al. 2021). With the assistance of large language models (Radford et al. 2021; Raffel et al. 2020), current research can generate videos from text, contributing to the prosperous of image generation (Ramesh et al. 2022; Rombach et al. 2022).

Recent works in video generation (Esser et al. 2023; Ho et al. 2022b; Wu et al. 2022b, 2021, 2022a; Hong et al. 2022; Wang et al. 2023b,c) aim to emulate the success of image diffusion models. Video Diffusion Models (Ho et al. 2022b) extends the UNet (Ronneberger, Fischer, and Brox 2015) to 3D and incorporates factorized spacetime attention (Bertasius, Wang, and Torresani 2021). Imagen Video (Saharia et al. 2022) scales this process up and achieves superior resolution. However, both approaches involve training from scratch, which is both costly and time-consuming.

Alternative methods explore leveraging pre-trained text-to-image models. Make-A-Video (Singer et al. 2023) facilitates text-to-video generation through an expanded unCLIP framework. Tune-A-Video (Wu et al. 2022b) employs a one-shot tuning pipeline to generate edited videos from input guided by text. However, these techniques still necessitate an optimization process.

Compared to these video generation methods, our training-free method can yield high-quality results more efficiently and effectively.

2.2 Conditioning Generation

While GAN-based methods have made considerable progress in conditional generation (Mirza and Osindero 2014; Wang et al. 2018; Chan et al. 2019; Wang et al. 2019; Liu et al. 2019; Siarohin et al. 2019; Zhou et al. 2022; WANG et al. 2020; Wang et al. 2020, 2022), research on the conditional generation of diffusion models is limited. For the diffusion model-based methods, T2I Adapter (Mou et al. 2023) and ControlNet (Zhang and Agrawala 2023) aim to enhance controllability through the use of extra annotations. T2IAdapter (Mou et al. 2023) proposes aligning internal knowledge with external control signals to facilitate image generation. On the other hand, ControlNet duplicates and fixes the original weight of the large pre-trained T2I model. Utilizing the cloned weight, ControlNet trains a conditional branch for task-specific image control.

For diffusion-based conditional video generation, recent works have centered on text-driven video editing (Molad et al. 2023; Esser et al. 2023; Ceylan, Huang, and Mitra 2023; Liu et al. 2023a; Wang et al. 2023a; Qi et al. 2023). These methods prioritize similarity with the input video rather than creating new content. In contrast, our method uses dynamic features from the input reference video for more creative generation and can add extra foreground movements. Concurrent works like Follow-Your-Pose (Ma et al. 2023) and Text2Video-Zero (Khachatryan et al. 2023)

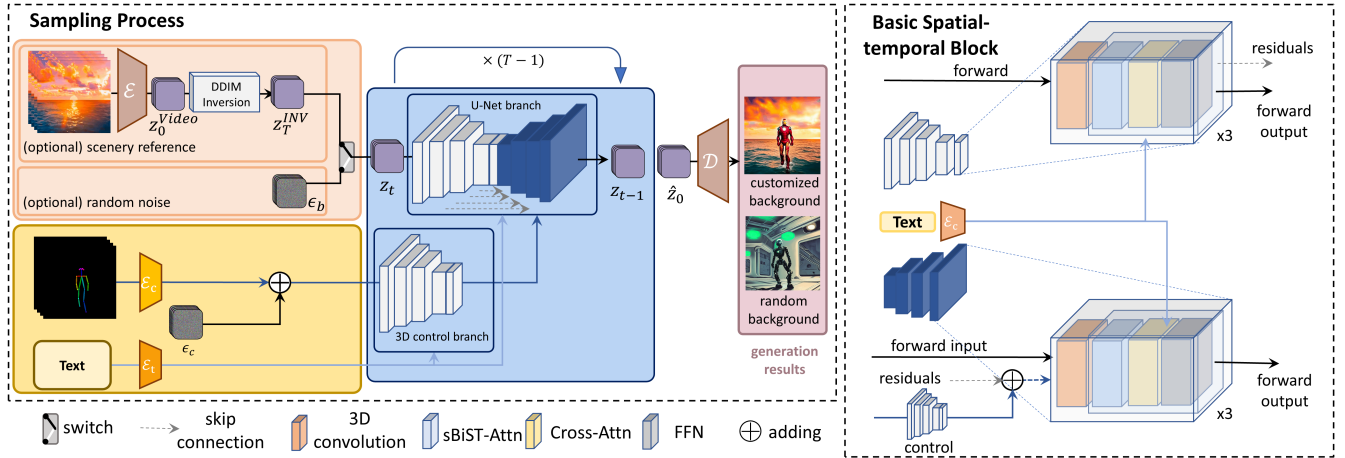


Figure 2: **Illustration of our proposed training-free pipeline.** (Left) Our framework consists of a UNet branch and a 3D control branch. The UNet branch receives either the inverted reference video z_T^{INV} or image-level noise ϵ_b for background generation. The 3D control branch receives an encoded condition for foreground generation. The text description is fed into both branches. (Right) Illustration of our basic spatial-temporal block. We employ our proposed sBiST-Attn module into the basic block between the 3D convolution block and the cross-attention block. The detail of sBiST-Attn module is shown in Fig. 3

generate videos based on given conditions. However, these methods still require training or have difficulty in generating realistic background movements (e.g., the flow of the waves in Fig. 1 (a). The dynamic version on our website can better show the advantages of our method.). Moreover, we propose additional techniques to improve time and conditioning consistency and introduce dynamic scene referencing, a novel approach in this field.

3 Preliminaries

Stable Diffusion. Stable Diffusion employs an autoencoder (Van Den Oord, Vinyals et al. 2017) to preprocess images. An image x in RGB space is encoded into a latent form by encoder \mathcal{E} and then decoded back to RGB space by decoder \mathcal{D} . The diffusion process operates with the encoded latent $z = \mathcal{E}(x)$.

For the diffusion forward process, Gaussian noise is iteratively added to latent z_0 over T iterations (Ho, Jain, and Abbeel 2020):

$$q(z_t | z_{t-1}) = \mathcal{N}\left(z_t; \sqrt{1 - \beta_t}z_{t-1}, \beta_t I\right), \quad (1)$$

$$t = 1, 2, \dots, T,$$

where $q(z_t | z_{t-1})$ denotes the conditional density function and β is given.

The backward process is accomplished by a well-trained Stable Diffusion model that incrementally denoises the latent variable \hat{z}_0 from the noise z_T . Typically, the T2I diffusion model leverages a UNet architecture, with text conditions being integrated as supplementary information. The trained diffusion model can also conduct a deterministic forward process, which can be restored back to the original z_0 . This deterministic forward process is referred to as

DDIM inversion (Song, Meng, and Ermon 2021; Dhariwal and Nichol 2021). We will refer to z_T as the noisy latent code and z_0 as the original latent in the subsequent section. Unless otherwise specified, the frames and videos discussed henceforth refer to those in latent space.

ControlNet. ControlNet (Zhang and Agrawala 2023) enhances pre-trained large-scale diffusion models by introducing extra input conditions. These inputs are processed by a specially designed conditioning control branch, which originates from a clone of the encoding and middle blocks of the T2I diffusion model and is subsequently trained on task-specific datasets. The output from this control branch is added to the skip connections and the middle block of the T2I model’s UNet architecture.

4 Methods

ConditionVideo leverages guided annotation, denoted as *Condition*, and optional reference scenery, denoted as *Video*, to generate realistic videos. We start by introducing our training-free pipeline in Sec. 4, followed by our method for modeling motion in Sec. 4.2. In Sec. 4.3, we present our sparse bi-directional spatial-temporal attention (sBiST-Attn) mechanism. Finally, a detailed explanation of our proposed 3D control branch is provided in Sec. 4.4.

4.1 Training-Free Sampling Pipeline

Fig. 2 depicts our proposed training-free sampling pipeline. Inheriting the autoencoder $\mathcal{D}(\mathcal{E}(\cdot))$ from the pre-trained image diffusion model (Sec. 3), we conduct video transformation between RGB space and latent space frame by frame. Our ConditionVideo model contains two branches: a UNet branch and a 3D control branch. A text description is fed

into both branches. Depending on the user’s preference for customized or random background, the UNet branch accepts either the inverted code z_T^{INV} of the reference background video or the random noise ϵ_b . The condition is fed into the 3D control branch after being added with random noise ϵ_c . We will further describe this disentanglement input mechanism and random noise ϵ_b, ϵ_c in Sec. 4.2.

Our branch uses the original weight of ControlNet (Zhang and Agrawala 2023). As illustrated on the right side of Fig. 2, we modify the basic spatial-temporal blocks of these two branches from the conditional T2I model by transforming 2D convolution into 3D with $1 \times 3 \times 3$ kernel and replacing the self-attention module with our proposed sBiST-Attn module (Sec. 4.3). We keep other input-output mechanisms the same as before.

4.2 Strategy for Motion Representation

Disentanglement for Latent Motion Representation In conventional diffusion models for generation (e.g., ControlNet), the noise vector ϵ is sampled from an i.i.d. Gaussian distribution $\epsilon \sim \mathcal{N}(0, I)$ and then shared by both the control branch and UNet branch. However, if we follow the original mechanism and let the inverse background video’s latent code to shared by two branches, we observe that the background generation results will be blurred (Experiments are shown in the Appx B.). This is because using the same latent to generate both the foreground and the background presumes that the foreground character has a strong relationship with the background. Motivated by this observation, we explicitly disentangle the video motion presentation into two components: the motion of the background and the motion of the foreground. The background motion is generated by the UNet branch whose latent code is presented as background noise $\epsilon_b \sim \mathcal{N}(0, I)$. The foreground motion is represented by the given conditional annotations while the appearance representation of the foreground is generated from the noise $\epsilon_c \sim \mathcal{N}(0, I)$.

Strategy for Temporal Consistency Motion Representation To attain temporal consistency across consecutively generated frames, We investigated selected noise patterns that facilitate the creation of cohesive videos. Consistency in foreground generation can be established by ensuring that the control branch produces accurate conditional controls. Consequently, we propose utilizing our control branch input for this purpose: $C_{cond} = \epsilon_c + \mathcal{E}_c(Condition), \epsilon_{c_i} \in \epsilon_c, \epsilon_{c_i} \sim \mathcal{N}(0, I) \subseteq \mathbb{R}^{H \times W \times C}, \forall i, j = 1, \dots, F, \epsilon_{c_i} = \epsilon_{c_j}$, where H, W , and C denote the height, width, and channel of the latent z_t , F represents the total frame number, C_{cond} denotes the encoded conditional vector which will be fed into the control branch and \mathcal{E}_c denotes the conditional encoder. Additionally, it’s important to observe that ϵ_{c_i} corresponds to a single frame of noise derived from the video-level noise denoted as ϵ_c . The same relationship applies to ϵ_{b_i} and ϵ_b as well.

When generating backgrounds, there are two approaches we could take. The first is to create the background using background noise ϵ_b : $\epsilon_{b_i} \in \epsilon_b, \epsilon_{b_i} \sim \mathcal{N}(0, I) \subseteq \mathbb{R}^{H \times W \times C}, \epsilon_{b_i} = \epsilon_{b_j}, \forall i, j = 1, \dots, F$. The second approach is to gen-

erate the background from an inverted latent code, z_T^{INV} , of the reference scenery video. Notably, we observed that the dynamic motion correlation present in the original video is retained when it undergoes DDIM inversion. So we utilize this latent motion correlation to generate background videos. Our ConditionVideo method is more user-friendly and cost-efficient compared to techniques that require motion training.

Algorithm 1: Sampling Algorithm

Input: *Condition, Text, Video*(Optional)

Parameter: T

Output: \hat{X}_0 : generated video

```

1: if Video is not None then
2:    $z_0^{Video} \leftarrow \mathcal{E}(Video)$  //encode video
3:    $z_T^{NV} \leftarrow \text{DDIM\_Inversion}(z_0^{Video}, T, \text{UNetBranch})$ 
4:    $z_T \leftarrow z_T^{INV}$  //customize background
5: else
6:    $z_T \leftarrow \epsilon_b$ , //random background
7: end if
8:  $C_{cond} \leftarrow \epsilon_c + \mathcal{E}_c(Condition)$  //encode condition
9:  $C_{text} \leftarrow \mathcal{E}_t(Text)$  //encode input prompt
10: for  $t = T \dots 1$  do
11:    $c_t \leftarrow \text{ConrtolBranch}(C_{cond}, t, C_{text})$ 
12:    $\hat{z}_{t-1} \leftarrow \text{DDIM\_Backward}(z_t, t, C_{text}, c_t, \text{UNetBranch})$ 
13: end for
14:  $\hat{X}_0 \leftarrow \mathcal{D}(\hat{z}_0)$ 
15: return  $\hat{X}_0$ 

```

During the sampling process, in the first forward step $t = T$, we feed the background latent code z_T^{INV} or ϵ_b into the UNet branch and the condition C_{cond} into our 3D control branch. Then, during the subsequent reverse steps $t = T - 1, \dots, 0$, we feed the denoised latent z_t into the UNet branch while still using C_{cond} for 3D control branch input. The details of the sampling algorithm are shown in Alg. 1

4.3 Sparse Bi-directional Spatial-Temporal Attention (sBiST-Attn)

Taking into account both temporal coherence and computational complexity, we propose a sparse bi-directional spatial-temporal attention (sBiST-Attn) mechanism, as depicted in Fig. 3. For video latent $z_t^i, i = 1, \dots, F$, the attention matrix is computed between frame z_t^i and its bi-directional frames, sampled with a gap of 3. This interval was chosen after weighing frame consistency and computational cost (see Appx C.1). For each z_t^i in z_t , we derive the query feature from its frame z_t^i . The key and value features are derived from the bi-directional frames $z_t^{3j+1}, j = 0, \dots, \lfloor (F - 1)/3 \rfloor$. Mathematically, our sBiST-Attn can be expressed as:

$$\begin{cases} \text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d}}\right) \cdot V \\ Q = W^Q z_t^i, K = W^K z_t^{[3j+1]}, V = W^V z_t^{[3j+1]}, \\ j = 0, 1, \dots, \lfloor (F - 1)/3 \rfloor \end{cases} \quad (2)$$

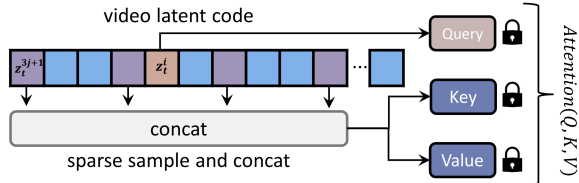


Figure 3: **Illustration of ConditionVideo’s sparse bi-directional attention.** The purple blocks signify the frame we’ve selected for concatenation, which can be computed for key and value. The pink block represents the current block from which we’ll calculate the query. The blue blocks correspond to the other frames within the video sequence. Latent features of frame z_t^i , bi-directional frames z_t^{3j+1} , $j = 0, \dots, \lfloor (F-1)/3 \rfloor$ are projected to query Q , key K and value V . Then the attention-weighted sum is computed based on key, query, and value. Notice that the parameters are the same as the ones in the self-attention module of the pre-trained image model.

where $[\cdot]$ denotes the concatenation operation, and W^Q, W^K, W^V are the weighted matrices that are identical to those used in the self-attention layers of the image generation model.

4.4 3D Control Branch

Frame-wise conditional guidance is generally effective, but there may be instances when the network doesn’t correctly interpret the guide, resulting in an inconsistent conditional output. Given the continuous nature of condition movements, ConditionVideo propose enhancing conditional alignment by referencing neighboring frames. If a frame isn’t properly aligned due to weak control, other correctly aligned frames can provide more substantial conditional alignment information. In light of this, we design our control branch to operate temporally, where we choose to replace the self-attention module with the sBiST-Attn module and inflate 2D convolution to 3D. The replacing attention module can consider both previous and subsequent frames, thereby bolstering our control effectiveness.

5 Experiments

5.1 Implementation Details

We implement our model based on the pre-trained weights of ControlNet (Zhang and Agrawala 2023) and Stable Diffusion (Rombach et al. 2022) 1.5. We generate 24 frames with a resolution of 512×512 pixels for each video. During inference, we use the same sampling setting as Tune-A-Video (Wu et al. 2022b).

5.2 Main results

In Fig. 1, we display the success of our training-free video generation technique. The generated results from ConditionVideo, depicted in Fig. 1 (a), imitate moving scenery videos and show realistic waves as well as generate the correct character movement based on posture. Notably, the

style of the backgrounds is distinct from the original guiding videos, while the motion of the backgrounds remains constant. Furthermore, our model can generate consistent backgrounds when sampling ϵ_b from Gaussian noise based on conditional information, as shown in Fig.1 (b),(c),(d). These videos showcase high temporal consistency and rich graphical content.

5.3 Comparison

Compared Methods We compare our method with Tune-A-Video (Wu et al. 2022b), ControlNet (Zhang and Agrawala 2023), and Text2Video-Zero (Khachatryan et al. 2023). For Tune-A-Video, we first fine-tune the model on the video from which the condition was extracted, and then sample from the corresponding noise latent code of the condition video.

Qualitative Comparison Our visual comparison conditioning on pose, canny, and depth information is presented in Fig. 4, 5, and 6. Tune-A-Video struggles to align well with our given condition and text description. ControlNet demonstrates improvement in condition-alignment accuracy but suffers from a lack of temporal consistency. Despite the capability of Text2Video to produce videos of exceptional quality, there are still some minor imperfections that we have identified and indicated using a red circle in the figure. Our model surpasses all others, showcasing outstanding condition-alignment quality and frame consistency.

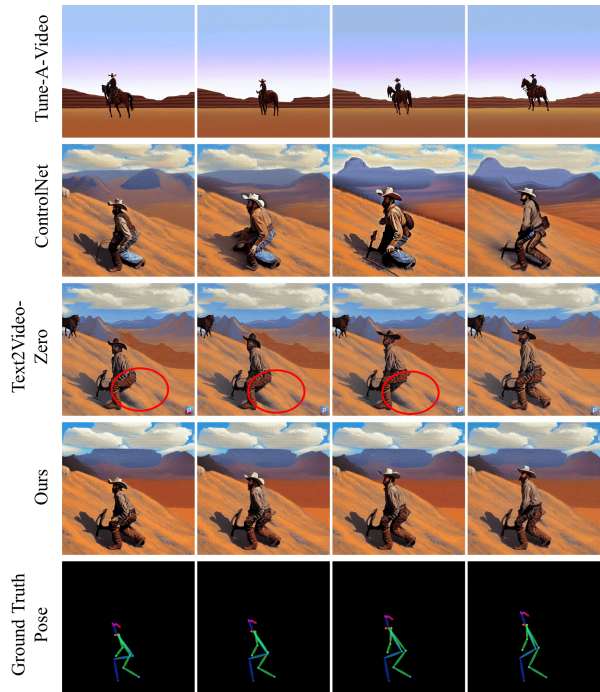


Figure 4: **Qualitative comparison condition on the pose.** “The Cowboy, on a rugged mountain range, Western painting style”. Our results outperform in both temporal consistency and pose accuracy, while others have difficulty in maintaining either one or both of the qualities.

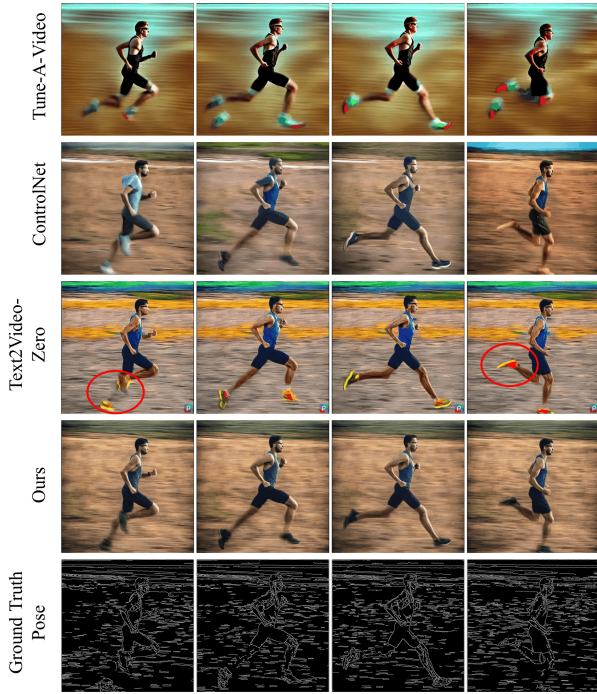


Figure 5: **Qualitative comparison condition on canny.** “A man is runnin”. Tune-A-Video fails in canny-alignment. ControlNet generates low temporal consistent frames. Although Text2Video outperforms the first two methods, it generates parts of the legs that do not correspond to the real structure of the human body and the shoes are not the same color.

Method	FC(%)	CS	PA (%)
Tune-A-Video	95.84	30.74	26.13
ControlNet	94.22	32.97	79.51
Text2Video-Zero	98.82	32.84	78.50
Ours	99.02	33.03	83.12

Table 1: Quantitative comparisons condition on pose. FC,CS,PA represent *frame consistency*, *clip score* and *pose accuracy*, respectively

Quantitative Comparison We evaluate all the methods using three metrics: *frame consistency* (Esser et al. 2023; Wang et al. 2023a; Radford et al. 2021), *clip score* (Ho et al. 2022a; Hessel et al. 2021; Park et al. 2021), and *pose accuracy* (Ma et al. 2023). As other conditions are hard to evaluate, we use pose accuracy for conditional consistency only. The results on different conditions are shown in Tab. 1 and 2. We achieve the highest frame consistency, and clip score in all conditions, indicating that our method exhibits the best text alignment. We also have the best pose-video alignment among the other three techniques of conditioning on the pose.

The conditions are randomly generated from a group of 120 different videos. For more information please see Appx D.2.

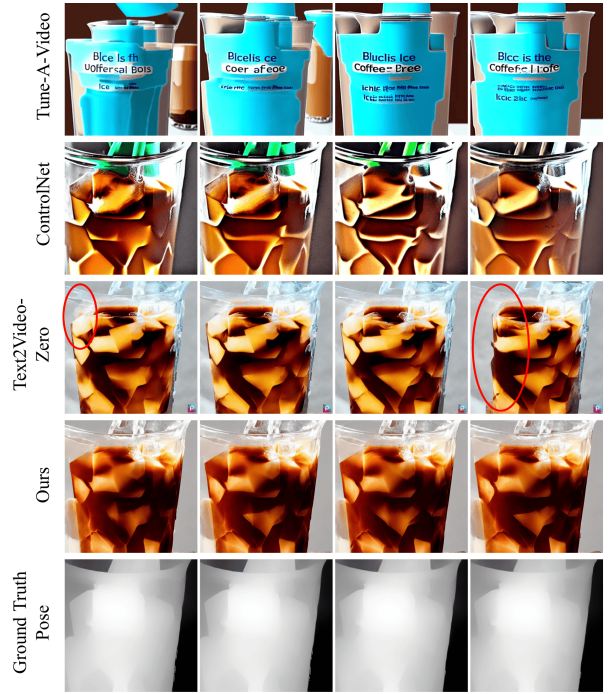


Figure 6: **Qualitative comparison condition on depth.** “ice coffee”. All three methods used for comparison have the problem of changing the appearance of the object when the viewpoint is switched, and only our method ensures the consistency of the appearance before and after.

5.4 Ablation Study

We conduct an ablation study on the pose condition, temporal module, and 3D control branch. Our qualitative results are visualized in Fig. 7. In this study, we alter each component for comparison while keeping all other settings the same.

Ablation on Pose Condition We evaluate performance with and without using pose, as shown in Fig. 7. Without pose conditioning, the video is fixed as an image, while the use of pose control allows for the generation of videos with certain temporal semantic information.

Ablation on Temporal Module Training-free video generation heavily relies on effective spatial-temporal modeling. In addition to comparing with a self-attention mechanism without temporal, we conduct an ablation study on three different spatial-temporal mechanisms. First, we remove our sBiST-attention mechanism and replaced it with Sparse-Causal attention (Wu et al. 2022b). Then, we compare our bi-directional attention mechanism with a dense attention mechanism (Wang et al. 2023a) which attends to all frames for key and value.

The results are presented in Tab. 3. A comparison of temporal and non-temporal attention underlines the importance of temporal modeling for generating time-consistent videos. By comparing our method with Sparse Causal attention, we demonstrate the effectiveness of ConditionVideo’s sBiST at-

Method	Condition	FC(%)	CS
Tune-A-Video	-	95.84	30.74
ControlNet	Canny	90.53	29.65
Text2Video-Zero	Canny	97.44	28.76
Ours	Canny	97.64	29.76
ControlNet	Depth	90.63	30.16
Text2Video-Zero	Depth	97.46	29.38
Ours	Depth	97.65	30.54
ControlNet	Segment	91.87	31.85
Ours	Segment	98.13	32.09

Table 2: Quantitative comparisons condition on canny, depth and segment.

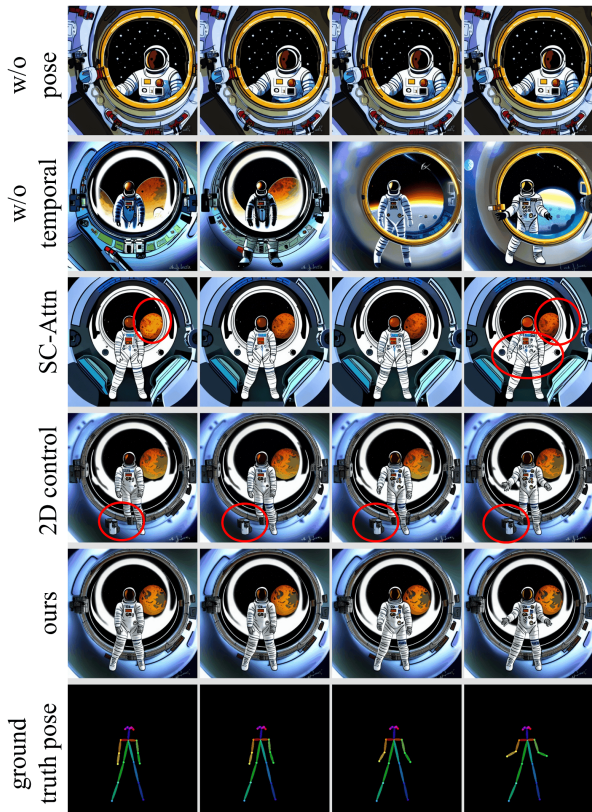


Figure 7: **Ablations of each component, generated from image-level noise.** “The astronaut, in a spacewalk, sci-fi digital art style”. 1st row displays the generation result without pose conditioning. 2nd and 3rd rows show the results after replacing our sBiST-Attn with self-Attn and SC-Attn (Wu et al. 2022b). 4th row presents the result with the 2D condition-control branch.

tention module, proving that incorporating information from bi-directional frames improves performance compared to using only previous frames. Furthermore, we observe almost no difference in frame consistency between our method and dense attention, despite the latter requiring more than double

Method	FC(%)	Time
w/o Temp-Attn	94.22	31s
S-C Attn	98.77	43s
sBiST-Attn	99.02	1m30s
Full-Attn	99.03	3m37s

Table 3: Ablations on temporal module. Time represents the duration required to generate a 24-frame video with a size of 512x512.

Method	FC(%)	CS (%)	PA (%)
2D control	99.03	33.11	81.26
3D control	99.02	33.03	83.12

Table 4: Ablation on 3D control branch. FC, CS, PA represent frame consistency, clip score, and pose-accuracy, respectively.

our generation duration.

Ablations on 3D control branch We compare our 3D control branch with a 2D version that processes conditions frame-by-frame. For the 2D branch, we utilize the original ControlNet conditional branch. Both control branches are evaluated in terms of frame consistency, clip score, and pose accuracy. Results in Tab. 4 show that our 3D control branch outperforms the 2D control branch in pose accuracy while maintaining similar frame consistency and clip scores. This proves that additional consideration of bi-directional frames enhances pose control.

6 Discussion and Conclusion

In this paper, we propose ConditionVideo, a training-free method for generating videos with reasonable motion. We introduce a method that generates motion representation conditioned on background video and conditional information. Our method additionally strengthens frame consistency and condition alignment through our sBiST-Attn mechanism and 3D control branch. Experimental results demonstrate that our method can generate high-quality videos, opening new avenues for research in video generation and AI-based content creation.

While the condition-based and enhanced temporal attention blocks contribute to enhancing the temporal coherence of the video, we have observed that using sparse conditions, such as pose information, can still lead to videos with noticeable flickering. To address this issue, a potential solution would involve incorporating more densely sampled control inputs and additional temporal-related structures.

References

- Bertasius, G.; Wang, H.; and Torresani, L. 2021. Is space-time attention all you need for video understanding? In *ICML*, volume 2.
- Blattmann, A.; Rombach, R.; Ling, H.; Dockhorn, T.; Kim, S. W.; Fidler, S.; and Kreis, K. 2023. Align your latents:

- High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Ceylan, D.; Huang, C.-H. P.; and Mitra, N. J. 2023. Pix2video: Video editing using image diffusion. *arXiv preprint arXiv:2303.12688*.
- Chan, C.; Ginosar, S.; Zhou, T.; and Efros, A. A. 2019. Everybody dance now. In *Proceedings of the IEEE/CVF international conference on computer vision*.
- Dhariwal, P.; and Nichol, A. 2021. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34.
- Esser, P.; Chiu, J.; Atighehchian, P.; Granskog, J.; and Germanidis, A. 2023. Structure and content-guided video synthesis with diffusion models. *arXiv preprint arXiv:2302.03011*.
- Gafni, O.; Polyak, A.; Ashual, O.; Sheynin, S.; Parikh, D.; and Taigman, Y. 2022. Make-a-scene: Scene-based text-to-image generation with human priors. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XV*. Springer.
- Hessel, J.; Holtzman, A.; Forbes, M.; Le Bras, R.; and Choi, Y. 2021. CLIPScore: A Reference-free Evaluation Metric for Image Captioning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics.
- Ho, J.; Chan, W.; Saharia, C.; Whang, J.; Gao, R.; Gritsenko, A.; Kingma, D. P.; Poole, B.; Norouzi, M.; Fleet, D. J.; et al. 2022a. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*.
- Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33.
- Ho, J.; and Salimans, T. 2022. Classifier-free diffusion guidance.
- Ho, J.; Salimans, T.; Gritsenko, A.; Chan, W.; Norouzi, M.; and Fleet, D. J. 2022b. Video diffusion models.
- Hong, W.; Ding, M.; Zheng, W.; Liu, X.; and Tang, J. 2022. CogVideo: Large-scale Pretraining for Text-to-Video Generation via Transformers.
- Khachatryan, L.; Movsisyan, A.; Tadevosyan, V.; Henschel, R.; Wang, Z.; Navasardyan, S.; and Shi, H. 2023. Text2Video-Zero: Text-to-Image Diffusion Models are Zero-Shot Video Generators. *arXiv preprint arXiv:2303.13439*.
- Kingma, D.; Salimans, T.; Poole, B.; and Ho, J. 2021. Variational diffusion models. *Advances in neural information processing systems*, 34.
- Liu, S.; Zhang, Y.; Li, W.; Lin, Z.; and Jia, J. 2023a. Video-P2P: Video Editing with Cross-attention Control. *arXiv preprint arXiv:2303.04761*.
- Liu, W.; Piao, Z.; Min, J.; Luo, W.; Ma, L.; and Gao, S. 2019. Liquid warping gan: A unified framework for human motion imitation, appearance transfer and novel view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*.
- Liu, Y.; Han, T.; Ma, S.; Zhang, J.; Yang, Y.; Tian, J.; He, H.; Li, A.; He, M.; Liu, Z.; et al. 2023b. Summary of chatgpt/gpt-4 research and perspective towards the future of large language models. *arXiv preprint arXiv:2304.01852*.
- Ma, Y.; He, Y.; Cun, X.; Wang, X.; Shan, Y.; Li, X.; and Chen, Q. 2023. Follow Your Pose: Pose-Guided Text-to-Video Generation using Pose-Free Videos. *arXiv preprint arXiv:2304.01186*.
- Mirza, M.; and Osindero, S. 2014. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*.
- MMPose Contributors. 2020. OpenMMLab Pose Estimation Toolbox and Benchmark.
- Molad, E.; Horwitz, E.; Valevski, D.; Acha, A. R.; Matias, Y.; Pritch, Y.; Leviathan, Y.; and Hoshen, Y. 2023. Dreamix: Video diffusion models are general video editors. *arXiv preprint arXiv:2302.01329*.
- Mou, C.; Wang, X.; Xie, L.; Zhang, J.; Qi, Z.; Shan, Y.; and Qie, X. 2023. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453*.
- Nicki Skafta Detlefsen; Jiri Borovec; Justus Schock; Ananya Harsh; Teddy Koker; Luca Di Liello; Daniel Stancl; Changsheng Quan; Maxim Grechkin; and William Falcon. 2022. TorchMetrics - Measuring Reproducibility in PyTorch.
- Park, D. H.; Azadi, S.; Liu, X.; Darrell, T.; and Rohrbach, A. 2021. Benchmark for compositional text-to-image synthesis. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*.
- Qi, C.; Cun, X.; Zhang, Y.; Lei, C.; Wang, X.; Shan, Y.; and Chen, Q. 2023. Fatezero: Fusing attentions for zero-shot text-based video editing.
- Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR.
- Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1).
- Ramesh, A.; Dhariwal, P.; Nichol, A.; Chu, C.; and Chen, M. 2022. Hierarchical text-conditional image generation with clip latents.
- Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Ronneberger, O.; Fischer, P.; and Brox, T. 2015. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III* 18. Springer.

- Saharia, C.; Chan, W.; Saxena, S.; Li, L.; Whang, J.; Denton, E. L.; Ghasemipour, K.; Gontijo Lopes, R.; Karagol Ayan, B.; Salimans, T.; et al. 2022. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35.
- Siarohin, A.; Lathuilière, S.; Tulyakov, S.; Ricci, E.; and Sebe, N. 2019. First order motion model for image animation. *Advances in Neural Information Processing Systems*, 32.
- Singer, U.; Polyak, A.; Hayes, T.; Yin, X.; An, J.; Zhang, S.; Hu, Q.; Yang, H.; Ashual, O.; Gafni, O.; Parikh, D.; Gupta, S.; and Taigman, Y. 2023. Make-A-Video: Text-to-Video Generation without Text-Video Data. In *The Eleventh International Conference on Learning Representations*.
- Sohl-Dickstein, J.; Weiss, E.; Maheswaranathan, N.; and Ganguli, S. 2015. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*. PMLR.
- Song, J.; Meng, C.; and Ermon, S. 2020. Denoising Diffusion Implicit Models. *CoRR*, abs/2010.02502.
- Song, J.; Meng, C.; and Ermon, S. 2021. Denoising Diffusion Implicit Models. In *International Conference on Learning Representations*.
- Song, Y.; Sohl-Dickstein, J.; Kingma, D. P.; Kumar, A.; Ermon, S.; and Poole, B. 2021. Score-Based Generative Modeling through Stochastic Differential Equations. In *International Conference on Learning Representations*.
- Sun, K.; Xiao, B.; Liu, D.; and Wang, J. 2019. Deep high-resolution representation learning for human pose estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*.
- Van Den Oord, A.; Vinyals, O.; et al. 2017. Neural discrete representation learning. *Advances in neural information processing systems*.
- Wang, T.-C.; Liu, M.-Y.; Tao, A.; Liu, G.; Kautz, J.; and Catanzaro, B. 2019. Few-shot Video-to-Video Synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Wang, T.-C.; Liu, M.-Y.; Zhu, J.-Y.; Liu, G.; Tao, A.; Kautz, J.; and Catanzaro, B. 2018. Video-to-Video Synthesis. In *Conference on Neural Information Processing Systems (NeurIPS)*.
- Wang, W.; Xie, K.; Liu, Z.; Chen, H.; Cao, Y.; Wang, X.; and Shen, C. 2023a. Zero-Shot Video Editing Using Off-The-Shelf Image Diffusion Models. *arXiv preprint arXiv:2303.17599*.
- Wang, Y.; Bilinski, P.; Bremond, F.; and Dantcheva, A. 2020. G3AN: Disentangling Appearance and Motion for Video Generation. In *CVPR*.
- WANG, Y.; Bilinski, P.; Bremond, F.; and Dantcheva, A. 2020. ImaGINator: Conditional Spatio-Temporal GAN for Video Generation. In *WACV*.
- Wang, Y.; Chen, X.; Ma, X.; Zhou, S.; Huang, Z.; Wang, Y.; Yang, C.; He, Y.; Yu, J.; Yang, P.; et al. 2023b. LAVIE: High-Quality Video Generation with Cascaded Latent Diffusion Models. *arXiv preprint arXiv:2309.15103*.
- Wang, Y.; Ma, X.; Chen, X.; Dantcheva, A.; Dai, B.; and Qiao, Y. 2023c. LEO: Generative Latent Image Animator for Human Video Synthesis. *arXiv preprint arXiv:2305.03989*.
- Wang, Y.; Yang, D.; Bremond, F.; and Dantcheva, A. 2022. Latent Image Animator: Learning to Animate Images via Latent Space Navigation. In *ICLR*.
- Wu, C.; Huang, L.; Zhang, Q.; Li, B.; Ji, L.; Yang, F.; Sapiro, G.; and Duan, N. 2021. Godiva: Generating open-domain videos from natural descriptions. *arXiv preprint arXiv:2104.14806*.
- Wu, C.; Liang, J.; Ji, L.; Yang, F.; Fang, Y.; Jiang, D.; and Duan, N. 2022a. Nüwa: Visual synthesis pre-training for neural visual world creation. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XVI*. Springer.
- Wu, J. Z.; Ge, Y.; Wang, X.; Lei, W.; Gu, Y.; Hsu, W.; Shan, Y.; Qie, X.; and Shou, M. Z. 2022b. Tune-A-Video: One-Shot Tuning of Image Diffusion Models for Text-to-Video Generation. *arXiv preprint arXiv:2212.11565*.
- Zhang, L.; and Agrawala, M. 2023. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*.
- Zhou, X.; Yin, M.; Chen, X.; Sun, L.; Gao, C.; and Li, Q. 2022. Cross attention based style distribution for controllable person image synthesis. In *European Conference on Computer Vision*, 161–178. Springer.

Appendix

In our appendix, we initially present an expanded set of generated results, along with more of the comparative content. Subsequently, we exhibit results to demonstrate the effectiveness of disentanglement in latent motion representation. Finally, we provide implementation details of our experiments, including model architecture, sampling strategy, and evaluation details.

A More Results

We’ve included more generated video results at <https://pengbo807.github.io/conditionvideo-website/>. By leveraging the large text-to-image diffusion model, our method can generate videos without any fine-tuning. Videos are created from condition annotation, various prompts, and optional reference scenery videos. When the reference scenery video is not provided, the background can be generated from random noise. Results show our model can generate dynamic and realistic videos by giving different prompts. In addition, our model can generate customized backgrounds of the video when given reference scenery videos. The multi-people half-body results are also demonstrated on the website, where we successfully generate three dancing Spiderman with the prompt “*three Spiderman at night*”.

B Disentanglement for Latent Motion Representation

We observe that the generated video’s foreground and background would be mixed if we feed the same latent noise into both the UNet-branch and control branch. Therefore, we propose to use disentanglement on latent input, which uses different inputs for two branches. As shown in Fig. 8, the method without disentanglement of latent motion representation tends to produce visuals where the foreground merges with the background. In contrast, our method successfully creates a distinct foreground.

C More Ablation and Comparison Results

C.1 more qualitative results on temporal attention

We conduct experimentation with various gap values, as depicted in Figure 9. Additionally, we present a broader range of outcomes concerning diverse attention mechanisms in Table 5, encompassing SC-Attn (Wu et al. 2022b), Full-Attn (Wang et al. 2023a), and Cross-Frame Attn from Text2Video-Zero (Khachatryan et al. 2023). Following an assessment of temporal and frame consistency, we determine the optimal sBiST-attention gap to be 3.

C.2 More quantitative results

Here, we provide more ablation study results focusing on the individual components of our model. Initially, our 3D control branch was replaced with a 2D control branch derived from ControlNet (abbr. w/o 3D control). Subsequently,

Method	FC(%)	Time
SC-Attn	98.77	43s
Cross-Frame-Attn	98.86	34s
Full-Attn(gap1)	99.03	3m37s
sBiST-Attn(gap2)	99.01	1m51s
sBiST-Attn(gap3)	99.02	1m30s
sBiST-Attn(gap4)	99.00	1m13s
sBiST-Attn(gap5)	99.00	1m7s
sBiST-Attn(gap6)	98.99	1m
sBiST-Attn(gap9)	98.99	54m
sBiST-Attn(gap12)	98.90	47s

Table 5: Ablation on temporal module

test the result after replacing our sBiST-Attn with Sparse-Causal Attention (SC-Attn) (abbr. w/o BiAttn). Fig.10 and 11 demonstrate the importance of each model component. When any of the components are replaced or removed, it leads to a noticeable decrease in performance. Note that we also find significant pose accuracy improvement with our sBiST-Attn against SC-Attn, as shown in Fig.10(a).

Figs.13,14,12 show more comparison results of our model with Text2Video-Zero condition on pose, canny and depth. From the many previous comparisons, it is clear that our model has a significant advantage over Tune-A-Video and ControlNet, so we omit the comparisons between the two models here. Guided by the more ambiguous conditions, our model is more adept at reasoning about the most plausible outlook from the bidirectional information.

D Experiment Details

D.1 Model Architecture and Sampling Strategy

The architecture of our model is shown in Tab. 7. Inspired by ControlNet (Zhang and Agrawala 2023), our model consists of two branches, namely the UNet branch and the control branch. The control branch contains 1 conditional encoder, 12 downsample blocks, and 1 middle block. The downsample blocks and the middle block share the same structure as the UNet branch’s corresponding block. The conditional encoder consists of 4 convolutional layers, which convert 512×512 conditions to 64×64 feature maps. The output of each block of the control branch is fed into the corresponding upsampler after being processed by a convolution layer. For each block within our model, we inflate the 2D convolution to a 3D convolution whose kernel size is from 3×3 to $1 \times 3 \times 3$. Additionally, all original self-attention is replaced with our proposed sBiST-attention. The convolutional layers of the condition encoder are also inflated to 3D. It’s worth noting that we retain the original Encoder and Decoder, identical to those used in the large image generation model, thereby processing videos frame by frame.

For customized background video generation, we first use the DDIM inversion (Song, Meng, and Ermon 2020) to obtain the latent noise of the reference scene video. For the random background generation, we sample random Gaussian noise. After obtaining latent noise, we perform DDIM



Figure 8: The effectiveness of disentanglement for latent motion representation. Prompt: “Fishers in the sea”.

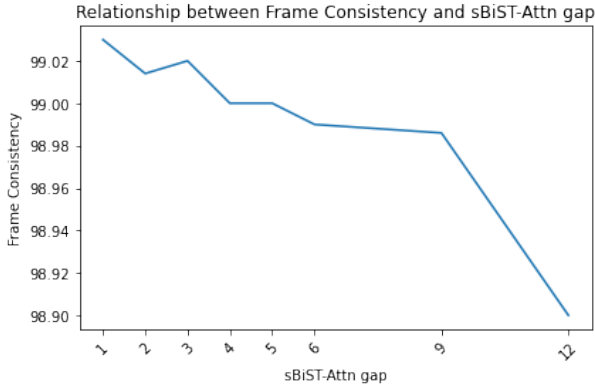


Figure 9: Ablation study of sBiST-Attn gap

sampling along with classifier-free guidance (Ho and Salimans 2022). We set the guidance scale to 12.5 and the number of time steps for DDIM sampling to 50. All the hyper-parameters of our method are shown in Tab.6.

hyper-parameters	value
DDIM inference step	50
classifier guidance scale	12.5
sBiST-Attn-gap	3

Table 6: Hyper-parameters

D.2 Evaluation Details

In our evaluation process, we carried out speed analysis experiments on a single NVIDIA A100 GPU equipped with 80GB of memory. For a fair comparison, We use the same random seed while sampling with different models. We produce 24 frames for each video, each with a resolution of 512×512 pixels, and the videos are created at a frame rate of 10 frames per second (fps). Throughout our comparative analyses with previous methods and the ablation study, each experiment was conducted with 100 reference conditions and 100 prompts. The reference conditions were extracted from a dancing video extracted from an online entertainment platform with 4800 frames. we use a frame rate of 2 and sample

48 frames continuously from the video. Then we get 100 videos with 24 frames each. We present the details of the prompt generation in the following subsection.

We assess all the methods using three measurement criteria: “frame consistency” (Esser et al. 2023; Wang et al. 2023a; Radford et al. 2021), “clip score” (Ho et al. 2022a; Hessel et al. 2021; Park et al. 2021), and “pose accuracy” (Ma et al. 2023). Next, we describe in detail how these three types of metrics are implemented.

Frame consistency Following the approaches of (Esser et al. 2023) and (Wang et al. 2023a), we calculate the consistency of the generated video. Specifically, we first calculate each frame’s CLIP image embedding (Radford et al. 2021). Then, we compute the cosine similarity of two consecutive frames’ CLIP image embeddings and take the average similarity of all frame pairs as the final score.

Clip Score We follow the approach in (Ho et al. 2022a) to compute the clip score (Hessel et al. 2021; Park et al. 2021). The clip score is computed on each frame and then averaged. We use the TorchMetrics (Nicki Skafte Detlefsen et al. 2022) package for computing clip score.

Pose Accuracy Similar to (Ma et al. 2023), we use HRNet (Sun et al. 2019) to compute our video-pose alignment quality. Specifically, we employ HRNet to predict the skeleton of the generated video, as well as extract the ground truth skeleton from the video where the original guiding pose heatmap was obtained. We normalize the extracted skeletons, and then compute the distance between each keypoint in the pose between the generated skeleton and the ground truth skeleton. We set a threshold of 0.05. If the distance is smaller than the threshold, we regard it as accurate. We compute the average accuracy across all key point skeletons to get the final score. As Tune-A-Video cannot do pose control, we consider the pose of the inverted video as the ground truth poses for computation. We use HRNet (Sun et al. 2019), available in the MMPose repository (MMPose Contributors 2020) as our pose detector.

D.3 Prompts

We employed ChatGPT (Liu et al. 2023b) to generate 100 random prompts and produce 100 arbitrary prompts with different scenery, styles, and characters such as Ironman, robot, and spaceman. We then employed these prompts to

UNet Branch	Feature Map Size	Control Branch	Feature Map Size
Encoder	(64,64)	Conditional Encoder	(24,64,64)
DownsampleBlock1 ($\times 3$)	(24,64,64)	DownsampleBlock1 ($\times 3$)	(24,64,64)
DownsampleBlock2 ($\times 3$)	(24,32,32)	DownsampleBlock2 $\times 3$	(24,32,32)
DownsampleBlock3 ($\times 3$)	(24,16,16)	DownsampleBlock3 ($\times 3$)	(24,16,16)
DownsampleBlock4 ($\times 3$)	(24,8,8)	DownsampleBlock4 ($\times 3$)	(24,8,8)
MiddleBlock	(24,8,8)	MiddleBlock	(24,8,8)
UpsampleBlock4 ($\times 3$)	(24,8,8)	convolution	(24,64,64)
UpsampleBlock3 ($\times 3$)	(24,16,16)	convolution	(24,64,64)
UpsampleBlock2 ($\times 3$)	(24,32,32)	convolution	(24,64,64)
UpsampleBlock1 ($\times 3$)	(24,64,64)	convolution	(24,64,64)
Decoder	(64,64)		

Table 7: Architecture of our model

create video clips. For the pose accuracy evaluation, we requested chatGPT to generate prompts exclusively containing human characters, because the utilized pose recognition model (HRNet (Sun et al. 2019)) would dramatically decrease its recognizing accuracy when given non-human objects (*i.e.*, robots). Specifically, we additionally required chatGPT to give the text that would allow the HRNet to easily recognize the generated video. We used the prompt “*give me 100 prompts to generate video, preferably in the form of [description of human character], in/on [background description], such as The Batman, in the sea’, The Iron man, in the forest, Makoto Shinkai style’, oil painting of a beautiful girl avatar style’. You should make sure a pose detector like HRNet can detect the human pose of the generated video*”. For frame consistency and clip score evaluations, we removed the constraints for human-like characters only, allowing us to assess model performances in a wider range of scenarios.

We use the following prompts to evaluate the proposed method.

List of 100 Prompts for Frame Consistency and Clip Score Evaluation

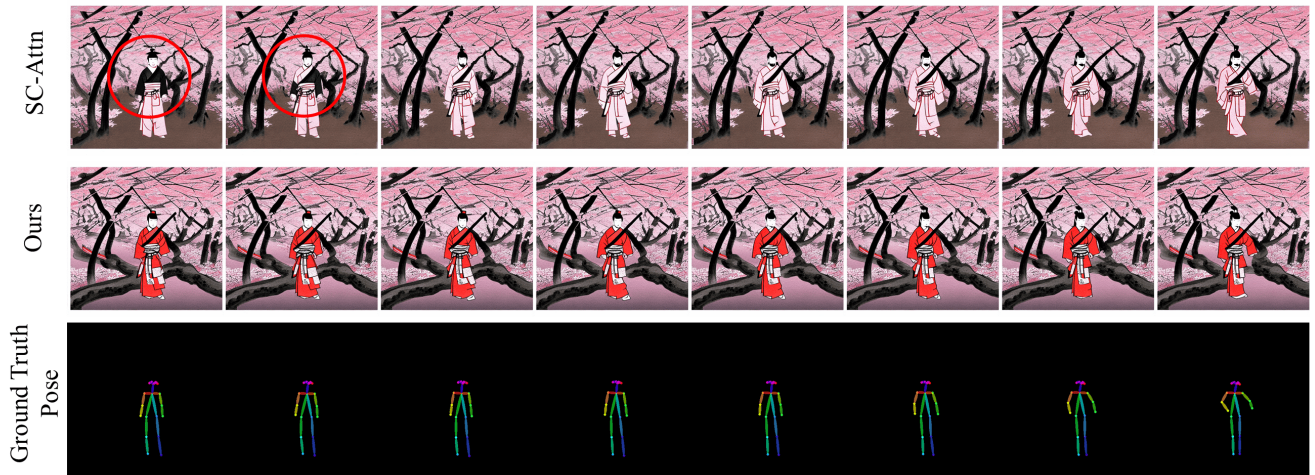
1. a man in a suit and tie
2. a man standing on top of a cliff
3. man on hill watching a meteor, cartoon artwork
4. The Samurai, on a mountaintop, traditional ink painting style
5. The Space Explorer, in an interstellar spaceship, futuristic anime style
6. The Pirate, on a sandy beach, digital painting style
7. The Cyborg, in a neon-lit alleyway, cyberpunk anime style
8. The Cowboy, on a dusty plain, oil painting style
9. The Ninja, in a bamboo forest, watercolor painting style
10. The Robot, on a post-apocalyptic city street, sci-fi comic book style
11. The Knight, in a medieval castle, oil painting style
12. The Alien, on a distant planet, surrealistic painting style
13. The Samurai, in a cherry blossom garden, sumi-e painting style
14. The Detective, in a noir cityscape, graphic novel style
15. The Mermaid, in an underwater kingdom, fantasy painting style
16. The Robot, in a high-tech laboratory, digital painting style
17. The Pirate, on a ship deck, traditional painting style
18. The Samurai, in a misty forest, manga-style
19. The Steampunk Inventor, in a clockwork laboratory, Victorian-era painting style
20. The Ghost, in a haunted house, horror comic book style
21. The Knight, on a battlefield, historical painting style
22. The Alien, in a futuristic city, sci-fi painting style
23. The Vampire, in a gothic castle, dark painting style
24. The Samurai, in a snowy landscape, impressionist painting style
25. The Cyborg, in a virtual reality world, digital art style
26. The Cowboy, on a rugged mountain range, Western painting style
27. The Robot, in a futuristic cityscape, sci-fi painting style
28. The Detective, in a crime scene investigation, graphic novel style
29. The Mermaid, in a coral reef, watercolor painting style
30. The Pirate, on a deserted island, tropical painting style
31. The Knight, in a mythical forest, medieval painting style
32. The Alien, in a mysterious underground base, surrealistic painting style
33. The Samurai, in a tranquil garden, traditional painting style
34. The Cyborg, in a post-apocalyptic wasteland, digital art style
35. The Cowboy, in a small-town saloon, Western painting style
36. The Robot, in a cyberpunk megacity, sci-fi painting style
37. The Detective, in a neon-lit alleyway, graphic novel style
38. The Mermaid, in a shipwreck, fantasy painting style

39. The Pirate, in a bustling port city, oil painting style
 40. The Knight, in a royal castle, historical painting style
 41. The Alien, in a deep space nebula, surrealist painting style
 42. The Samurai, in a cherry blossom forest, sumi-e painting style
 43. The Cyborg, in a virtual reality cyberworld, digital art style
 44. The Cowboy, in a desert ghost town, Western painting style
 45. The Robot, in a futuristic laboratory, sci-fi painting style
 46. The Detective, in a city rooftop chase, graphic novel style
 47. The Mermaid, in a moonlit lagoon, fantasy painting style
 48. The Pirate, in a naval battle, traditional painting style
 49. The Knight, in a medieval joust, historical painting style
 50. The Alien, in an abandoned spaceship, surrealist painting style
 51. The Samurai, in a bamboo forest at night, sumi-e painting style
 52. The Explorer, in a dense jungle, documentary style
 53. The Assassin, in a city street at night, noir painting style
 54. The Monster, in a dark cave, horror movie style
 55. The Superhero, in a city skyline, comic book style
 56. The Samurai, in a foggy bamboo forest, manga-style
 57. The Dancer, on a stage, performing arts style
 58. The Alien, in a laboratory, sci-fi digital art style
 59. The Detective, in a murder investigation, crime drama style
 60. The Pirate, on a treasure island, adventure painting style
 61. The Knight, in a medieval siege, historical painting style
 62. The Samurai, in a peaceful pond, sumi-e painting style
 63. The Robot, in a factory, sci-fi painting style
 64. The Cowboy, on a rodeo, Western painting style
 65. The Witch, in a haunted forest, dark painting style
 66. The Astronaut, in a spacewalk, sci-fi digital art style
 67. The Superhero, in a city rooftop, comic book style
 68. The Samurai, in a fiery battlefield, manga-style
 69. The Cyborg, in a futuristic city, sci-fi painting style
 70. The Monster, in a graveyard, horror movie style
 71. The Pirate, on a ship deck in a storm, traditional painting style
 72. The Knight, in a medieval banquet, historical painting style
 73. The Samurai, in a cherry blossom tree, sumi-e painting style
 74. The Robot, in a futuristic transport system, sci-fi digital art style
 75. The Cowboy, on a horseback ride, Western painting style
 76. The Detective, in a chase sequence, crime drama style
 77. The Witch, in a spell-casting ritual, dark painting style
 78. The Astronaut, in a lunar colony, sci-fi digital art style
 79. The Superhero, in a city alleyway, comic book style
 80. The Samurai, in a sacred temple, manga-style
 81. The Cyborg, in a post-human society, sci-fi painting style
 82. The Monster, in a creepy mansion, horror movie style
 83. The Pirate, on a deserted beach, adventure painting style
 84. The Knight, in a medieval battle, historical painting style
 85. The Samurai, in a peaceful waterfall, sumi-e painting style
 86. The Robot, in a futuristic military base, sci-fi digital art style
 87. The Cowboy, on a cattle drive, Western painting style
 88. The Detective, in a spy thriller, crime drama style
 89. The Witch, in a dark forest at night, dark painting style
 90. The Astronaut, in a spaceship launch, sci-fi digital art style
 91. The Superhero, in a city skyscraper, comic book style
 92. The Samurai, in a snow-covered landscape, manga-style
 93. The Cyborg, in a robotic dystopia, sci-fi painting style
 94. The Monster, in a creepy attic, horror movie style
 95. The Pirate, on a high seas adventure, traditional painting style
 96. The Knight, in a royal coronation, historical painting style
 97. The Samurai, in a dragon fight, sumi-e painting style
 98. The Robot, in a futuristic cityscape at night, sci-fi digital art style
 99. The Cowboy, in a Western shootout, Western painting style
 100. The Detective, in a crime thriller, crime drama style
- List of 100 Prompts for Pose Accuracy Evaluation**
1. The adventurous explorer, in a dense jungle
 2. The daring race car driver, on a winding race track
 3. The graceful ice skater, on a glittering ice rink
 4. The skilled skateboarder, on a sunny skate park
 5. The mysterious magician, in a dark and dusty magic shop
 6. The brave firefighter, in a raging inferno
 7. The daring surfer, on a towering ocean wave
 8. The energetic breakdancer, in a busy city street
 9. The elegant ballroom dancer, in a grand ballroom
 10. The playful street artist, in a vibrant graffiti alleyway
 11. The graceful figure skater, on a frozen lake
 12. The bold parkour enthusiast, in an urban park
 13. The mysterious ninja, in a misty mountain landscape
 14. The daring rock climber, on a steep rock face
 15. The passionate flamenco dancer, in a bustling city square
 16. The skilled pole dancer, in a dimly lit club
 17. The determined marathon runner, on a scenic country road
 18. The elegant equestrian, on a sprawling ranch

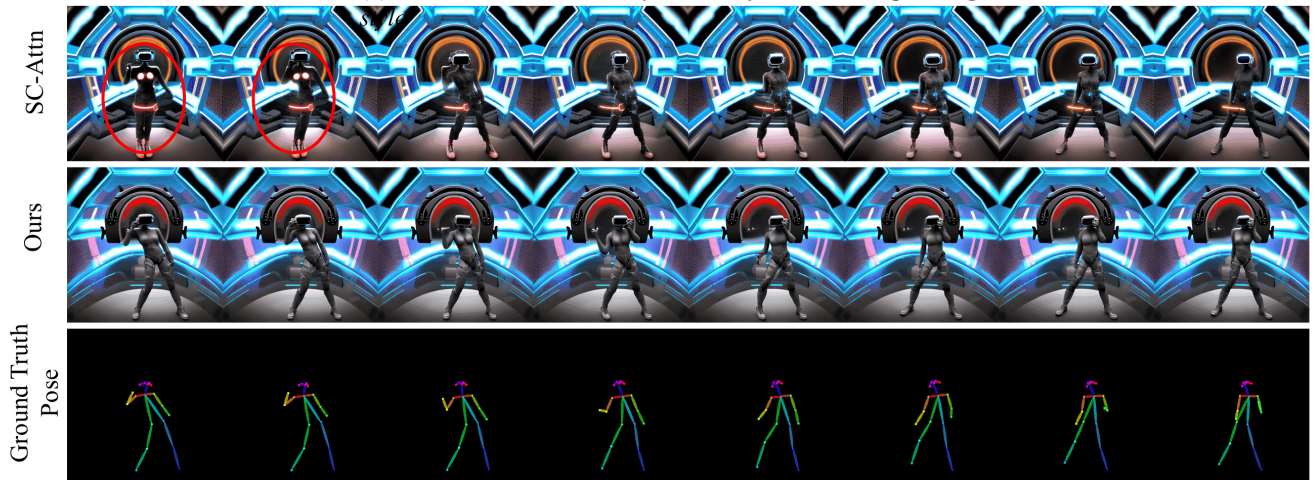
19. The adventurous hot air balloonist, in a sky full of clouds
20. The talented street performer, on a crowded pedestrian walkway
21. The daring skydiver, in a clear blue sky
22. The graceful gymnast, in a grand gymnastics arena
23. The skilled fencer, in a historic castle courtyard
24. The mystical sorceress, in a dimly lit forest
25. The daring mountain climber, on a towering mountain peak
26. The playful acrobat, in a colorful circus tent
27. The noble knight, in a medieval castle
28. The mysterious detective, in a dark alleyway
29. The glamorous fashion model, on a bustling city street
30. The fierce warrior, on a desolate mountaintop
31. The charming street musician, on a busy street corner
32. The adventurous mountaineer, on a snow-covered peak
33. The energetic zumba dancer, in a vibrant dance studio
34. The daring parkour enthusiast, on a futuristic city rooftop
35. The skilled trapeze artist, in a circus high wire act
36. The determined triathlete, on a grueling triathlon course
37. The graceful synchronized swimmer, in a luxurious pool
38. The playful breakdancer, in an urban street art scene
39. The bold cliff diver, on a rocky ocean cliff
40. The talented mime artist, in a quiet park
41. The daring rodeo cowboy, in a rodeo arena
42. The adventurous kayaker, on a fast-flowing river
43. The skilled BMX rider, on a dirt bike track
44. The passionate tango dancer, in a grand dance hall
45. The mischievous pirate, on a stormy sea
46. The bold bungee jumper, on a tall bridge
47. The daring ski jumper, on a snowy mountain slope
48. The glamorous Hollywood actor, on a red carpet event
49. The mystical mermaid, in a deep ocean abyss
50. The playful gymnastics coach, in a grand gymnasium
51. The mysterious vampire, in a dimly lit castle
52. The determined kickboxer, in a bustling gym
53. The adventurous surfer, on a secluded beach
54. The playful street dancer, in a crowded urban street
55. The skilled juggler, in a busy street fair
56. The daring base jumper, on a towering cliff
57. The noble samurai, in a peaceful cherry blossom garden
58. The passionate belly dancer, in a colorful Arabian bazaar
59. The daring wing walker, on a biplane soaring high above the clouds
60. The mystical fairy, in a magical forest
61. The talented magician, on a grand stage
62. The bold snowboarder, on a snowy mountain slope
63. The playful roller skater, in a bustling skate park
64. The daring ice climber, on a frozen waterfall
65. The graceful pole vaulter, on a busy track and field
66. The adventurous kayaker, in a breathtaking canyon
67. The skilled martial artist, in a quiet temple
68. The elegant ballet dancer, in a grand theater
69. The daring free runner, in an urban obstacle course
70. The graceful rhythmic gymnast, in a grand arena
71. The adventurous dirt bike rider, on a rough terrain
72. The bold street racer, on a busy city street
73. The talented ice sculptor, in a snowy park
74. The playful hula hooper, in a colorful park
75. The daring mountain biker, on a rocky mountain trail
76. The noble gladiator, in a grand coliseum
77. The mysterious ghost, in a dark haunted mansion
78. The graceful figure skater, on a frozen river
79. The adventurous rock climber, on a steep rock face
80. The daring cliff diver, on a tall ocean cliff
81. The passionate salsa dancer, in a vibrant dance hall
82. The bold wingsuit flyer, soaring through a majestic canyon
83. The talented street magician, in a crowded street market
84. The playful parkour athlete, on a busy city rooftop
85. The determined long distance runner, on a scenic trail
86. The skilled contortionist, in a grand circus tent
87. The bold rodeo bull rider, in a rodeo arena
88. The adventurous sky glider, in a sky full of clouds
89. The graceful equestrian vaulting performer, on a sprawling farm
90. The daring sky surfer, on a windy day
91. The talented ice skater, on a frozen pond
92. The passionate flamenco dancer, in a vibrant city square
93. The mysterious street magician, in a dark alleyway
94. The skilled cross-country skier, on a snowy mountain trail
95. The daring motocross rider, on a dirt track
96. The playful street performer, on a busy pedestrian walkway
97. The noble horseback rider, on a majestic countryside
98. The daring snowmobiler, on a snowy terrain
99. The adventurous white water rafter, on a fast-flowing river
100. The passionate tap dancer, on a grand theater stage



(a) "The Cowboy, on a dusty plain, oil painting style"

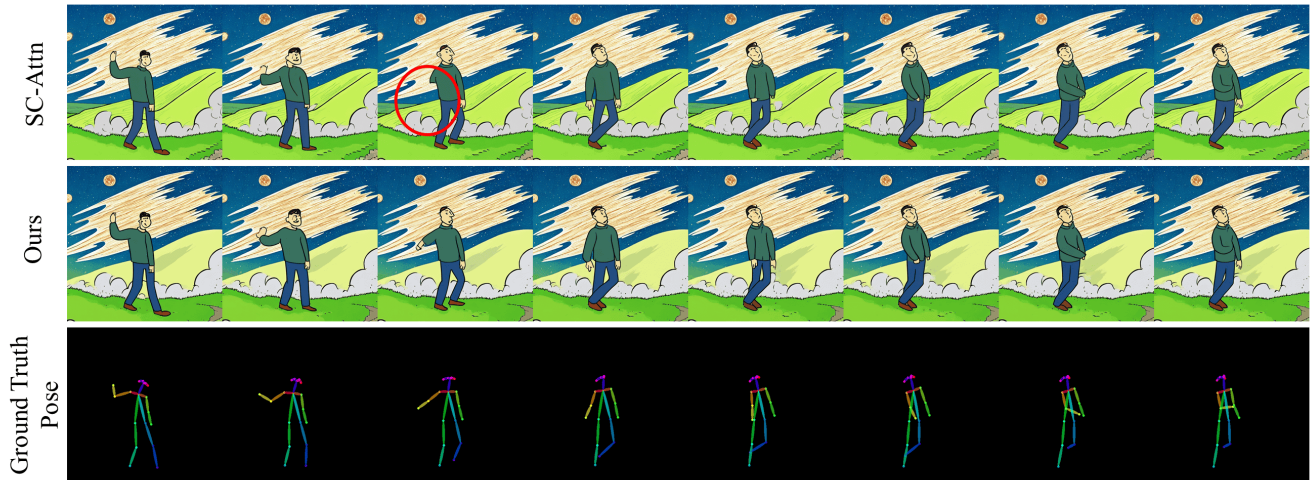


(b) "The Samurai, in a cherry blossom forest, sumi-e painting"

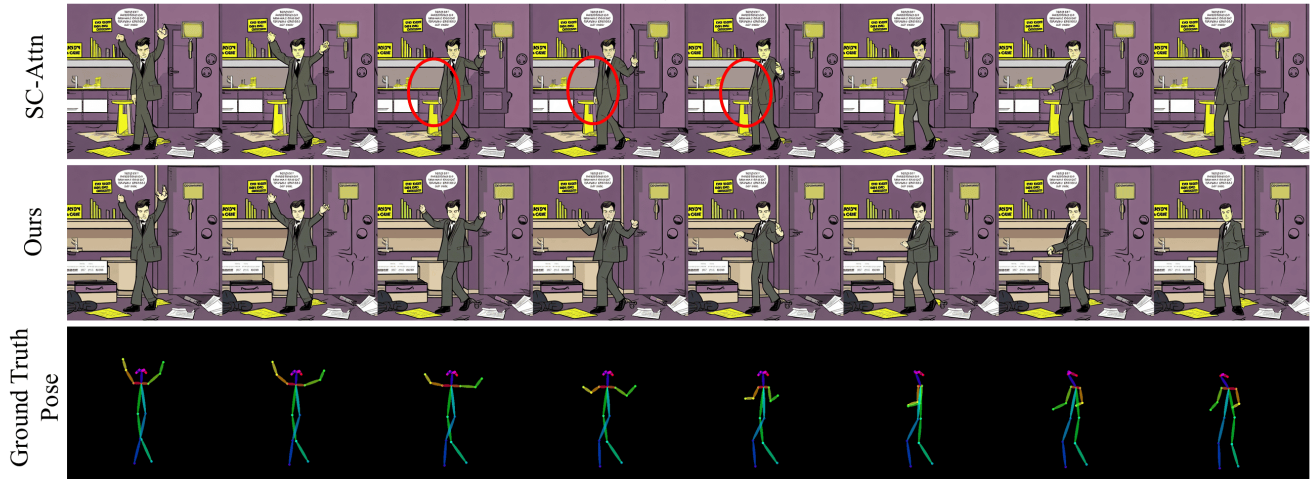


(c) "The Cyborg, in a virtual reality cyberworld, digital art style"

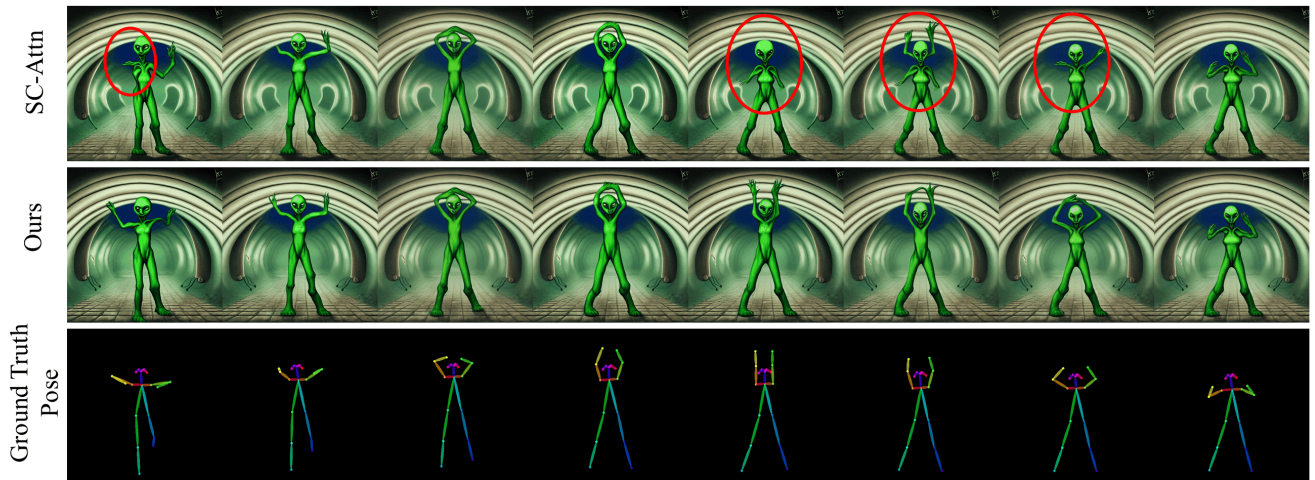
Figure 10: More ablation results on pose: SC-Attn vs Ours



(a) "man on hill watching a meteor; cartoon artwork"

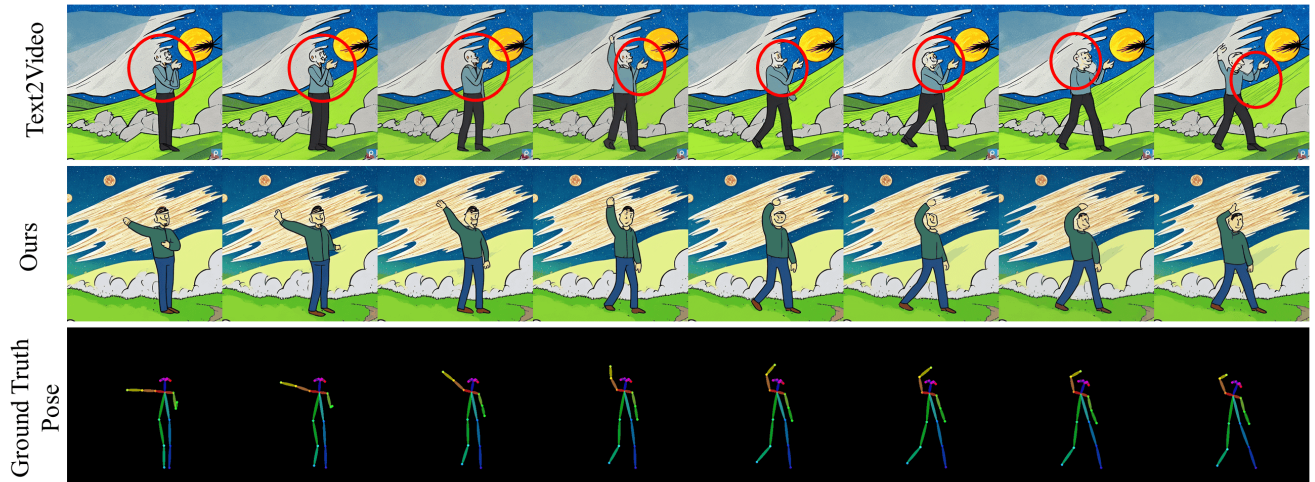


(b) "The Detective, in a crime scene investigation, graphic novel style"

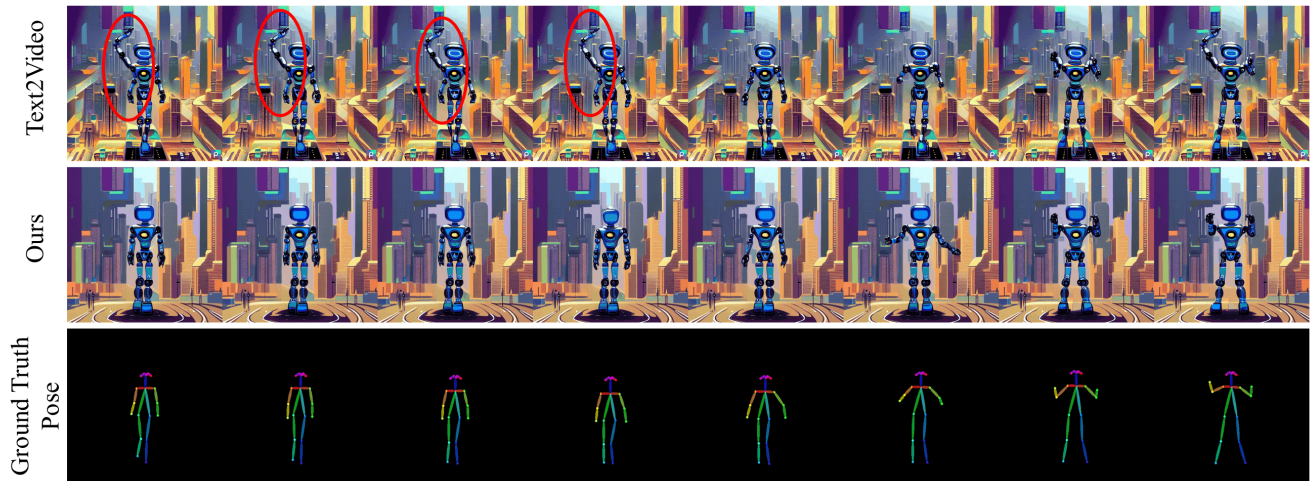


(c) "The Alien, in a mysterious underground base, surrealistic painting style"

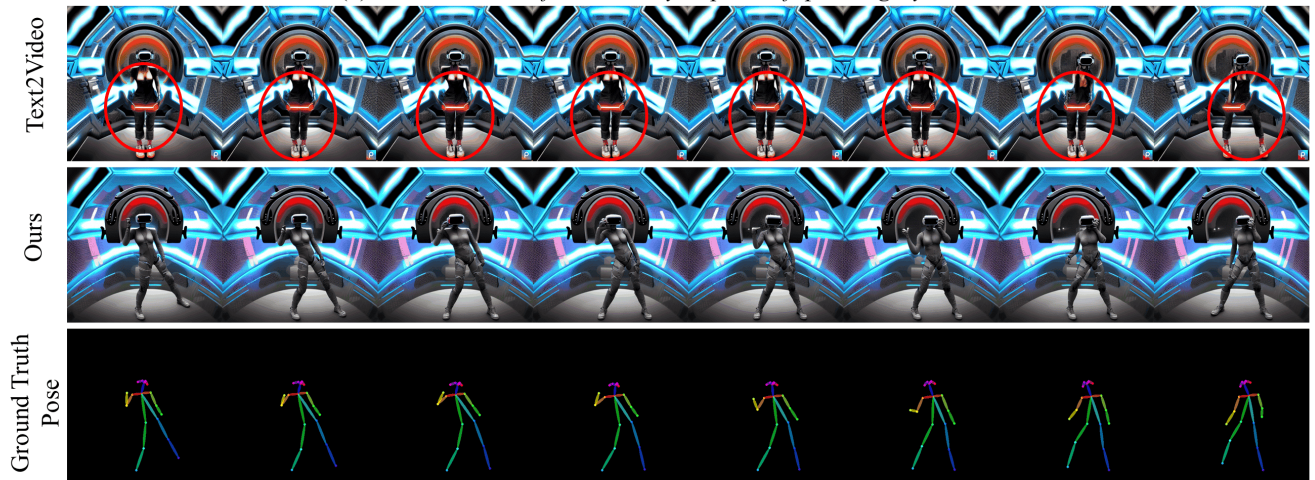
Figure 11: More ablation results on pose: 2D Control vs Ours



(a) "man on hill watching a meteor, cartoon artwork"



(b) "The Robot, in a futuristic cityscape, sci-fi painting style"



(c) "The Cyborg, in a virtual reality cyberworld, digital art style"

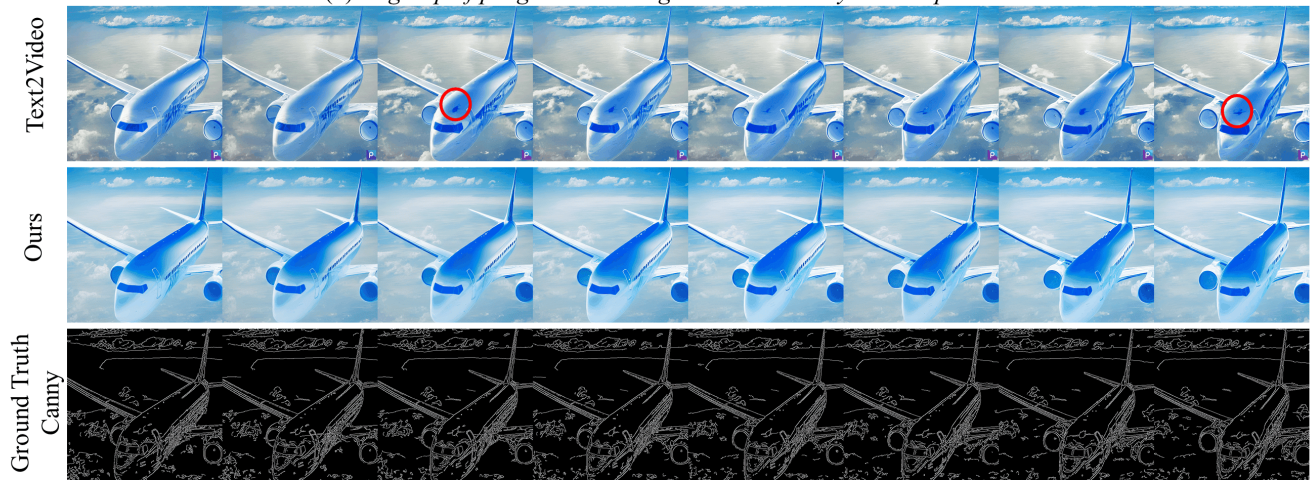
Figure 12: More comparison results on pose



(a) "The Steampunk Inventor, in a clockwork laboratory, Victorian-era painting style"

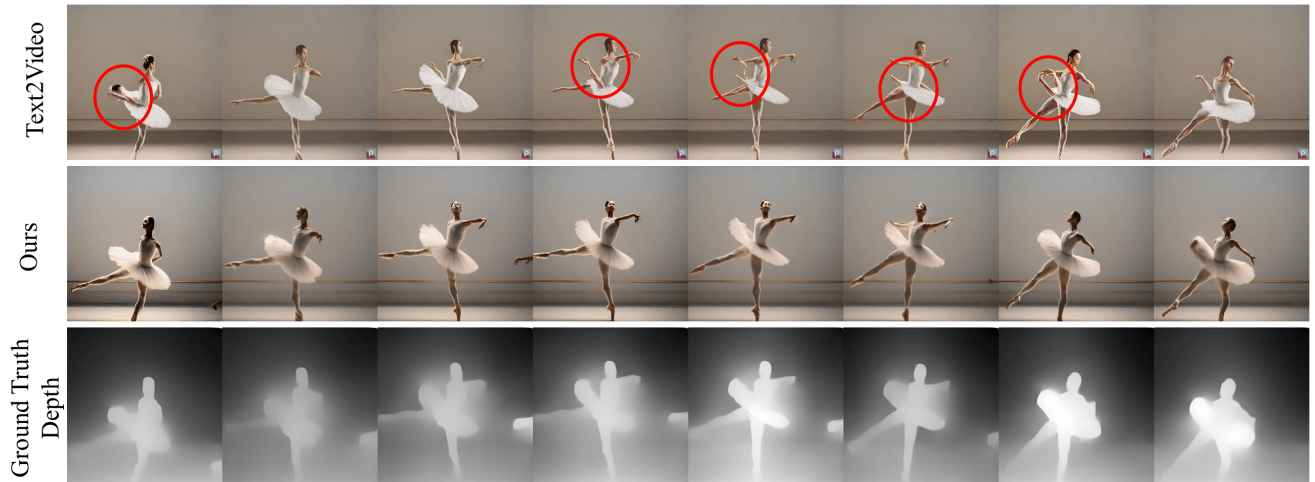


(b) "A group of penguins waddling across the snowy landscape"

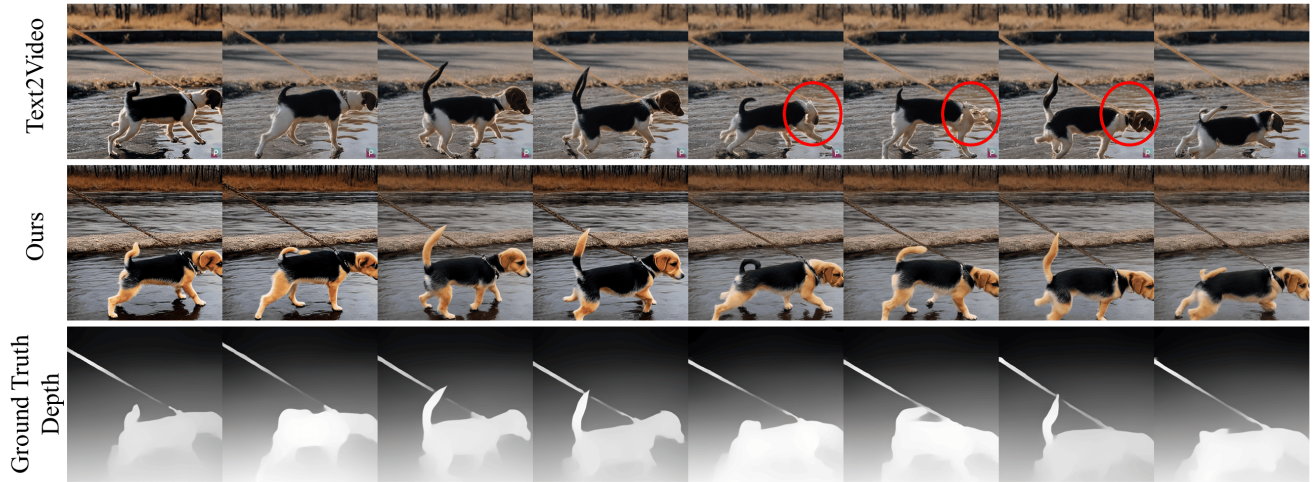


(c) "Passenger jet banking gracefully in the azure sky"

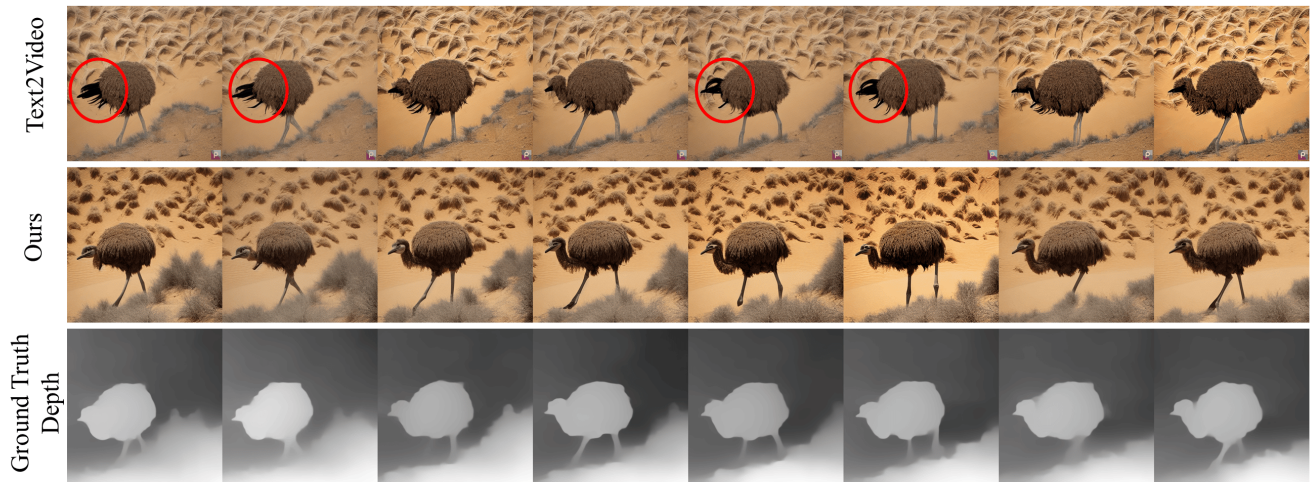
Figure 13: More comparison results on canny



(a) "A graceful ballet dancer on a moonlit stage."



(b) "A puppy with bright eyes and a happy demeanor walking by a river"



(c) "A tall ostrich taking deliberate steps through the desert sands"

Figure 14: More comparison results on canny